

Titre: Prédiction dynamique de durées de trajets de camions de transport
de minerai dans les mines souterraines

Auteur: Victor Simon
Author:

Date: 2024

Type: Mémoire ou thèse / Dissertation or Thesis

Référence: Simon, V. (2024). Prédiction dynamique de durées de trajets de camions de
transport de minerai dans les mines souterraines [Mémoire de maîtrise,
Citation: Polytechnique Montréal]. PolyPublie. <https://publications.polymtl.ca/59213/>

 **Document en libre accès dans PolyPublie**
Open Access document in PolyPublie

URL de PolyPublie: <https://publications.polymtl.ca/59213/>
PolyPublie URL:

**Directeurs de
recherche:** Robert Pellerin, & Michel Gamache
Advisors:

Programme: Maîtrise recherche en génie industriel
Program:

POLYTECHNIQUE MONTRÉAL

affiliée à l'Université de Montréal

**Prédiction dynamique de durées de trajets de camions de transport de minerais
dans les mines souterraines**

VICTOR SIMON

Département de mathématiques et de génie industriel

Mémoire présenté en vue de l'obtention du diplôme de *Maîtrise ès sciences appliquées*
Génie industriel

Août 2024

POLYTECHNIQUE MONTRÉAL

affiliée à l'Université de Montréal

Ce mémoire intitulé :

**Prédiction dynamique de durées de trajets de camions de transport de minerai
dans les mines souterraines**

présenté par **Victor SIMON**

en vue de l'obtention du diplôme de *Maîtrise ès sciences appliquées*
a été dûment accepté par le jury d'examen constitué de :

Bruno AGARD, président

Robert PELLERIN, membre et directeur de recherche

Michel GAMACHE, membre et codirecteur de recherche

Souheil-Antoine TAHAN, membre

DÉDICACE

*À Carolane,
le meilleur soutien que j'aurais pu trouver durant les trois dernières années.*

*À ma famille,
qui a réussi à rester derrière moi, même outre-Atlantique.*

*À Laurianne et Simon,
qui m'ont sorti des périodes de fortes turbulences.*

REMERCIEMENTS

En premier lieu, je tiens à remercier du fond du cœur mes directeurs de recherche Robert PELLERIN et Michel GAMACHE, professeurs à Polytechnique Montréal, qui ont su déployer des trésors de patience pour me laisser le temps de création nécessaire à l'élaboration du présent mémoire. Leur encadrement, leurs enseignements et les expériences stimulantes qu'ils m'ont permis de vivre ont été d'une valeur inestimable et ont rendu mon expérience de recherche véritablement exaltante. En outre, je souhaite leur adresser ma gratitude pour le soutien financier remarquable dont ils ont pu me faire bénéficier tout au long de mes travaux de recherche, me permettant d'avoir quotidiennement le cœur plus léger et de pouvoir concentrer mes efforts sur le plus important.

Je tiens à exprimer ma gratitude envers mon partenaire industriel anonyme et à tous ses membres concernés de près ou de loin par mon projet de recherche. En plus de m'avoir donné accès à sa précieuse base de données opérationnelles et aux plans des sites miniers étudiés, ce partenaire m'a permis de réaliser un stage Mitacs enrichissant de six mois et de vivre une visite mémorable d'une mine exceptionnelle.

Je souhaite remercier tous les responsables du Parcours Canada d'Arts et Métiers ParisTech, qui m'ont permis de prendre connaissance de l'existence de ce parcours de recherche en double-diplôme à Polytechnique Montréal il y a cinq ans, et qui m'ont informé de la démarche administrative à suivre. Je leur dois aujourd'hui la chance d'écrire ces lignes.

Mes remerciements s'adressent par ailleurs aux organismes qui ont subventionné ma recherche ou qui ont rendu possible les expériences dont j'ai pu profiter. Je citerai d'abord le Conseil de Recherche en Sciences Naturelles et en Génie du Canada (CRNSG) dont le programme « Mine Intelligente et Autonome » (MIA), obtenu par mon codirecteur de recherche Michel GAMACHE, m'a offert une « Formation Orientée vers la Nouveauté, la Collaboration et l'Expérience en Recherche » (FONCER) constituée d'une école d'été formidablement enrichissante, d'un stage Mitacs et de cours complémentaires ayant participé à solidifier mes connaissances en génie minier et à me développer personnellement. Je remercierai aussi le groupe MISA, pôle d'excellence minier, qui a participé au financement de l'école d'été susmentionnée et qui m'a, de surcroît, offert l'occasion de présenter mes travaux dans le cadre de la conférence annuelle « Québec Mines + Énergie » via son quatrième colloque « Mission Mine Autonome-2030 » (MMA-2030) intitulé « En route vers l'autonomisation ». Naturellement, je remercierai aussi Mitacs qui m'a permis de réaliser un stage rémunéré de six mois en parallèle de ma maîtrise via son programme Accélération.

Finalement, j'exprime ma reconnaissance envers toutes les personnes qui ont participé de près ou de loin à la réalisation et à la réussite de ce projet de maîtrise.

RÉSUMÉ

Caractérisé par ses conditions opérationnelles difficiles et sa versatilité multifactorielle, l'environnement minier souterrain complexifie grandement la tâche de prédiction de durées de temps de trajet de camions de transport de minerai. Bien que cette tâche soit nécessaire à une planification précise des successions d'opérations du quart de travail à venir, les techniques conventionnelles sont souvent restreintes à des mesures statistiques simplistes, lesquelles manquent cruellement de précision, ou bien à une expertise humaine. Naturellement, cette dernière est circonscrite à un site minier donné voire à certains itinéraires uniquement, sa disponibilité est limitée et la qualité de ses performances est difficilement quantifiable. Diamétralement opposée à ces techniques, notre approche inédite, robuste et sophistiquée se base sur les données de détections issues de systèmes émergents de balises souterraines d'identification de véhicules. Elle combine un algorithme complexe de préparation de données générant de multiples variables opérationnelles à partir de ces détections, et un modèle intégré d'apprentissage automatique capable d'exploiter les séries temporelles observées pour formuler ses prédictions de temps de trajet. Quel que soit le site minier souterrain étudié, si tant est qu'il dispose d'un système de balises de détection, notre succession de modèles permet de prédire avec une fiabilité accrue les temps de trajet de camions de transport de minerai sur tous les itinéraires prédominants du site, et ce même lorsque ledit système de balises présente des défaillances sérieuses. Mieux encore et complètement inédit au vu de la littérature actuelle, le modèle proposé ne nécessite pas d'autres sources de données pour être correctement entraîné et les valeurs des variables explicatives sont toutes connues à l'avance de chaque quart de travail par les planificateurs des opérations minières. Notre approche a été testée et validée non seulement sur le site minier sur lequel il a été développé mais également sur un autre site dissemblable dont le système de détection présentait une défaillance majeure, démontrant la robustesse et la capacité de généralisation de notre modèle à d'autres contextes miniers. Des propositions d'amélioration critiques des systèmes de balises de détection existants ont été formulées et de multiples pistes de recherche ont pu être proposées au vu de l'ampleur du travail qu'il reste à accomplir pour perfectionner les techniques de prédiction. Par toutes ses contributions, ce mémoire pave la voie à une amélioration substantielle de la planification et de l'efficacité opérationnelle des mines souterraines, marquant une avancée significative dans le domaine de l'ingénierie minière.

ABSTRACT

Characterized by its challenging operational conditions and multifactorial versatility, the underground mining environment greatly complicates the task of predicting ore haul truck travel times. While this task is necessary for precise planning of the sequence of operations for the upcoming shift, conventional techniques are often restricted to simplistic statistical measures, which are sorely lacking in precision, or to human expertise. Naturally, this expertise is confined to a given mining site or even to specific routes only, its availability is limited, and the quality of its performance is difficult to quantify. In stark contrast to these techniques, our novel, robust, and sophisticated approach relies on data from emerging underground vehicle identification beacon systems. It combines a complex data preparation algorithm that generates multiple operational variables from this data and an integrated machine learning model capable of leveraging the observed time series to make its travel time predictions. Regardless of the underground mining site studied, provided it has a beacon detection system, our succession of models can predict ore haul truck travel times on all predominant routes with increased reliability, even when said beacon system exhibits serious failures. Even more uniquely, given the current literature, the proposed comprehensive model does not require other data sources for proper training, and the prediction variable values are all known in advance of each shift by mining operation planners. Our approach has been tested and validated not only on the mining site where it was developed but also on another dissimilar site whose detection system had a major fault, demonstrating its robustness and generalizability to other mining contexts. Critical improvements to existing beacon detection systems were proposed, and numerous research avenues were suggested in light of the vast amount of work that remains to be done to perfect prediction techniques. With all these contributions, this thesis paves the way for substantial improvements in planning and operational efficiency in underground mines, marking a significant advance in the field of mining engineering.

TABLE DES MATIÈRES

DÉDICACE	iii
REMERCIEMENTS	iv
RÉSUMÉ	v
ABSTRACT	vi
TABLE DES MATIÈRES	vii
LISTE DES TABLEAUX	xi
LISTE DES FIGURES	xiii
LISTE DES SIGLES ET ABRÉVIATIONS	xvii
CHAPITRE 1 INTRODUCTION	1
CHAPITRE 2 REVUE DE LITTÉRATURE	4
2.1 Définitions et concepts de base	4
2.1.1 Mines	4
2.1.2 Camions de transport de minerai	4
2.1.3 Trajets des HT	5
2.1.4 TT	6
2.1.5 Prédiction de TT	6
2.2 Protocole de recherche d'articles scientifiques	6
2.3 Analyse des résultats	10
2.3.1 Approches basées sur des modèles de ML	10
2.3.2 Approches basées sur des modèles de simulation	16
2.3.3 Autres approches	18
2.4 Revue critique	19
2.5 Conclusion	21
CHAPITRE 3 MÉTHODOLOGIE DE RECHERCHE	22
3.1 Objectifs de recherche	22
3.2 Méthodologie de recherche	23

3.3 Conclusion	23
CHAPITRE 4 DÉVELOPPEMENT D'UN MODÈLE DE PRÉPARATION DE DON-	
NÉES	24
4.1 Description du cas d'étude	24
4.2 Contraintes et spécificités du cas d'étude	26
4.2.1 Activités des HT	26
4.2.2 Étude des trajets historiques des HT	28
4.3 Spécifications des requis du modèle	29
4.4 Identification des variables d'intérêt et des données correspondantes	29
4.4.1 Variables d'intérêt	29
4.4.2 Données correspondantes	33
4.5 Préparation des données	35
4.5.1 Extraction du jeu de données d'intérêt	35
4.5.2 Identification du niveau de profondeur de chaque balise	36
4.5.3 Correction d'incohérences entre les périodes de collecte de données	39
4.5.4 Identification des itinéraires prédominants	40
4.5.5 Listage des balises de changement de niveau	42
4.5.6 Création d'un algorithme de reconnaissance de trajets sans détours	47
4.5.7 Analyse d'histogrammes de TT sur les trajets prédominants et améliorations résultantes	53
4.5.8 Étude de déplacements erratiques de HT sur de faibles distances	71
4.5.9 Génération de variables d'intérêt temporelles	75
4.5.10 Estimation du nombre de HT se déplaçant activement sur un itinéraire donné durant un quart de travail donné	76
4.5.11 Évaluation de l'influence de variables d'intérêt sur les TT	77
4.6 Inventaire des pistes d'amélioration relatives aux capteurs, à l'acquisition des données et à leur gestion	84
4.6.1 Pistes d'amélioration relatives aux capteurs	85
4.6.2 Pistes d'amélioration relatives à l'acquisition des données	86
4.6.3 Piste d'amélioration relative à la gestion des données : élimination des périodes de confusion d'identifiants de véhicules	88
4.7 Prétraitement des données préparées	89
4.8 Conclusion	91
CHAPITRE 5 DÉVELOPPEMENT D'UN MODÈLE INTÉGRÉ DE PRÉDICTION	
DE TEMPS DE TRAJETS	92

5.1	Spécifications du support informatique utilisé	92
5.2	Spécification des requis du modèle	92
5.3	Architecture du modèle, sélection des sous-modèles et fonctionnement théorique	94
5.3.1	Sous-modèle de partitionnement	95
5.3.2	Sous-modèles de prédiction de TT	96
5.4	Sélection de l'itinéraire de test et identification du seuil de filtrage à privilégier	104
5.4.1	Sélection de l'itinéraire de test	104
5.4.2	Identification du seuil de filtrage de TT à privilégier	106
5.5	Partitionnement des données	109
5.6	Prédiction de TT	118
5.6.1	Prédictions initiales de TT sur un itinéraire donné	118
5.6.2	Prédiction finale du TT par un modèle d'empilement	137
5.7	Évaluation de la qualité des prédictions de TT	140
5.8	Discussion	147
5.9	Conclusion	152
CHAPITRE 6 TEST DE GÉNÉRALISATION DU MODÈLE GLOBAL		156
6.1	Description du cas d'étude	156
6.2	Contraintes et spécificités du cas d'étude	156
6.3	Attributs correspondants aux variables d'intérêt	157
6.4	Application de notre modèle de préparation de données	159
6.4.1	Extraction du jeu de données d'intérêt et remarques générales	159
6.4.2	Identification du niveau de profondeur de chaque balise	160
6.4.3	Recherche d'incohérences entre les périodes de collecte de données	160
6.4.4	Identification de l'itinéraire prédominant à étudier	167
6.4.5	Listage des balises de changement de niveau	168
6.4.6	Utilisation de notre algorithme de reconnaissance de trajets	169
6.4.7	Analyse d'histogrammes de TT sur l'itinéraire prédominant et amélioration résultante	170
6.4.8	Génération de variables d'intérêt additionnelles	175
6.4.9	Évaluation de l'influence de variables d'intérêt sur les TT	176
6.4.10	Inventaire des pistes d'amélioration relatives aux capteurs, à l'acquisition des données et à leur gestion	181
6.4.11	Prétraitement des données	182
6.4.12	Conclusion de la préparation de données	182
6.5	Application de notre modèle intégré de prédiction de TT	183

6.5.1	Sélection de l'itinéraire de test	183
6.5.2	Partitionnement des données	184
6.5.3	Prédiction de TT	189
6.6	Conclusion	200
CHAPITRE 7 CONCLUSION		202
7.1	Synthèse des travaux	202
7.1.1	Préparation des données	202
7.1.2	Modèles de prédiction	202
7.1.3	Application à deux sites miniers	203
7.2	Limitations de la solution proposée	203
7.3	Améliorations futures	204
RÉFÉRENCES		206

LISTE DES TABLEAUX

Tableau 2.1 – Plan de concepts	7
Tableau 2.2 – Présentation des articles retenus (ordre antichronologique)	11
Tableau 4.1 – Caractérisation des attributs sélectionnés dans la BDD de la mine 1	34
Tableau 4.2 – Stratégie naïve d’identification de balises de changement de niveau en descente	46
Tableau 5.1 – Spécifications techniques de l’appareil utilisé dans le présent mémoire.	93
Tableau 5.2 – Ensembles de sélection des hyperparamètres du MLP et classe associée	115
Tableau 5.3 – Combinaison optimale des hyperparamètres du MLP	116
Tableau 5.4 – Ensembles de sélection des hyperparamètres pour XGBoost, RF, GBR et SVM	119
Tableau 5.5 – Combinaisons optimales d’hyperparamètres de XGBoost, RF, GBR et SVM	121
Tableau 5.6 – Ensembles de sélection des hyperparamètres du BRNN et classe associée	122
Tableau 5.7 – Combinaison optimale des hyperparamètres de notre BRNN.	123
Tableau 5.8 – Ensembles de sélection des hyperparamètres du LSTM et classe associée	127
Tableau 5.9 – Combinaison optimale des hyperparamètres de notre LSTM.	128
Tableau 5.10 – Combinaison optimale des hyperparamètres du MLP	130
Tableau 5.11 – Combinaisons optimales d’hyperparamètres pour nos six modèles appliqués à un segment inter-niveaux	131
Tableau 5.12 – Combinaison optimale des hyperparamètres du MLP	135
Tableau 5.13 – Combinaisons optimales d’hyperparamètres pour nos six modèles appliqués à un segment majeur	137
Tableau 5.14 – Ensembles de sélection des hyperparamètres pour les modèles utilisés	139
Tableau 5.15 – Combinaison optimale des hyperparamètres du GBR, modèle d’empilement optimal	139
Tableau 5.16 – Performances de nos six modèles de ML et du modèle de référence (« Référence ») pour l’itinéraire complet	141
Tableau 5.17 – Performances finales de nos sous-modèles de prédiction initiale de TT, du modèle d’empilement et du modèle de référence.	144
Tableau 6.1 – Caractérisation des attributs sélectionnés dans la BDD de la mine 2	158

Tableau 6.2 – Bornes de la période d’activité et nombre de détections de chacun des HT	160
Tableau 6.3 – Reproduction simplifiée des séquences indésirables de données de détection observées	162
Tableau 6.4 – Combinaison optimale des hyperparamètres du MLP pour les trajets conventionnels	186
Tableau 6.5 – Combinaison optimale des hyperparamètres du MLP pour les trajets autonomes	188
Tableau 6.6 – Combinaison optimale des hyperparamètres de notre LSTM pour les trajets conventionnels sur l’itinéraire complet Niv3→Niv30	189
Tableau 6.7 – Combinaison optimale des hyperparamètres de notre LSTM pour les trajets autonomes sur l’itinéraire complet Niv3→Niv30	190
Tableau 6.8 – Combinaison optimale des hyperparamètres de notre LSTM pour les trajets conventionnels sur l’itinéraire sectionné Niv3→Niv30.	191
Tableau 6.9 – Performances successives de nos sous-modèles de prédiction initiale, du modèle d’empilement et du modèle de référence appliqués aux TT conventionnels	199

LISTE DES FIGURES

Figure 2.1 – Protocole d'identification des articles.....	9
Figure 4.1 – Architecture fondamentale de la mine 1	25
Figure 4.2 – Aperçu des activités de transport de minerai des HT dans la mine 1 ...	25
Figure 4.3 – Graphe d'évolution de la profondeur d'un HT de la mine 1 sur 22 jours	42
Figure 4.4 – Emplacement des balises de changement de niveau.....	43
Figure 4.5 – Histogramme des TT observés sur l'itinéraire surface→niveau200 dans la mine 1	53
Figure 4.6 – Modélisation d'une configuration de trajet qui rallonge le TT A→A', pouvant provoquer un second mode.....	56
Figure 4.7 – Histogramme des TT A→A' repérés par notre algorithme.....	57
Figure 4.8 – Histogramme des TT observés sur l'itinéraire A'→Niveau150, faisant figurer des TT Q-I	59
Figure 4.9 – Histogramme des TT observés sur l'itinéraire A'→Niveau200 dans la mine 1	66
Figure 4.10 – Histogramme des TT observés sur l'itinéraire Niveau200→A' dans la mine 1	67
Figure 4.11 – Histogramme des TT observés sur l'itinéraire Surface→Niveau300 via la rampe 2 dans la mine 1	69
Figure 4.12 – Histogramme des TT observés sur l'itinéraire Niveau300→Surface via la rampe 2 dans la mine 1	69
Figure 4.13 – Histogramme des TT observés sur l'itinéraire Niveau300→Surface via la rampe 1 dans la mine 1	70
Figure 4.14 – Diagramme à barres du TT médian selon le quart de travail sur l'itinéraire Surface→Niveau300 de la rampe 2 et barres d'erreur pour un niveau de confiance de 99,7%, seuil de filtrage établi via Tukey	79
Figure 4.15 – Diagramme à barres du TT moyen selon le quart de travail sur l'itinéraire Surface→Niveau300 de la rampe 2 et barres d'erreur pour un niveau de confiance de 99,7%, filtrage des TT supérieurs 30 minutes ...	80
Figure 4.16 – Diagramme à barres du TT moyen selon le nombre de HT actifs sur l'itinéraire Surface→Niveau300 de la rampe 2 et barres d'erreur pour un niveau de confiance de 95%, filtrage des TT supérieurs 30 minutes..	82

Figure 4.17 – Diagramme à barres du TT moyen selon le nombre de HT actifs sur l’itinéraire Surface→Niveau300 de la rampe 2 et barres d’erreur pour un niveau de confiance de 95%, seuil de filtrage établi via Tukey	83
Figure 4.18 – Diagramme à barres du TT moyen selon le nombre de HT actifs sur l’itinéraire Surface→Niveau300 de la rampe 2 et barres d’erreur pour un niveau de confiance de 95%, filtrage des 5% des TT les plus longs . .	84
Figure 4.19 – Diagramme à barres du TT moyen selon l’identifiant du HT correspondant sur l’itinéraire Surface→Niveau300 de la rampe 2 et barres d’erreur pour un niveau de confiance de 95%, filtrage des 5% des TT les plus longs	85
Figure 5.1 – Organigramme décrivant l’architecture de notre modèle de prédiction de TT au sein de notre méthode globale	94
Figure 5.2 – Architecture fondamentale de la mine 1	99
Figure 5.3 – Histogramme des TT observés sur l’itinéraire surface→niveau350 via la rampe 2 dans la mine 1	106
Figure 5.4 – Illustration des résultats du GMM lors de l’inférence de cinq groupes dans X_{appr} pour l’itinéraire surface→niveau300 via la rampe 2 dans la mine 1	111
Figure 5.5 – Évolution de la valeur de l’AIC en fonction du nombre de groupes inférés par le GMM dans y_{appr} pour l’itinéraire surface→niveau350 via la rampe 2 dans la mine 1	112
Figure 5.6 – Inférence de cinq groupes par le GMM dans y_{appr} pour l’itinéraire surface→niveau350 via la rampe 2 dans la mine 1	113
Figure 5.7 – Code Python de la classe personnalisée Capsule	115
Figure 5.8 – Diagramme à barres des vrais groupes annoncés par le GMM et des groupes prédits par le MLP pour l’itinéraire surface→niveau350 via la rampe 2 dans la mine 1	117
Figure 5.9 – Inférence de six groupes par le GMM dans y_{appr} pour le segment majeur surface→niveau300 via la rampe 2 dans la mine 1	135
Figure 5.10 – Diagramme à barres des vrais groupes annoncés par le GMM et des groupes prédits par le MLP pour l’itinéraire surface→niveau350 via la rampe 2 dans la mine 1	136
Figure 5.11 – Histogrammes des TT réels et des TT prédits par le SVM pour l’itinéraire complet	142
Figure 5.12 – Histogrammes des TT réels et des TT prédits par notre modèle d’empilement	145

Figure 6.1 – Tracé continu des profondeurs successives occupées par le HT1 de la mine 2 sur une période de cinq mois	161
Figure 6.2 – Graphe des profondeurs ponctuelles successives occupées par l'un des HT de la mine 2 sur une période de cinq mois	163
Figure 6.3 – Visualisation matricielle des absences de données de balises dans la mine 2	165
Figure 6.4 – Histogramme des TT conventionnels sur Niv3→Niv30 dans la mine 2 ..	171
Figure 6.5 – Histogramme à deux modes suspect des TT théoriquement autonomes sur Niv3→Niv30 dans la mine 2	171
Figure 6.6 – Histogramme des TT autonomes sur Niv3→Niv30 dans la mine 2	172
Figure 6.7 – Histogramme des TT autonomes sur Niv30→Niv3 dans la mine 2	174
Figure 6.8 – Histogramme des TT conventionnels sur Niv30→Niv3 dans la mine 2 ..	174
Figure 6.9 – Diagramme à barres du TT conventionnel moyen selon le quart de travail sur Niv3→Niv30 dans la mine 2 et barres d'erreur pour un niveau de confiance de 99,7%	177
Figure 6.10 – Diagramme à barres du TT conventionnel médian selon le quart de travail sur Niv3→Niv30 dans la mine 2 et barres d'erreur pour un niveau de confiance de 99,7%	177
Figure 6.11 – Diagramme à barres du TT conventionnel moyen selon le nombre de HT actifs sur Niv3→Niv30 dans la mine 2 et barres d'erreur pour un niveau de confiance de 95%	179
Figure 6.12 – Diagramme à barres du TT conventionnel médian selon le nombre de HT actifs sur Niv3→Niv30 dans la mine 2 et barres d'erreur pour un niveau de confiance de 95%	179
Figure 6.13 – Diagramme à barres du TT conventionnel moyen selon le numéro du HT sur Niv3→Niv30 dans la mine 2 et barres d'erreur pour un niveau de confiance de 99,7%	180
Figure 6.14 – Diagramme à barres du TT conventionnel médian selon le numéro du HT sur Niv3→Niv30 dans la mine 2 et barres d'erreur pour un niveau de confiance de 99,7%	180
Figure 6.15 – Diagramme à barres du TT conventionnel moyen selon le jour de la semaine sur Niv3→Niv30 dans la mine 2 et barres d'erreur pour un niveau de confiance de 95%	181
Figure 6.16 – Évolution de l'AIC et du BIC lorsque le GMM est directement appliqué à X_{train} conventionnel sur Niv3→Niv30 dans la mine 2	184

Figure 6.17 – Évolution de l'AIC et du BIC lorsque le GMM est appliqué à y_{train} conventionnel sur Niv3→Niv30 dans la mine 2.	185
Figure 6.18 – Partitionnement en quatre groupes de y conventionnel sur Niv3→Niv30 dans la mine 2 via le GMM.	185
Figure 6.19 – Diagramme à barres des groupes prédits par le MLP par rapport à ceux prévus par le GMM pour les TT conventionnels sur Niv3→Niv30 dans la mine 2	186
Figure 6.20 – Évolution de l'AIC et du BIC lorsque le GMM est appliqué à y_{train} autonome sur Niv3→Niv30 dans la mine 2.	187
Figure 6.21 – Partitionnement en quatre groupes de y autonome sur Niv3→Niv30 dans la mine 2 via le GMM.	187
Figure 6.22 – Diagramme à barres des groupes prédits par le MLP par rapport à ceux prévus par le GMM pour y autonome sur Niv3→Niv30 dans la mine 2 .	188
Figure 6.23 – Comparaison des histogrammes des TT conventionnels réels et prédits par le LSTM seul sur Niv3→Niv30 dans la mine 2	193
Figure 6.24 – Comparaison des histogrammes de TT autonomes réels et prédits par le LSTM seul sur Niv3→Niv30 dans la mine 2.	194
Figure 6.25 – Comparaison des histogrammes des TT conventionnels réels et prédits par notre modèle global de prédiction sur Niv3→Niv30 dans la mine 2 .	198

LISTE DES SIGLES ET ABRÉVIATIONS

AIC	Critère d'information d'Akaike (<i>Akaike Information Criterion</i>)
ANN	Réseau de neurones artificiels (<i>Artificial Neural Network</i>)
API	Interface de programmation d'application (<i>Application Programming Interface</i>)
BDD	Base De Données
BIC	Critère d'information bayésien (<i>Bayesian Information Criterion</i>)
BPNN	Modèle de réseau de neurones à rétropropagation du gradient (<i>Back-Propagation Neural Network</i>)
BRNN	Réseau de neurones régularisé bayésien (<i>Bayesian Regularized Neural Network</i>)
CPU	Unité centrale de traitement (<i>Central Processing Unit</i>)
CSV	Valeurs séparées par des virgules (<i>Comma-Separated Values</i>)
DRM	Méthodologie de recherche en conception (<i>Design Research Methodology</i>)
DT	Modèle d'arbres de décision (<i>Decision Trees</i>)
ELM	Modèle d'apprentissage automatique extrême (<i>Extreme Learning Machine</i>)
ESG	engagements Environnementaux, Sociaux et de Gouvernance
GBR	Modèle de régression du <i>boosting</i> du gradient (<i>Gradient Boosting Regressor</i>)
GLM	Modèle linéaire généralisé (<i>Generalized Linear Model</i>)
GMM	Modèle de mélange gaussien (<i>Gaussian Mixture Model</i>)
GPU	Unité de traitement graphique (<i>Graphics Processing Unit</i>)
HT	Camion de transport de minerai (<i>Haul Truck</i>)
IoT	Internet des objets (<i>Internet of Things</i>)
JSON	Notation objet de JavaScript (<i>JavaScript Object Notation</i>)
kNN	Modèle des k -plus proches voisins (<i>k-Nearest Neighbours</i>)
KPI	Indicateur clé de performance (<i>Key Performance Indicator</i>)
LHD	Chargeuse sur pneus (<i>Load-Haul-Dump</i>)
LightGBM	Modèle d'apprentissage automatique léger de <i>boosting</i> du gradient (<i>Light Gradient Boosting Machine</i>)
LSTM	Réseau de neurones à mémoire à court et long terme (<i>Long Short-Term Memory</i>)

MAD	DéviatiOn absolue moyenne (<i>Mean Absolute Deviation</i>)
MAE	Erreur absolue moyenne (<i>Mean Absolute Error</i>)
MAPE	Pourcentage d'écart absolu moyen (<i>Mean Absolute Percentage Error</i>)
MILP	Modèle linéaire de programmation en nombres entiers mixtes (<i>Mixed-Integer Linear Programming</i>)
ML	Apprentissage automatique (<i>Machine Learning</i>)
MLP	Perceptron multicouche (<i>Multilayer Perceptron</i>)
MLR	Régression linéaire multiple (<i>Multiple Linear Regression</i>)
MSE	Erreur quadratique moyenne (<i>Mean Squared Error</i>)
PCA	Analyse en composantes principales (<i>Principal Component Analysis</i>)
PDG	Président Directeur Général
Q-I	Quasi-Instantané
RAM	Mémoire vive (<i>Random Access Memory</i>)
RF	Modèle de forêt aléatoire (<i>RandomForest</i>)
RFID	Identification par radiofréquence (<i>Radio Frequency IDentification</i>)
RNN	Réseau de neurones récurrent (<i>Recurrent Neural Network</i>)
RMSE	Racine de l'erreur quadratique moyenne (<i>Root Mean Squared Error</i>)
R^2	Coefficient de détermination
SVM	Machine à vecteurs de support (<i>Support Vector Machine</i>)
TT	Temps de Trajet
XGBoost	Modèle de <i>boosting</i> extrême du gradient (<i>eXtreme Gradient Boosting</i>)

CHAPITRE 1 INTRODUCTION

De même que l’agriculture, la sylviculture et la pêche, l’exploitation des mines est l’une des quatre activités composant le secteur économique primaire [1]. À ce titre, le secteur minier constitue incontestablement l’un des piliers économiques de notre société. Pour autant, des événements récents ont pu montrer dans quelles mesures la stabilité de ce pilier reste relative, les facteurs de risques ne manquant pas.

Ce constat est justement celui qui ressort de l’enquête annuelle menée par le cabinet d’audit financier et de conseil EY entre juin et septembre 2022 portant sur les principaux risques et opportunités business pour le secteur minier en 2023 [2]. S’y l’on se fie à cette enquête, parmi les enjeux ayant brutalement gagné de l’envergure au classement des risques et opportunités de 2023 par rapport à 2022, on retrouve des enjeux géopolitiques, de perturbation des approvisionnements, de main d’oeuvre, mais aussi et surtout des critères de coût et de productivité, passés de la quinzième à la cinquième place du classement en l’espace d’une année. Toujours d’après cette enquête, l’inflation brutale récente, qui impacte significativement le coût des opérations minières, serait la raison pour laquelle ces derniers critères devraient constituer pour le secteur minier un facteur de risques nettement plus important pour l’année 2023 que pour l’année précédente.

La situation délicate provoquée par la conjonction de tous ces facteurs de risques ajoute inmanquablement une certaine pression aux industriels du secteur qui doivent se réinventer pour gagner en productivité et réduire leurs coûts d’exploitation. Ces innovations ne doivent toutefois pas mettre en péril leurs engagements Environnementaux, Sociaux et de Gouvernance (ESG) qui doivent figurer parmi leurs priorités en tout temps puisqu’ils figuraient à la première place du classement de l’enquête d’EY pour 2023 et pour 2022 [2].

Pour parvenir à marier ces objectifs potentiellement divergents, les industriels miniers commencent à s’appuyer sur une évolution du secteur minier parfois qualifiée de révolution technologique : l’avènement des « mines 4.0 » [3]. Ce terme, dérivé du concept d’Industrie 4.0, désigne littéralement des mines du futur pourvues de l’ensemble des évolutions habituellement associées à l’Industrie 4.0. On citera en particulier des innovations technologiques de rupture récentes qui modifient et modifieront en profondeur l’industrie minière : systèmes de communication sans fil souterrains, équipements connectés et pilotables à distance voire purement autonomes, géolocalisation souterraine ou encore implantation massive de capteurs fixes ou embarqués permettant la collecte de mégadonnées dites « de télémétrie » [4].

Parmi les grands axes de développement permis par ces nouvelles technologies, l’amélioration

potentielle de la planification court terme des activités minières est tout à fait notable. La planification court terme, ou planification opérationnelle, consiste à ordonnancer avec précision l'ordre dans lequel s'effectueront les opérations *a minima* durant le prochain quart de travail, parfois les quarts suivants voire jusqu'à toute la semaine à venir [5]. Ce processus de gestion est flexible et fréquemment corrigé [6], il peut s'appuyer sur des données opérationnelles obtenues en temps réel et sur d'autres informations concernant des modifications du contexte de fonctionnement. Proche des opérations à venir, la planification court terme donne des résultats plus fiables que les planifications à plus long terme. Elle permet aussi de réduire les durées opérationnelles dédiées à des activités non productrices de valeur ajoutée (en particulier les temps d'attente), ce qui améliore la productivité globale des activités d'extraction. Par ailleurs, une bonne planification possède des avantages supplémentaires : elle favorise le maintien d'un environnement de travail sécuritaire durant les opérations minières et permet généralement de diminuer la consommation de carburant des équipements miniers par tonne de minerai extrait, ce qui s'avère intéressant d'un point de vue financier et environnemental. L'amélioration de la planification court terme d'activités minières est donc une stratégie prometteuse dans le contexte qui affecte aujourd'hui le secteur minier.

Toutefois, sa réalisation impose aux industriels miniers de relever certains défis, à commencer par la maîtrise de la prédiction de la durée des activités minières à ordonner en fonction du contexte opérationnel à venir. Parmi celles-ci, l'activité de transport de minerai constitue un enjeu majeur en termes de productivité et de réduction des coûts de fonctionnement [7]. Cette activité étant généralement assurée par des camions de transport (*Haul Trucks* (HT)) de minerai, on pourra considérer qu'une bonne capacité d'estimation des Temps de Trajet (TT) de ces HT est donc nécessaire à l'amélioration de la planification court terme de toute exploitation minière. La criticité de cette tâche nous porte à croire qu'elle mérite d'être étudiée de manière approfondie.

Aussi, dans le présent mémoire, nous nous fixerons pour objectif d'améliorer la prédiction de TT de HT.

À cet effet, voici l'organisation retenue pour la réalisation du présent mémoire. Dans ce premier chapitre, le sujet du mémoire, le contexte global qui s'y rattache et les objectifs de recherche ont été introduits. Au chapitre 2, les travaux de recherche existants en lien avec notre sujet seront explorés dans une revue de littérature systématique et les lacunes dans la recherche actuelle seront soulignées au cours d'une revue critique. Dans le chapitre 3, nous expliciterons notre objectif spécifique ainsi que la démarche que nous adopterons pour l'atteindre. Nous profiterons ensuite du chapitre 4 pour décrire notre préparation de données. Le chapitre 5 sera dédié au développement et à l'application de notre modèle de prédiction

de TT. Le chapitre 6 nous permettra de vérifier la capacité de généralisation de notre modèle complet. Enfin, nous conclurons sur les résultats de recherche présentés dans ce mémoire au cours du septième et dernier chapitre.

CHAPITRE 2 REVUE DE LITTÉRATURE

Ce chapitre présente d’abord les définitions et concepts fondamentaux relatifs à notre sujet au cours de la section 2.1. Une revue de littérature systématique traitant de la prédiction de TT de HT, établie sous forme de protocole standardisé et répliquable, est ensuite dressée à la section 2.2. L’analyse des documents obtenus se fera au cours de la section 2.3. Ultimement, une revue d’analyse critique aura pour but d’examiner la qualité des documents pour déterminer la validité des résultats qui y sont présentés à la section 2.4. Elle permettra également de pleinement situer les travaux de recherche du présent mémoire dans le courant de la recherche antérieure et contemporaine dans ce domaine précis. Enfin, nous concluons.

2.1 Définitions et concepts de base

2.1.1 Mines

Selon l’Office québécois de la langue française, on pourra définir une mine comme une « zone où l’on exploite des substances utiles (autres que des matériaux rocheux) sous forme de gisement ou de filon, soit à ciel ouvert, soit par puits et galeries » [8]. Dans la suite du présent mémoire, on prendra donc bien soin de différencier les mines à ciel ouvert et les mines souterraines, dont les contextes opérationnels respectifs peuvent faire amplement varier les problématiques rencontrées.

2.1.2 Camions de transport de minerai

Les HT, qui peuvent aussi être évoqués par les termes anglais « *haulers* » (littéralement « transporteurs ») et parfois « *dump trucks* » ou « *dumpers* » (« tombereaux » ou « camions à benne basculante ») [9, 10, 11], sont couramment employés dans les mines de surface et dans les mines souterraines pour transporter la roche fragmentée. Ils sont aussi particulièrement utilisés pour transporter le mort-terrain devant obligatoirement être extraits pour permettre le développement de la mine.

Ils doivent être distingués des chargeuses sur pneus à godet, que l’on peut aussi rencontrer sous les termes anglais suivants : « *Load-Haul-Dump vehicles/loaders* » (LHD), « *loaders* », « *scooptrams* » ou encore « *muckers* » [12, 13]. Ces véhicules sont abondamment utilisés dans les mines souterraines pour déblayer efficacement les roches et charger les HT. Les LHD sont en revanche nettement moins performants pour le transport intensif de minerai sur des distances appréciables que les HT, spécialisés pour cette tâche et qui sont donc au centre du

sujet du présent mémoire.

2.1.3 Trajets des HT

Lors du transport de minerai, les trajets effectués par les HT relient une zone d'excavation active à une zone de déchargement (chute à minerai, aire de stockage ou concasseur) avec parfois un dénivelé positif considérable [14] lorsque le trajet implique une rampe. Une rampe est une route pentue généralement unique et pouvant être empruntée dans les deux sens de circulation reliant les différents niveaux d'une mine souterraine et parfois la surface. Naturellement, les HT font aussi le trajet inverse, en descente, habituellement à vide. Les camions chargés, qui montent, ont habituellement la priorité de passage [15] dans les rampes. Ce dernier point est surtout important en souterrain, car les tunnels offrent un espace de croisement et de manœuvre bien plus étroit que les larges routes pouvant être construites dans les mines de surface. Des feux de signalisation sont parfois installés pour indiquer à un éventuel HT descendant qu'un autre HT remonte ; le premier se stationne alors dans une cavité creusée dans les parois spécialement pour cette situation (appelée « chambre »), laissant au second l'espace suffisant pour passer, après quoi il peut reprendre son trajet.

Dans le cas de mines souterraines dont les rampes débouchent en surface, les HT changent notablement d'environnement durant chacun de leur trajet : lors d'une remontée classique, ils franchissent d'abord une succession de tunnels horizontaux rectilignes appelés « tunnels d'exploitation » ou « *drifts* » puis remontent par une rampe, habituellement spiralee et partagée avec tous les autres véhicules de la mine, avant de déboucher en plein air et d'effectuer le reste de leur trajet en surface jusqu'à la zone de traitement du minerai [15]. Les conditions opérationnelles de ces sections de trajets sont notablement variées. On peut citer en particulier des problématiques liées aux conditions météorologiques en surface, qui peuvent aussi modifier l'état de la chaussée, tandis qu'en souterrain la visibilité limitée rend plus délicat l'évitement de piétons et de véhicules. Enfin, les trajets des HT peuvent être déviés par de nécessaires ravitaillements de carburant, mais aussi par d'autres détours éventuels par un garage ou par un refuge (dans le cas des mines souterraines) dépendamment des circonstances et des besoins du conducteur.

Il convient de mentionner que, dans le présent mémoire, le terme « trajet » désignera le déplacement d'un HT d'un point à un autre (i.e. un aller simple) ; tandis que le terme « itinéraire » désignera un tracé géographique orienté reliant ces points. Un nombre illimité de trajets peuvent donc être observés sur un même itinéraire.

2.1.4 TT

Les TT peuvent amplement varier, même sur un trajet identique [15]. En effet, les nombreuses variations de contexte opérationnel évoquées dans la sous-section 2.1.3 sont naturellement à l'origine d'importantes variations de TT. De surcroît, le modèle et l'état du HT utilisé ainsi que le niveau de compétence du conducteur, son type de conduite et sa présence mentale (vigilance, niveau de fatigue, etc.) peuvent aussi avoir une influence considérable sur les TT.

Pour mesurer les TT dans les mines de surface, on utilise habituellement le système de positionnement par satellites GPS pour repérer les déplacements des HT et en déduire automatiquement la durée de chacun des trajets reliant deux zones pré-établies entre elles [16]. En revanche, dans le cas des mines souterraines, la tâche est bien plus ardue puisque le GPS n'est pas disponible en souterrain. Bien qu'il soit possible de procéder manuellement pour y mesurer les TT, par chronométrage, il devient nettement plus efficace de nos jours d'utiliser les informations collectées par certains capteurs nouvellement implantés dans de nombreuses mines souterraines, en particulier les balises de localisation souterraines [4]. Il va sans dire que les données générées automatiquement par les capteurs peuvent être stockées continuellement et durablement dans des Bases De Données (BDD), ce qui permet de les utiliser ultérieurement afin de retrouver théoriquement la totalité des TT de HT ayant eu lieu durant une période d'activités donnée sur une section donnée de la mine.

2.1.5 Prédiction de TT

La prédiction de TT est rendue naturellement plus difficile par les multiples incertitudes liées à chaque trajet et aux amples variations de TT qu'elles provoquent.

En comparant les critères de performances adaptés, on pourra juger de la qualité d'un modèle de prédiction en comparaison d'un autre.

2.2 Protocole de recherche d'articles scientifiques

Cette section vise à présenter l'élaboration d'une stratégie de recherche destinée à être appliquée, d'une part, à la BDD de publications d'articles scientifiques Scopus et d'autre part, au moteur de recherche scientifique Google Scholar. Cette stratégie doit nous permettre de rassembler les articles les plus pertinents concernant notre sujet.

En décomposant notre sujet de recherche, on distingue cinq concepts-clés, qui correspondent respectivement aux cinq sous-sections de la section 2.1 : les **mines**, les **HT** utilisés dans ces mines, les **trajets** qu'ils effectuent, la **durée** de ces derniers, et la **prédiction** de ces durées.

On liste alors tous les mots et expressions permettant d'exprimer pertinemment chacun de ces concepts en anglais pour ériger notre plan de concepts. On représente ce dernier dans le tableau 2.1 (les mots-clés « *Productivity* », « *Distance** » et « *Cycle** » ont été ajoutés à posteriori comme expliqué ci-après).

TABLEAU 2.1 Plan de concepts

Concept n°1	Concept n°2	Concept n°3	Concept n°4	Concept n°5
Site minier	Prédiction	Durée	Trajet	HT
Mine Mines	Predict* Forecast*	Time Times Timing Duration* Productivity	Route* Course* Itinerar* Road* Track Tracks Travel Travels Distance* Haul Hauls Cycle*	Truck* Hauler* HT Dumper*

On recherchera les occurrences des mots-clés de notre plan de concepts exclusivement dans le titre, le résumé et les mots-clés des articles répertoriés dans les résultats Scopus pour s'assurer d'une meilleure pertinence des articles qui seront rencontrés.

On décide aussi de limiter la langue des articles recherchés à la seule langue anglaise.

Par ailleurs, la découverte d'articles pouvant être considérés comme pertinents hors de Scopus, écrits en anglais et qui n'auraient pas été détectés sur la base de notre plan de concepts initial nous amène à ajouter à ce dernier les mots-clés suivants :

- « *Cycle** » (« Cycle(s) » en français) : appliqué aux HT, ce terme désigne la succession des six étapes classiques du transport de minerai par HT : le trajet à vide, une attente éventuelle à la zone de chargement, la phase de chargement, le trajet chargé, une attente éventuelle à la zone de déchargement et la phase de déchargement. Il évoque ainsi la notion de trajet, sans directement la mentionner, dans les résumés de nombreux articles [17] ;

- « *Productivity* » (« Productivité » en français) : ce terme remplace parfois la notion de durée dans le résumé des articles. On cherche en effet à prédire dans certains articles la productivité horaire des HT en tonnes de minerai extraites par période temporelle, sur la base de la prédiction de TT de HT (mentionnés uniquement dans le corps de l'article, et non dans leur titre, leur résumé ou leurs mots-clés) ; et
- « *Distance** » (« Distance(s) » en français) : ce terme peut parfois exprimer la notion de trajet quelconque dans les résumés d'articles bien qu'il soit nettement moins précis que les autres mots-clés mentionnés dans la quatrième colonne du tableau 2.1.

On choisit enfin d'exclure les six domaines d'études suivants : « *Medicine* », « *Biochemistry, Genetics and Molecular Biology* », « *Agricultural and Biological Sciences* », « *Social Sciences* », « *Environmental Science* » et « *Energy* ».

En combinant l'ensemble des éléments précédents, on obtient la requête suivante :

« **TITLE-ABS-KEY ((mine OR mines) AND (predict* OR forecast*) AND (time OR times OR timing OR duration* OR productivity) AND (route* OR course* OR itinerar* OR road* OR track OR tracks OR travel OR travels OR distance* OR haul OR hauls OR cycle*) AND (truck* OR hauler* OR ht OR dumper*)) AND (LIMIT-TO (LANGUAGE , "English")) AND (EXCLUDE (SUBJAREA , "MEDI")) AND (EXCLUDE (SUBJAREA , "BIOC")) AND (EXCLUDE (SUBJAREA , "AGRI")) AND (EXCLUDE (SUBJAREA , "SOCI")) AND (EXCLUDE (SUBJAREA , "ENVI")) AND (EXCLUDE (SUBJAREA , "ENER")) ».**

En date du **21 avril 2023**, l'exploration via **Scopus** des BDD **Compendex et Inspec** avec la requête précédente renvoyait **67 documents**.

Par ailleurs, le **4 octobre 2023**, on saisit une requête dans le moteur de recherche **Google Scholar** reprenant les mots-clés de la requête Scopus précédente. On collecte ainsi **50 documents**.

On représente à la figure 2.1 les étapes successives de notre processus de sélection d'articles sous la forme d'un organigramme de programmation.

On mentionne le nombre d'articles totalisé à chaque étape de l'organigramme par l'expression « (n=*) » où «*» désigne le nombre d'articles. Ainsi, conformément aux étapes indiquées dans cet organigramme, les 67 documents issus de Scopus trouvés précédemment sont d'abord comparés aux 50 documents issus de la recherche Google Scholar pour dépister d'éventuels doublons, ici au nombre de 14. Après avoir supprimé ces derniers, on se retrouve avec une collection de 103 documents distincts. On lit les titres et les résumés de chacun de ces docu-

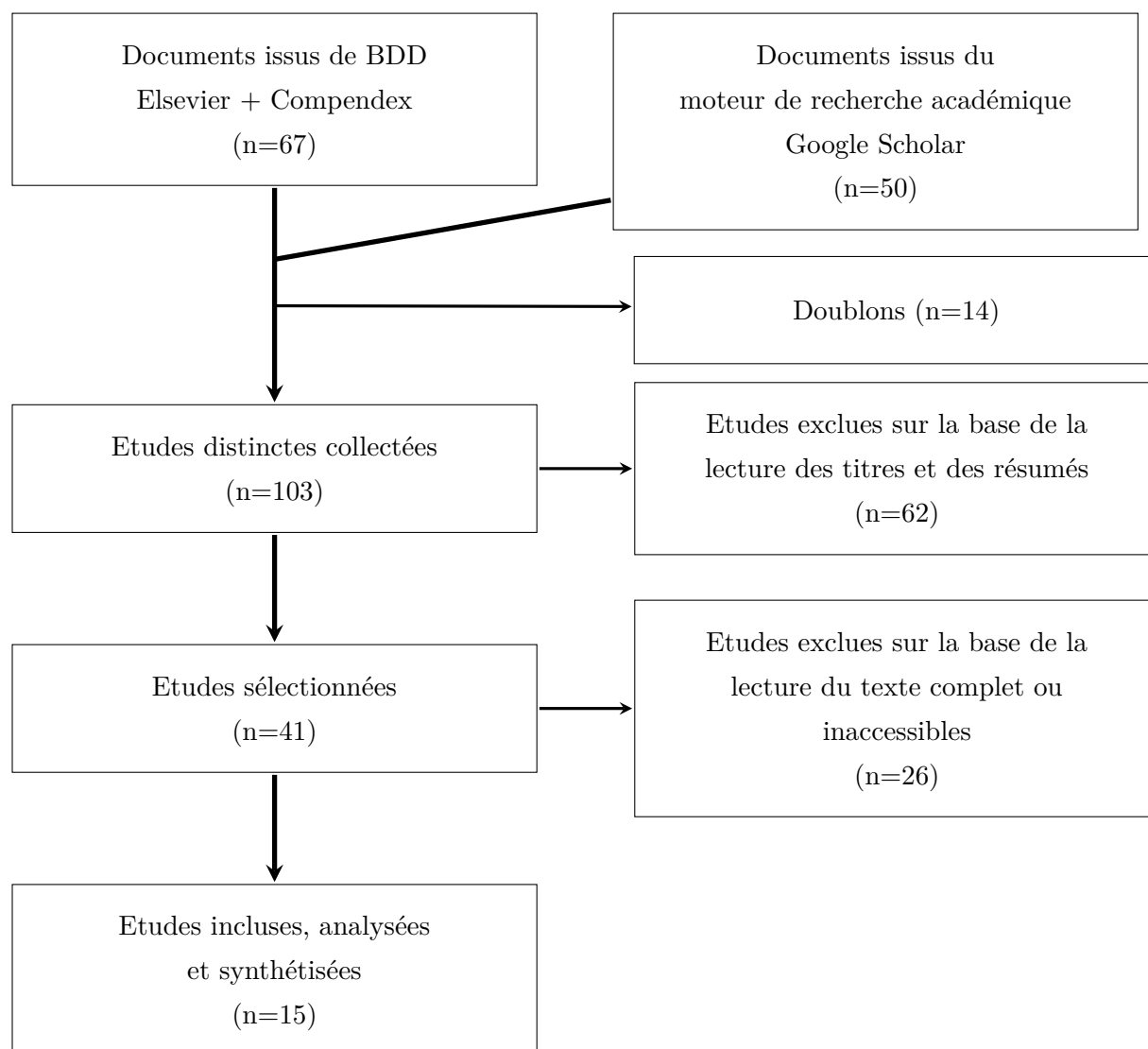


FIGURE 2.1 Protocole d'identification des articles

ments et on exclut tous ceux dont le titre ou le résumé s'éloigne clairement de notre sujet. On exclut ainsi 62 documents, soit successivement :

- 11 comptes-rendus de conférence dont la succession de titres des sujets reprend les mots-clés de la requête sans les relier entre eux ;
- quatre documents qui ne traitent pas de cas miniers ;
- un document qui ne traite pas d'une activité se déroulant directement sur un site minier ;
- 44 documents qui sont centrés sur des sujets éloignés des TT de HT ou de la productivité des activités de transport de minerai, ces documents sont chacun centré sur l'un des sujets suivants : la consommation de carburant et d'énergie, l'usure des pneus des HT,

les routes empruntées, les propriétés du matériau à excaver et à transporter, la santé des opérateurs, les problèmes techniques et les anomalies, la maintenance, le relief minier, la ventilation, les déchets miniers, la prédiction de coûts d'exploitation, le volume sonore, l'acquisition de données, l'évaluation des risques et des accidents, la modélisation de l'environnement minier ou la détection de la fatigue des opérateurs ;

- un document centré sur de la productivité clairement sans lien avec les TT ; et
- un document traitant de TT de convoyeurs et non de HT.

A la suite de ces exclusions, 41 études sont sélectionnées.

On tente alors d'accéder au texte complet de chacun de ces documents. Deux documents sont exclus en raison de la complète inaccessibilité de leur texte complet. L'un des documents est exclu car son texte est rédigé dans la langue écrite chinoise. Six documents ne sont pas centrés sur les TT, et 16 autres ne font pas mention de prédiction de TT bien qu'ils utilisent des TT. Enfin, un article est exclu, car il s'intéresse à prédire les TT des convoyeurs. Au final, 26 résultats supplémentaires sont exclus.

Notre protocole d'identification des articles nous mène donc à retenir 15 documents, qui seront analysés et synthétisés dans la section 2.3.

On présente ces documents, triés par année de publication, avec le type d'exploitation minière auxquels ils se rapportent (à ciel ouvert ou en souterrain) et la méthode de prédiction qui y est privilégiée pour la prédiction de TT (apprentissage automatique (*Machine-Learning* (ML)), simulation, modèle mathématique ou algorithme d'optimisation) dans le tableau 2.2.

2.3 Analyse des résultats

Pour examiner en détail les approches existantes dans la prédiction des TT de HT, notre analyse des résultats sera structurée en trois sous-sections principales, correspondant aux différentes techniques de prédiction de TT identifiées dans la littérature. Ces trois types de techniques sont respectivement le ML, la simulation et les autres techniques.

2.3.1 Approches basées sur des modèles de ML

Dans leurs trois articles identifiés comme pertinents par le protocole précédent, Fan et al. [18, 21, 22] cherchent à prédire la productivité horaire, directement reliée aux TT, de HT opérant dans des mines de sables bitumineux en surface. Pour ce faire, ils appliquent successivement un modèle de mélange gaussien (*Gaussian Mixture Model* (GMM)) puis un ou plusieurs modèles de ML aux données opérationnelles des HT. Les GMM sont des modèles de mélange

TABLEAU 2.2 Présentation des articles retenus (ordre antichronologique)

N°	Nom du document	Date	Exploitation	Méthode
1	Weighted ensembles of artificial neural networks based on Gaussian mixture modeling for truck productivity prediction at open-pit mines [18]	2023	À ciel ouvert	ML
2	Underground mine truck travel time prediction based on stacking integrated learning [15]	2023	Souterraine	ML
3	Prediction Method of Truck Travel Time in Open Pit Mines Based on LSTM Model [19]	2023	À ciel ouvert	ML
4	Coordinating Multiple Cooperative Vehicle Trajectories on Shared Road Networks [20]	2023	À ciel ouvert	Optimisation
5	Prediction of truck productivity at mine sites using tree-based ensemble models combined with Gaussian mixture modelling [21]	2022	À ciel ouvert	ML
6	Preprocessing Large Datasets Using Gaussian Mixture Modelling to Improve Prediction Accuracy of Truck Productivity at Mine Sites [22]	2022	À ciel ouvert	ML
7	Use of Machine Learning Algorithm Models to Optimize the Fleet Management System in Opencast Mines [23]	2022	À ciel ouvert	ML
8	A Simulation Model for Estimation of Mine Haulage Fleet Productivity [24]	2020	À ciel ouvert	Simulation
9	Simulation of truck haulage operations in an underground mine using big data from an ICT-based mine safety management system [25]	2019	Souterraine	Simulation
10	The Use of a Machine Learning Method to Predict the Real-Time Link Travel Time of Open-Pit Trucks [26]	2018	À ciel ouvert	ML
11	Dispatch with Confidence : Integration of Machine Learning, Optimization and Simulation for Open Pit Mines [27]	2017	À ciel ouvert	ML
12	Simulation and optimization in open pit mining [28]	2015	À ciel ouvert	Simulation
13	A comparative study of truck cycle time prediction methods in open-pit mining [17]	2010	À ciel ouvert	ML & simulation
14	Modelling performance and retarder chart of off-highway trucks by cubic splines for cycle time estimation [29]	2005	À ciel ouvert	Modèle mathématique
15	A computerized model for truck dispatching in open pit mines [30]	1997	À ciel ouvert	ML

visant à approximer une distribution observée par un mélange de plusieurs lois normales latentes, dans des proportions calculées par le modèle et dont les paramètres sont eux aussi calculés par le modèle [31]. Dans chacun des trois articles, l'application préalable d'un GMM aux données vise à lui apprendre à partitionner les données, qui correspondent aux conditions opérationnelles de chaque cycle de transport, en trois classes qui permettraient d'expliquer partiellement les basses, moyennes et hautes productivités horaires. Un ou plusieurs modèles de ML sont ensuite entraînés à prédire la productivité horaire en fonction des conditions opérationnelles et selon la classe latente pouvant être pré-identifiée par le GMM d'après ces conditions opérationnelles. Pour évaluer la qualité des prédictions de ces combinaisons de modèles, Fan et al. comparent les coefficients de détermination (R^2) ou les R^2 ajustés des modèles à ceux de modèles de ML entraînés sans l'aide d'un GMM. Ainsi, dans leur premier article [22], les auteurs constatent que le seul modèle prédictif qu'ils utilisent, le modèle de régression linéaire multiple (*Multiple Linear Regression* (MLR)), obtient un R^2 ajusté de 0,23 sans GMM et 0,75 avec GMM. Dans leur deuxième article [21], ils comparent les performances de modèles d'arbres de décision (*Decision Trees* (DT)), de forêts aléatoires (*Random Forest* (RF)), de MLR et de régression par *boosting* du gradient (*Gradient Boosting Regression* (GBR)) tous avec ou sans GMM et constatent que le modèle qui obtient la meilleure qualité de prédiction avec et sans GMM parmi les quatre modèles testés est le RF, qui obtient un R^2 d'environ 0,48 sans GMM et qui atteint un R^2 élevé d'environ 0,90 avec GMM. Enfin, dans leur troisième article [18], Fan et al. utilisent sans GMM quatre modèles prédictifs similaires à ceux utilisés dans leur deuxième article (ici DT, RF, GBR et un modèle de *boosting* extrême du gradient (*eXtreme Gradient Boosting* (XGBoost))), et comparent leurs résultats avec ceux de trois modèles prédictifs avec GMM basés sur les réseaux de neurones artificiels (*Artificial Neural Networks* (ANNs)) : un modèle de réseau de neurones régularisé bayésien (*Bayesian Regularized Neural Network* (BRNN)), un modèle de réseau de neurones à rétropropagation du gradient (*BackPropagation Neural Network* (BPNN)) et un modèle d'apprentissage automatique extrême (*Extreme Learning Machine* (ELM)) ; ils observent que le modèle sans GMM qui obtient le meilleur R^2 est XGBoost) avec un R^2 de près de 0,42 tandis que le meilleur modèle avec GMM est le BRNN avec un R^2 d'environ 0,86. Dans les trois articles, Fan et al. concluent à une qualité de prédiction de la productivité horaire des HT considérablement accrue par l'utilisation d'un GMM en amont de l'utilisation des modèles prédictifs proposés.

Li et al. [15] ont quant à eux cherché à démontrer la pertinence de l'apprentissage intégré par empilement pour prédire avec précision la durée des cycles de transport de HT dans une mine souterraine à accès par rampe. Pour ce faire, dans les cycles de transport des HT reliant un dépôt de minerai en surface à un unique niveau de la mine, ils ont différencié un total de trois

sections de trajet (i.e. la surface, la rampe et les galeries de développement/exploitation) et deux situations de chargement de HT différentes (i.e. HT chargé et vide), soit un total de six sections de trajets différentes par cycle complet de transport. Ils ont alors collecté les données opérationnelles correspondant à 200 trajets distincts sur chacune de ces sections. Ils ont ensuite entraîné trois modèles prédictifs i.e. un séparateur à vaste marge (*Support Vector Machine* (SVM)) par les moindres carrés, un modèle d'apprentissage automatique léger de *boosting* du gradient (*Light Gradient Boosting Machine* (LightGBM)) et un algorithme de RF à prédire le TT des HT sur ces six types de sections de trajet différentes, avec des variables d'entrée adaptées à chacun de ces environnements. Pour recombinaison les six sections de trajet constituant chacun des cycles de transport, Li et al. ont commencé par empiler les prédictions des trois modèles précédents concernant le TT des HT sur chacune des six sections constituant chacun des cycles de transport. Ils ont ensuite fait prédire à leur méta-modèle i.e. le modèle XGBoost, la durée totale du cycle de transport d'après les six prédictions précédentes. Les prédictions de cette combinaison de modèles sur les données de test sont associées à un pourcentage d'écart absolu moyen (*Mean Absolute Percentage Error* (MAPE)) de 2,31% pour les galeries souterraines, de 4,39% pour la rampe et de 4,56% pour la surface. Li et al. comparent alors ces résultats de prédiction avec ceux obtenus par des modèles seuls parmi ceux cités ci-avant ou des combinaisons moins sophistiquées de ces modèles. Ils détaillent particulièrement les résultats du modèle XGBoost seul, dont les résultats de prédiction sont associés à des MAPE tous supérieures à celles du XGBoost seul : 2,72% pour les galeries souterraines, 5,22% pour la rampe et 4,67% pour la surface. Li et al. concluent alors que ces résultats démontrent la meilleure performance de prédiction du modèle d'apprentissage intégré par empilement, combinaison des trois modèles de base et du méta-modèle.

Choudhury et Naik [23] se sont concentrés sur l'amélioration de la productivité et la réduction des coûts opérationnels dans une mine à ciel ouvert en minimisant le nombre total de HT à utiliser dans cette mine ; cet objectif leur a demandé de prédire en amont les TT de HT. Tout d'abord, ils ont segmenté le cycle de transport des HT en phases distinctes i.e. les temps de chargement/déchargement, les temps de positionnement et les TT. Ils ont ensuite entraîné un algorithme de RF, un algorithme des k -plus proches voisins (k -Nearest Neighbours (kNN)) et un SVM à prédire les TT des HT pour un unique itinéraire. L'algorithme de RF ayant surpassé les deux autres modèles prédictifs sur cet itinéraire avec un MAPE de 0,06% comparativement aux MAPE de 3,47% pour le SVM et 4,30% pour le kNN, il est sélectionné pour faire les prédictions de TT sur deux autres itinéraires, pour lesquels il obtient respectivement un MAPE de 0,95% et de 1,1%. Les résultats de prédiction de l'algorithme de RF sont utilisés dans la fonction de minimisation du nombre total de camions. Celle-ci liste le nombre de camions à affecter à chacun des trajets ainsi qu'aux zones de chargement et de

déchargement pour une répartition optimale. Choudhury et Naik concluent entre autres que les résultats expérimentaux obtenus démontrent l'intérêt de l'utilisation du ML pour prédire les TT avant d'utiliser ces derniers pour optimiser le problème de répartition des HT.

Sun et al. [26] ont développé une méthode visant à prédire les TT par section des HT dans les mines à ciel ouvert. Dans ce but, ils divisent les trajets en sections fixes i.e. des sections de trajet qui resteront identiques au cours du développement de la mine, et en sections provisoires, lesquelles seront modifiées. Ils entraînent ensuite trois modèles de prédiction différents i.e. un algorithme de kNN, un SVM, et un algorithme de RF. La comparaison de la déviation absolue moyenne (*Mean Absolute Deviation* (MAD)) et du MAPE des résultats de prédiction de TT de ces trois modèles prédictifs sur trois sections fixes et trois sections temporaires montre que les modèles de SVM et de RF sont toujours meilleurs que l'algorithme de kNN. Sun et al. comparent alors les meilleurs résultats obtenus pour chacune des sections aux méthodes traditionnelles de calcul de moyenne des TT pour chaque section, et une réduction considérable du MAPE est observée puisqu'elle diminue en moyenne de 12,54 points de pourcentage pour les sections fixes et en moyenne de 19,30 points de pourcentage pour les sections temporaires. De plus, l'intégration de caractéristiques météorologiques dans les modèles par les auteurs s'est aussi traduite par une amélioration des prédictions, diminuant encore le MAPE des modèles prédictifs de 5,13 points de pourcentage toutes sections confondues. La prédiction de TT par sections de trajet, en choisissant la combinaison réputée optimale de modèles prédictifs pour chaque section, est enfin comparée par les auteurs à la prédiction de TT sur le trajet tout entier (non sectionné) par le SVM et par le RF. La MAD et le MAPE de la prédiction par sections sont inférieures aux MAD et aux MAPE respectifs des prédictions sur le trajet tout entier, en particulier le MAPE est de 20,8% pour la prédiction du TT complet par le SVM, de 23,5% pour la prédiction du TT complet par le RF et de 9,0% pour la prédiction du TT par sections par l'assemblage réputé optimal de RF et SVM. Sun et al. concluent à une meilleure efficacité des prédictions de TT par des modèles prédictifs comparément aux méthodes de calcul de moyenne de TT, ils soulignent les performances supérieures des modèles prédictifs de SVM et de RF par rapport au kNN, mais aussi l'amélioration des résultats de prédiction grâce à l'utilisation de données météorologiques et grâce à la segmentation de trajets permettant la prédiction par section.

Chanda et Gardiner [17] ont mené une étude comparative sur les méthodes de prédiction du temps de cycle complet des camions dans les mines à ciel ouvert afin de déterminer la meilleure méthode selon eux. Le cycle complet inclut, sans s'y restreindre, le chargement des camions, leurs trajets aller et retour, leur déchargement, mais aussi leurs temps de positionnement et leurs temps d'attente. Les auteurs ont testé et comparé les performances de trois méthodes différentes : la simulation, un ANN et la MLR. Leur simulation se déroule via le logiciel

TALPAC, basé sur une méthode de simulation de Monte-Carlo [32] et couramment utilisé selon eux par l'industrie minière australienne pour simuler les activités minières incluant des interactions entre HT et chargeuses. La simulation vise ici à estimer le temps de cycle des HT. Les auteurs fournissent au logiciel et aux deux modèles de ML les données d'entrée supposées pertinentes pour chacun d'eux dont, pour les modèles prédictifs, les temps de chargement des HT, l'existence d'une file d'attente à la pelle ou d'une attente quelconque à la pelle et, pour le logiciel de simulation, la durée estimée de positionnement des camions et le temps de chargement des HT. Chanda et Gardiner obtiennent ainsi les prédictions de temps de cycle de chacune de ces méthodes pour cinq itinéraires à comparer, puis calculent le pourcentage d'erreur des prédictions pour chacun de ces cas avant de les représenter collectivement sur différents graphiques. Ils concluent que les modèles d'ANN et de MLR sont plus précis que le logiciel de simulation TALPAC quel que soit le trajet étudié. Selon eux, les résultats obtenus montrent clairement que le logiciel sous-estime la durée des cycles courts et sur-estime la durée des cycles longs.

Ristovski et al. [27] examinent une approche intégrée de répartition des HT basée sur le ML, l'optimisation et la simulation dans deux mines à ciel ouvert. Concernant la prédiction de TT incluse dans l'approche globale employée, c'est le ML qui est utilisé dans cet article. Le modèle employé est un modèle linéaire généralisé (*Generalized Linear Model* (GLM)) à loi de distribution Gamma. Les auteurs préparent ce modèle à différentes situations inédites en lui ajoutant des sous-modèles pour le rendre capable de réaliser ses prédictions pour de nouveaux HT dont l'identifiant et/ou les dimensions lui seraient inconnus. Ce modèle est comparé à deux modèles de référence pour la prédiction de TT que les auteurs affirment souvent utilisés dans la pratique : la prédiction via un calcul de la moyenne des TT et la prédiction par marche aléatoire qui consiste à prédire que le prochain TT devrait être identique au dernier TT observé. Ces modèles de référence disposent aussi de sous-modèles sur lesquels se reposer en cas de données inédites. Les résultats de prédiction des différents modèles testés sont comparés en calculant leur racine carrée de l'écart quadratique moyen (*Root Mean Square Error* (RMSE)). Ce RMSE est de 49,7 secondes pour le GLM, de 81,7 secondes pour le modèle de marche aléatoire et de 62 secondes pour le modèle de calcul de moyennes des TT. Le GLM a par ailleurs pu prédire 100% des TT durant la phase de test, alors que les modèles basés sur la marche aléatoire et les calculs de moyenne ont été seulement capables d'en prédire 95%, en raison de données insuffisantes liées à des itinéraires inédits. Les auteurs concluent que le GLM a une meilleure précision que les autres modèles, et qu'il constitue aussi un meilleur choix dans le contexte de leur étude du fait de sa flexibilité face à des scénarios inédits.

Ao, Li et Yang [19] prédisent les TT des HT dans les mines à ciel ouvert en utilisant un modèle

de réseau de neurones à longue mémoire à court terme (*Long Short-Term Memory* (LSTM)). Il s'agit d'une classe de réseau de neurones récurrent (*Recurrent Neural Networks* (RNN)). Ce modèle dispose de 11 variables d'entrée variées concernant des caractéristiques des HT, des caractéristiques des itinéraires, le style de conduite de l'opérateur et des conditions météorologiques. Ils entraînent et comparent les performances de ce modèle à celles de modèles de référence i.e. un réseau de neurones à rétropropagation du gradient (*BackPropagation Neural Network* (BPNN)) et un SVM de régression, et obtiennent un écart absolu moyen (*Mean Absolute Error* (MAE)) de 0,85 secondes et un MAPE de 2,6% pour le LSTM, un MAE de 2,86 secondes et un MAPE de 11,6% pour le BPNN, et un MAE de 3,58 secondes et un MAPE de 16,2% pour le SVR. Ces résultats permettent aux auteurs de constater une amélioration significative de la précision des prédictions de TT de HT par le LSTM à comparer aux autres méthodes, qualifiées de traditionnelles, dans les mines à ciel ouvert. Les auteurs concluent aussi que la prise en compte des données météorologiques, du style de conduite de l'opérateur et du chargement ou non du HT peut permettre d'améliorer la qualité des prédictions. Ils suggèrent enfin une amélioration potentielle consistant à prendre en compte l'intégralité des informations du réseau de la mine et non seulement les données opérationnelles liées au trajet dont on tente de prédire le TT.

Dans sa thèse, Temeng [30] propose un modèle de répartition des HT dans les mines à ciel ouvert en trois étapes, qui visent respectivement à minimiser les TT des HT, maximiser la production et maintenir la qualité du minerai. Concernant la prédiction de TT, un programme de simulation de véhicules commence par générer des TT pour de très nombreuses combinaisons de capacité de HT, de distance, de pente et de résistance au roulement différentes grâce aux données issues de graphes de performance et de ralentissement. La vitesse maximale des HT est fixée à 35 miles par heure pour la simulation. Un MLR est alors entraîné à prédire les TT obtenus par la simulation précédente. Cette prédiction se fait en fonction de différentes combinaisons des variables de distance, de résistance totale et de leurs interactions ; l'objectif étant de voir si le modèle gagne en robustesse grâce à toutes ces données. Une analyse en composantes principales permet de trouver les variables les plus influentes pour ce modèle. Temeng analyse les résultats de prédiction de son MLR en les comparant aux résultats des simulations basées sur les graphes de performance et de ralentissement. Il conclut que le MLR prédit de manière satisfaisante les TT générés par les simulations.

2.3.2 Approches basées sur des modèles de simulation

Baek et Choi [25] ont simulé des opérations de transport de minerai par HT dans une mine souterraine en s'appuyant sur trois mois de données opérationnelles issues d'un système de

reconnaissance des HT par balises. Ces balises sont essentiellement réparties dans le réseau minier souterrain, et détectent les étiquettes de reconnaissance des HT puis envoient l’horodatage de la détection et l’identifiant du HT à un serveur qui identifie la balise émettrice. En analysant la succession de balises ayant reconnu un HT en particulier ainsi que les horodatages, les auteurs sont alors capables de déterminer les TT réels pour chaque segment inter-balises et ce dans les deux sens. Baek et Choi en déduisent le TT moyen et l’écart-type de la distribution des TT historiques, qui leur servent de paramètres d’entrée pour leurs simulations subséquentes. En effet, les auteurs ont développé un modèle de simulation par événements discrets se basant sur ces grandeurs statistiques pour simuler les trajets et ainsi obtenir une prédiction de la productivité des HT. Leur simulation n’utilise aucune autre variable issue directement ou dérivée des données de télémétrie de la mine, elle utilise uniquement des variables contextuelles comme le nombre d’heures de travail journalier ou la charge utile des HT. Ils ont comparé les résultats de ces simulations aux résultats opérationnels relevés sur deux journées de travail réelles, relevant une haute précision de leur modèle de simulation.

Dans leur premier document identifié comme pertinent par le protocole précédent, Upadhyay et al. [28] explorent l’utilisation d’un modèle de simulation d’événements discrets pour une mine à ciel ouvert en se concentrant particulièrement sur la simulation du transport par HT. Les auteurs présentent d’abord un modèle de programmation linéaire mixte (MILP) qui doit fournir à un modèle de simulation, en aval, la répartition quasi-optimale des HT aux pelles afin de relier les objectifs opérationnels et la planification à court terme des tâches d’extraction. Les auteurs ont ensuite utilisé le logiciel de simulation Arena de Rockwell Automation sur les données réelles d’une mine de sables bitumineux à ciel ouvert en lui fournissant aussi les résultats du MILP pour obtenir des prévisions de TT sous forme d’histogrammes. La validation du modèle de simulation s’est basée sur la comparaison de nombreux indicateurs clés de performance (*Key Performance Indicators* (KPI)) avec des valeurs réelles, ainsi que sur l’allure de graphes de quantiles normalisés illustrant les TT des HT à vide et chargés et leurs vitesses de déplacement. D’après ces résultats, les auteurs confirment l’applicabilité du sous-modèle de simulation de transport pour modéliser les TT des HT. Ils évoquent par ailleurs la nécessité de combiner le MILP au processus de répartition du modèle de simulation pour vérifier qu’il reste applicable dans les simulations à plus long terme.

Dans leur second document, Upadhyay et al. [24] présentent un modèle de simulation de Monte-Carlo pour estimer la productivité horaire d’une flotte de HT dans les opérations minières à ciel ouvert. Dans un premier temps, ils utilisent la carte du réseau routier de la mine étudiée pour déterminer les chemins les plus courts entre les sites de chargement et de déchargement fixés par la séquence d’extraction du calendrier de production de cette mine. Ils

utilisent ensuite les distributions empiriques ajustées de la vitesse des HT en charge et à vide, la capacité de traction des HT ainsi que la pente et la résistance au roulement des segments d'itinéraires empruntés pour en déduire une distribution de probabilité théorique des TT. Ils lancent ensuite leur simulation de Monte-Carlo sur les itinéraires prédéterminés avec pour variables d'entrée le tonnage de minerai à transférer (issu du calendrier de production), la distribution théorique des TT et la distribution observée des temps de chargement et de déchargement. Les performances du modèle sont validées en comparant les résultats simulés à l'historique réel d'une mine de sables bitumineux du nord de l'Alberta au Canada. Les auteurs concluent que le modèle de simulation de Monte-Carlo est plus précis que les méthodes existantes pour estimer la productivité horaire dans un contexte similaire.

2.3.3 Autres approches

Pour sa part, Erarslan [29] a tâché d'estimer le temps de cycle des camions dans les mines à ciel ouvert via la généralisation de graphes de performance et de ralentissement à n'importe quel profil de route. Ces graphes sont ici fournis par le constructeur des différents modèles de HT et donnent la vitesse théorique à laquelle chacun des modèles de HT peut opérer selon la résistance totale rencontrée, qui est la somme de la résistance de la pente et de la résistance au roulement rencontrées. Les données obtenues de ces graphes sont traitées par une méthode d'interpolation par splines cubiques pour modéliser la résistance totale et la vitesse correspondante des HT pour des profils de route spécifiques, quelle que soit la résistance totale rencontrée. La spline cubique est une fonction définie par morceaux de fonctions cubiques connectés entre eux à des points de raccord prédéfinis par des contraintes de continuité allant jusqu'au second degré de dérivation (courbure) [33]. La prédiction de la vitesse des HT sur les profils de route rencontrés sur leur trajet permet ensuite de calculer leurs TT aller et retour d'après la distance parcourue sur chacun des profils de route rencontrés pour n'importe quel modèle de HT. L'auteur ajoute ensuite à ces TT les autres durées incluses dans un cycle complet de transport de minerai pour obtenir la durée totale d'un tel cycle en fonction du modèle de HT étudié. Enfin, il déduit de ces temps de cycle un nombre optimal de HT à allouer à chaque excavatrice selon le modèle de HT. L'utilisation d'un logiciel de simulation et la comparaison des résultats obtenus lui permet de confirmer que son approche est valide et qu'elle permet effectivement de déterminer le nombre de HT à utiliser. Erarslan conclut par ailleurs que l'utilisation de l'interpolation par spline cubique permet de généraliser les graphes de performance et de ralentissement à des conditions de route spécifiques et permet ainsi de déterminer la taille optimale de la flotte de camions pour les opérations minières.

Gun et al. [20] présentent leur approche de planification des trajectoires simultanées de HT

dans une mine à ciel ouvert. L’objectif qu’ils se fixent est de minimiser le TT total de la flotte de HT et ils utilisent pour cela un MILP. Dans ce modèle, le réseau de routes de la mine est modélisé comme un graphe orienté, où les nœuds représentent des intersections ou des lieux particuliers tandis que les arcs représentent des segments de route reliant ces nœuds. On attribue à chaque HT un point de départ et une destination et l’itinéraire le plus court est calculé via le graphe orienté. La dynamique des HT est caractérisée dans le MILP par différentes paramètres dont le poids du HT, son accélération et sa décélération maximales, sa vitesse maximale et la résistance au roulement de la route empruntée. Les interactions entre les HT sont prises en compte, et pour chaque paire de HT partageant un segment de route ou se croisant à une intersection, on modélise les points de conflit potentiels qui doivent être évités puis cet évitement se fait en déterminant l’ordre de passage optimal des HT. Pour réduire la complexité du MILP, les auteurs réduisent le nombre de variables, et utilisent des techniques itératives pour satisfaire les contraintes d’évitement des collisions au dernier moment, c’est-à-dire les ajouter uniquement lorsqu’elles sont nécessaires au fur et à mesure de l’optimisation au lieu de toutes les calculer dès le départ. Les auteurs ont testé leur approche sur un réseau de routes inspiré d’une véritable mine à ciel ouvert et ont observé une réduction du TT total de la flotte de HT par rapport à leur simulation d’une approche réactive basée sur des pratiques de conduite supposées courantes. Cette réduction est particulièrement importante lorsque le nombre de HT actifs simultanément est élevé : dans les cas les plus extrêmes où 60 HT coopèrent, les TT total de flotte sont d’environ 17 000 secondes pour le MILP et d’environ 36 000 secondes pour l’approche réactive. Les auteurs ont par ailleurs observé que l’application de leur MILP demandait un temps de calcul total de 120 secondes pour les problèmes les plus complexes, i.e. avec 60 HT simultanément actifs.

2.4 Revue critique

L’état de l’art des travaux de recherche portant sur la prédiction de TT de HT et de grandeurs analogues ont montré l’existence d’une grande variété de méthodologies et de modèles, mais aussi d’une grande variété d’objectifs de prédiction possibles. Aussi, il convient de faire une analyse globale et critique des 14 documents retenus.

Parmi les constatations les plus flagrantes figure sûrement la très grande majorité d’articles traitant de mines à ciel-ouvert. On ne trouve en effet que deux articles [15, 25] qui cherchent à prédire des TT de HT dans les mines souterraines. Cette faible représentation des documents portant sur les mines souterraines pourrait être due aux difficultés de prédiction évoquées par les auteurs de ces articles et inhérentes au contexte minier souterrain comme la forte instabilité et hétérogénéité des itinéraires de transport, la complexité et la faible visibilité

du réseau routier qui peuvent fortement modifier le comportement des conducteurs, et les manœuvres de croisement des HT dans les galeries étroites. On pourra aussi citer l'absence de GPS comme mentionné à la section 2.1.4 qui complexifie immanquablement la collecte de données historiques de référence sur les TT en souterrains. Ces données sont pourtant nécessaires à l'entraînement et à la validation des modèles de ML. L'un des articles [15] ne possède d'ailleurs qu'un échantillon de 200 trajets pour entraîner et tester ses modèles prédictifs, ce qui peut très vite limiter la qualité des prédictions de ces derniers.

La seconde remarque générale concerne la nature des modèles et des méthodes choisies par les auteurs pour prédire les TT ou des grandeurs analogues. En effet, au total, une grande majorité des 14 documents, soit 10 au total font appel à des modèles de ML tandis que trois documents seulement font appel à des méthodes de simulation. Un des documents fait appel aux deux méthodes en même temps et démontre la supériorité de deux modèles de ML basiques comparativement au logiciel de simulation TALPAC, couramment utilisé dans l'industrie minière [17, 34], menant à penser que les modèles de ML devraient être privilégiés dans un contexte similaire. Deux autres techniques, l'une reposant sur un modèle mathématique et l'autre sur un algorithme d'optimisation, sont employées.

Parmi les modèles de ML, une synthèse de la revue de littérature indique que les meilleurs modèles pour la prédiction de TT de HT sur les sites miniers sont généralement, selon la situation, le XGBoost, le RF, le LSTM, le GBR ou, potentiellement, le SVM et le BRNN. Concernant ce dernier modèle, uniquement utilisé dans l'un des articles de Fan et al. [18] et qui y donne les meilleurs résultats en présence de nombreux autres modèles, il est impossible de conclure à sa supériorité puisque seuls ses résultats avec GMM sont évoqués tandis que les autres modèles non basés sur des ANNs réalisent tous leurs prédictions sans GMM, partant donc sûrement avec un handicap.

Concernant les types de prédiction, certaines méthodes estiment la productivité horaire moyenne ou sa distribution, tandis que d'autres tâchent de prédire les temps de cycle et d'autres encore les TT. Ces dernières méthodes estiment parfois la moyenne ou la distribution de probabilité des TT de différents HT sur différents itinéraires, tandis que de nombreux modèles de ML prédisent pour chaque trajet une unique valeur de TT, même dans le cas où tous les trajets prédits ont lieu sur un itinéraire identique, pour s'adapter aux variations de conditions opérationnelles. Ces TT uniques peuvent être prédits pour des itinéraires fixes au complet, pour une succession de segments d'itinéraires ou encore pour des itinéraires quelconques, même inédits, si le modèle a été conçu et préparé pour généraliser. Dans ce dernier cas, de nombreuses données d'entrée additionnelles peuvent être nécessaires.

Il semble important de noter que de nombreux articles s'appuyant sur des modèles de ML

utilisent parfois des variables de prédiction qui semblent très difficilement réutilisables par les planificateurs miniers, voire inutilisables. Parmi ces variables discutables, on peut citer en particulier le nombre de piétons évités durant le trajet, le niveau d'usure des pneus du HT et l'attente ou non du HT à la pelleteuse et au site de déchargement. La connaissance de ces variables peut bien sûr améliorer la qualité des prédictions des modèles de ML mais elles semblent bien peu exploitables concrètement, en amont du trajet. Si les modèles requiert toutes ces variables en entrée, il leur serait même impossible de formuler les mêmes prédictions que celles obtenues par les auteurs.

Malgré la grande diversité des méthodologies rencontrées dans la littérature scientifique, une lacune reste criante. Aucun des 14 articles issus de la revue de littérature systématique n'aborde substantiellement son processus de préparation des données relatives aux TT au-delà de quelques notions très basiques de normalisation des données et de retrait des valeurs aberrantes. Cette étape est pourtant habituellement rapportée comme cruciale pour l'obtention de résultats de prédiction de bonne qualité [27].

2.5 Conclusion

A travers cette revue de littérature systématique, nous avons pu explorer une diversité considérable de méthodes et de techniques relatives à la prédiction de TT de HT sur les sites miniers dans un nombre de documents pourtant faible. Pour autant, l'absence totale de méthodologie substantielle de préparation de données relatives aux TT est frappante. La très faible proportion d'articles de prédiction de TT dans les mines souterraines est notable et nombre d'articles n'ont pas rendu leur solution pleinement exploitable par les planificateurs miniers. Aussi, ces observations nous amènent à proposer une nouvelle méthodologie pour la prédiction de TT de HT dans les mines souterraines qui pourra être effectivement utilisée par les planificateurs miniers. Sa structure sera présentée au chapitre suivant.

CHAPITRE 3 MÉTHODOLOGIE DE RECHERCHE

Ce chapitre présente l’approche méthodologique adoptée dans le présent mémoire. Nous définirons dans un premier temps nos objectifs de recherche, basés sur notre revue de littérature systématique, à la section 3.1. Dans un second temps, nous spécifierons le cadre méthodologique adopté en évoquant les études de cas concrets qui ont contribué au développement et à la validation de notre méthode à la section 3.2.

3.1 Objectifs de recherche

Notre revue de littérature systématique a mis en lumière une carence significative en méthodologies dédiées à la préparation de données relatives aux TT de HT. De même, la prédiction de TT de HT dans les mines souterraines a été très peu abordée dans la littérature. Les articles rencontrés utilisaient des variables non disponibles lors de la planification à court terme. Sur cette base, les objectifs principaux de cette recherche sont les suivants :

1. Développer une méthodologie spécifique et robuste pour la préparation de données liées aux trajets de HT dans les mines souterraines afin de traiter et d’exploiter efficacement une partie des multiples capteurs fixes et embarqués qui s’y trouvent ;
2. Sur la base de ces données préparées, concevoir et valider une nouvelle méthode de prédiction performante de TT de HT adaptée aux particularités opérationnelles des mines souterraines et s’appuyant uniquement sur des variables disponibles aux planificateurs miniers en amont des quarts de travail ; et
3. Tester la généralisabilité et l’efficacité de la méthode proposée sur différents sites miniers dissemblables afin d’assurer sa robustesse et sa flexibilité face aux variations du contexte opérationnel et à de nouvelles problématiques spécifiques.

Pour évaluer les performances de notre modèle de prédiction, qui correspond à notre second objectif, nous nous appuierons essentiellement sur le RMSE. Nous comparerons ainsi ses performances avec celles du modèle de référence de l’industrie minière : la moyenne des TT historiques. Le MAE sera aussi utilisé pour tempérer notre interprétation des performances observées. Ces mêmes critères permettront de vérifier que nous satisfaisons à notre troisième objectif.

3.2 Méthodologie de recherche

Le cadre méthodologique retenu est la *Design Research Methodology* (DRM) [35], qui combine des approches empiriques et expérimentales et qui est appropriée pour aborder les problématiques d’environnements industriels complexes.

Nous articulons la DRM autour de quatre phases principales :

1. **Exploration du sujet** : traité dans notre revue de littérature, au chapitre 2, ce premier volet a permis d’identifier les lacunes des méthodes actuelles de prédiction de TT et de sélectionner les techniques de préparation de données les plus prometteuses ;
2. **Première étude descriptive** : cette phase d’exploration d’un premier site minier consiste à recueillir des données opérationnelles, à comprendre les conditions opérationnelles des HT et à identifier les besoins spécifiques de prédiction de TT sur ce site. On intègre cette étude au chapitre 4 ;
3. **Étude normative** : sur la base des connaissances acquises lors de l’étude descriptive, on développe un modèle de préparation de données (chapitre 4) puis un modèle de prédiction de TT (chapitre 5). On veille à modéliser les différents cas de figures pour lesquels notre méthodologie serait mise en échec et on en déduit les adaptations nécessaires le cas échéant. Cette phase permet d’intégrer les meilleures pratiques et innovations technologiques permettant de pallier les déficiences identifiées ; et
4. **Seconde étude descriptive** : l’application pratique du modèle sur un autre site minier souterrain permet de tester sa performance et de l’ajuster en fonction des résultats obtenus. Cette phase permet d’évaluer la capacité de généralisation du modèle à d’autres sites miniers. Elle sera présentée au chapitre 6.

3.3 Conclusion

Ce chapitre a permis de définir clairement les objectifs et la méthodologie de notre recherche, qui vise à approcher avec rigueur la préparation de données de trajets de HT dans les mines souterraines ainsi que la prédiction des TT correspondants. Le prochain chapitre présentera la première étude descriptive qui détaille le contexte opérationnel et les spécificités du premier site minier rencontré avant de préparer les données qui y auront été extraites.

CHAPITRE 4 DÉVELOPPEMENT D’UN MODÈLE DE PRÉPARATION DE DONNÉES

Le présent chapitre présente le développement de notre modèle de préparation de données. Il se constitue d’abord d’une description du cas d’étude, présentée à la section 4.1, d’une analyse des contraintes et des spécificités du cas d’étude à la section 4.2 et d’une spécification des requis du modèle de préparation de données, basée sur le cas d’étude et décrite dans la section 4.3. Ces trois premières étapes sont suivies de l’identification des variables d’intérêt et des données qui s’y rattachent à la section 4.4. Viennent ensuite de multiples étapes de préparation des données, rassemblées à la section 4.5. Sur la base des étapes précédentes, la section 4.6 présente des pistes d’amélioration concernant la collecte et la gestion des données. Pour terminer, nous prétraiterons les données préparées au cours de la section 4.7, en vue de leur utilisation par les modèles de prédiction de TT du chapitre 5.

4.1 Description du cas d’étude

La présente section vise à esquisser un portrait global du site industriel auquel on s’intéresse. Dans l’intégralité du présent chapitre, nous nous intéressons à une mine souterraine canadienne aurifère à accès par rampes et exploitée via la méthode d’abattage par longs trous ; nous l’appellerons « mine 1 ».

L’architecture globale de la mine est présentée à la figure 4.1.

Depuis la surface, où se trouve la zone de déchargement de minerai, deux rampes d’accès sont disponibles : la rampe principale ou « rampe 1 », qui dessert tous les niveaux de la mine et une rampe secondaire ou « rampe 2 », connectée directement à un niveau intermédiaire se trouvant à 300 mètres de profondeur ou « niveau 300 » (cette abréviation sera aussi appliquée aux autres niveaux), créant ainsi un accès plus direct aux niveaux les plus profonds de la mine via la partie inférieure de la rampe 1. Le niveau 300 a une importance particulière puisque le garage souterrain de cette mine y est par ailleurs implanté. Le niveau le plus profond de cette mine est le niveau 500 tandis que le niveau 125 est le plus proche de la surface. Les niveaux immédiatement successifs sont constamment espacés d’un dénivelé de 25 mètres, ce que l’on nommera par la suite « écart inter-niveaux ».

Décrivons maintenant les activités de transport de minerai qui se déroulent dans cette mine. Au total, 11 HT peuvent opérer dans l’ensemble du site minier. À travers la figure 4.2, nous présentons un aperçu des activités de transport de minerai par HT sur la rampe principale

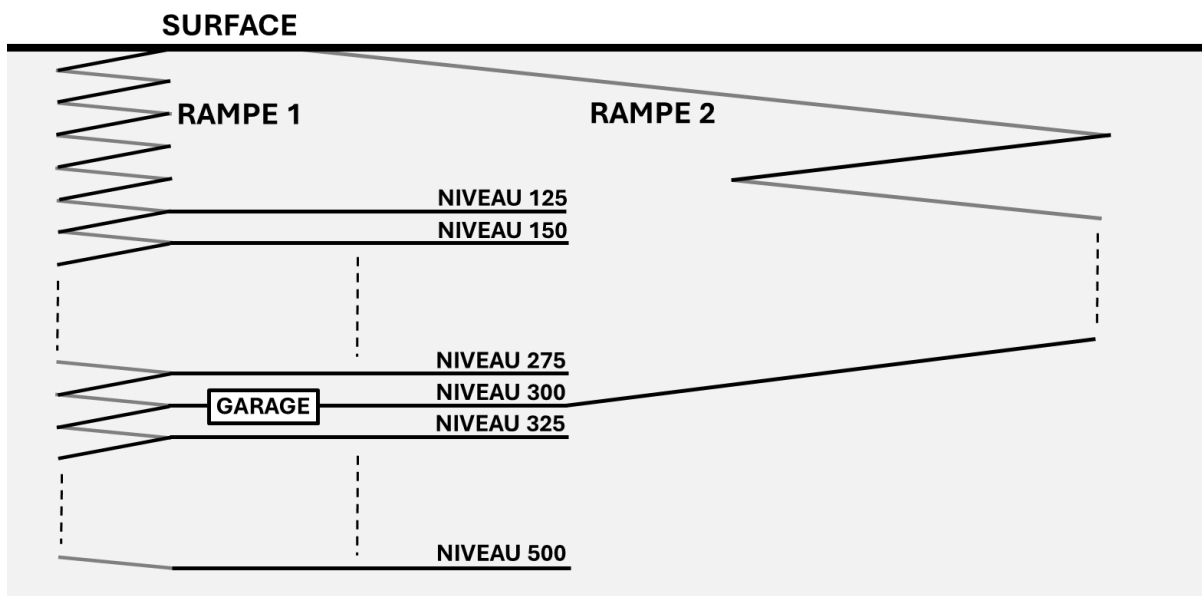


FIGURE 4.1 Architecture fondamentale de la mine 1

de la mine 1. Sur cette figure, tous les éléments de couleur verte font partie du système de connectivité de la mine 1, décrit à la section 4.2.

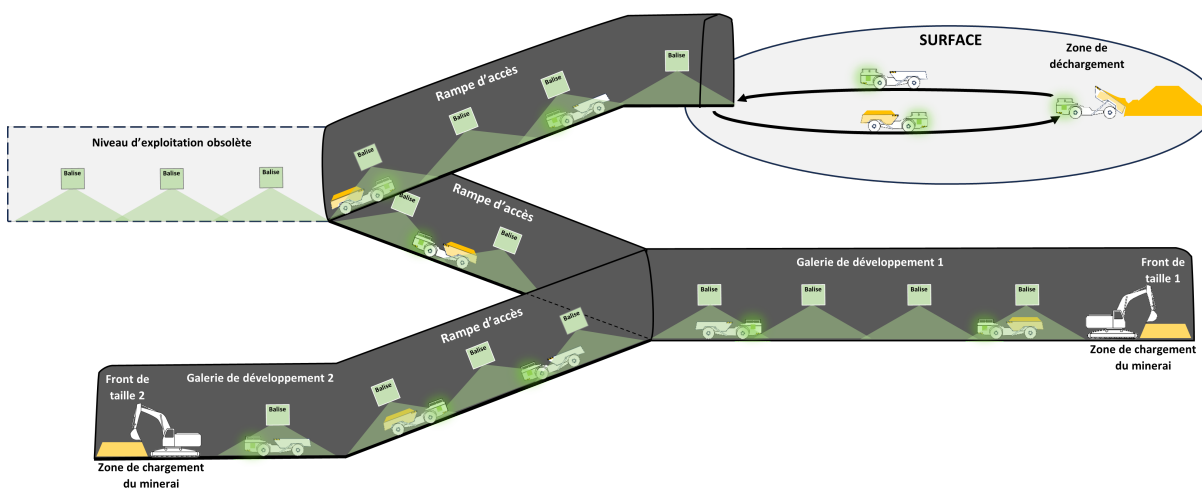


FIGURE 4.2 Aperçu des activités de transport de minerai des HT dans la mine 1

Dans cette mine, le cycle classique de transport de minerai d'un HT en six étapes se décompose comme suit :

1. Trajet aller : depuis la surface, le HT entre dans la mine via une rampe d'accès puis s'enfonce en cheminant sur celle-ci jusqu'au niveau de profondeur qu'il cible ; il se rend ensuite au front de taille auquel il est affecté via les galeries de développement de ce

niveau ;

2. Attente et positionnement (chargement) : si le site est déjà occupé, le HT doit patienter, après quoi il se positionne adéquatement dans la zone de chargement du minerai.
3. Chargement : le HT est chargé en quelques pelletées par l’engin de chargement, que ce soit une LHD ou une excavatrice de petite taille ;
4. Trajet retour : le HT suit généralement un itinéraire inverse à celui emprunté à l’aller pour rejoindre la surface et le site de déchargement qui s’y trouve ;
5. Attente et positionnement (déchargement) : le HT attend éventuellement que le site de déchargement se libère puis il se positionne ; et
6. Déchargement : le HT déverse sa cargaison de minerai au site de déchargement avant de réitérer le cycle complet.

Naturellement, ces cycles de transport de minerai ne sont pas immuables. Ils peuvent être régulièrement interrompus au cours d’un même quart de travail, comme lors du transport de roches stériles par le HT ou lors de son ravitaillement en carburant.

4.2 Contraintes et spécificités du cas d’étude

Au sein de cette section, on cherche à caractériser plus particulièrement l’environnement industriel auquel on s’intéresse et à identifier les contraintes qui s’y rattachent pour mieux adapter par la suite notre processus de préparation des données. Au regard de notre sujet d’étude, on analyse en particulier le fonctionnement du site minier relativement aux activités des HT, ainsi que les contraintes et spécificités liées à l’étude des trajets historiques de ces derniers.

4.2.1 Activités des HT

Les contraintes et les spécificités relatives aux activités des HT proviennent essentiellement du fait que l’on s’intéresse ici à une mine souterraine, autrement plus complexe et versatile que son homologue à ciel ouvert. Le réseau minier souterrain entraîne en effet les problématiques suivantes :

- **Conflits de trajectoires** : les rampes et les galeries étroites, qui ne disposent généralement que d’une seule voie, peuvent être simultanément partagées par des HT, des véhicules utilitaires et des piétons ; et ce, dans les deux sens de circulation. En plus d’un certain ralentissement des HT à proximité des piétons, les stratégies d’évitement laborieuses explicitées dans notre revue de la littérature peuvent être nécessaires lors

des croisements de véhicules. Dans la mine 1, ce sont effectivement les HT chargés qui sont prioritaires, i.e. les HT qui se dirigent vers la surface en remontant via une rampe. En outre, les intersections de galeries sont à même de créer des situations additionnelles de conflits de trajectoires ;

- **Manque de visibilité** : à l’obscurité ambiante des galeries souterraines s’ajoutent naturellement les nombreux angles morts des HT dûs aux parois rocheuses environnantes mais aussi à leur propre carrosserie. En effet, leur profil bas et leur cabine enchâssée restreignent notablement le champ de vision du conducteur. Ce manque de visibilité rend inévitablement la conduite plus difficile, et tend à faire ralentir d’autant plus les conducteurs des HT pour éviter les accidents avec les autres usagers du réseau minier ; et
- **Variété des sections de trajets** : comme mentionné par Li et al. [15] et comme évoqué dans la sous-section 2.3.1, les mines souterraines dont l’accès se fait via une rampe ont la particularité supplémentaire d’exposer leurs HT à des environnements de fonctionnement plus variés au cours de chacun de leurs trajets de transport de minerai. En effet, ces HT doivent à la fois se déplacer dans deux milieux souterrains bien distincts (i.e. à la fois dans des galeries souterraines horizontales rectilignes à multiples intersections et à la fois dans des rampes fortement pentues dont la courbure est parfois très marquée et qui peuvent être utilisées simultanément par tous les autres HT actifs dans les deux sens de circulation), mais aussi en surface. Dans ce dernier milieu, les conflits de trajectoires sont nettement moins problématiques mais les conditions météorologiques affectent tout à coup les trajets des HT, et potentiellement leurs TT [26]. En effet, la présence ou non de précipitations pluvieuses ou neigeuses est à même de changer la nature du terrain, et la température, l’humidité et la pression en surface varient à diverses échelles temporelles comparativement à la température des galeries souterraines, ajoutant de la variabilité aux trajets des HT. De plus, la visibilité est habituellement fortement accrue du fait de l’absence d’angles morts dûs aux parois rocheuses et grâce à la lumière du jour, mais elle varie amplement entre les quarts de travail diurnes et nocturnes et peut être particulièrement affectée par les conditions de brouillard voire de blizzard. Toutefois, comme détaillé ultérieurement dans les sous-sections 4.2.2 et 4.4.2, les portions de trajet en surface devront être exclues de notre étude. Les itinéraires auxquels nous nous intéresserons traverseront donc uniquement les rampes et les galeries.

Par ailleurs, des problématiques liées aux activités des HT et non inhérentes aux mines souterraines à accès par rampe concernent la mine 1, notamment de potentielles pauses

et détours des HT durant les périodes de transport de minerai (ravitaillement en carburant, pauses-repas, visite des installations sanitaires, objet à passer chercher au garage ou au refuge souterrain, etc.), ainsi que des problématiques de prévisibilité liées à l’existence de deux rampes et donc d’une plus grande variété de trajets et de conditions de trajets. Concernant les pauses et détours que l’on vient de mentionner, il ne faut pas négliger un point en particulier : la rémunération des conducteurs de HT de la mine 1 se faisant au prorata du nombre de tonnes de minerai extraites, leurs intérêts sont alignés avec une productivité élevée ce qui devrait limiter théoriquement les pauses et détours superflus.

4.2.2 Étude des trajets historiques des HT

En souterrain, l’absence d’accès au GPS et aux systèmes de triangularisation qui peuvent être utilisés dans les mines de surface complexifie significativement l’étude des déplacements des HT. Toutefois, la mine 1 dispose d’un réseau récent de connectivité souterraine auquel sont connectées un total de 135 balises souterraines de localisation de véhicules, lesquelles permettent d’étudier les déplacements historiques des HT y opérant. Cette technologie de balises est très similaire à celle installée dans la mine étudiée par Baek et Choi [25]. Les balises, représentées à la figure 4.2, peuvent reconnaître dans leur périmètre de détection respectif les puces d’identification par radiofréquence (*Radio Frequency IDentification* (RFID)) fixées sur chacun des HT de la mine 1, permettant leur association à un identifiant unique. La nature de ces données de détection sera détaillée à la sous-section 4.4.2. Les balises ne sont pas équitablement réparties sur l’intégralité du site minier : elles sont plutôt concentrées autour de points d’intérêt comme des intersections importantes ou le garage souterrain, tandis qu’elles sont bien plus rares dans les extrémités des ramifications de galeries et qu’elles sont presque complètement absentes de la surface de la mine. Ce dernier point est important car, dès lors qu’ils sont en surface, les HT évoluent dans un espace ouvert (i.e. les déplacements des HT sont peu contraints dans toutes les directions du plan). Il n’est donc plus possible de retracer leurs trajets historiques via les balises si elles ne couvrent pas la quasi-totalité du périmètre de déplacement des HT, puisque l’on ignore sinon l’existence d’éventuels détours ou pauses. Aussi, au vu du très faible nombre de balises capables de détecter les HT se déplaçant à la surface, l’étude des trajets historiques des HT de la mine 1 via les données de détection des balises doit se limiter aux trajets souterrains tandis que l’étude des trajets en surface doit passer par d’autres technologies.

4.3 Spécifications des requis du modèle

Afin de garantir le fonctionnement normal de notre modèle de préparation de données et de le rendre applicable au plus grand nombre possible de mines souterraines similaires, on spécifie les deux requis suivants :

1. **Réalisme des TT** : notre modèle devra permettre d'identifier et d'éliminer les éventuels trajets considérés comme étant impossibles physiquement, i.e. les trajets pour lesquels le TT observé implique nécessairement que le HT ait dû se déplacer à une vitesse excédant ses capacités réelles. Ces dernières lui permettent théoriquement d'approcher les 50 kilomètres par heure sur terrain plat, et l'environnement minier souterrain devrait considérablement limiter les pointes de vitesse en descente. Nous devons donc être attentifs à la cohérence des TT les plus courts détectés par le modèle ; et
2. **Capacité de généralisation** : notre modèle doit pouvoir permettre de nettoyer les données d'autres mines souterraines similaires sans qu'il soit nécessaire de modifier sa structure. Naturellement, la diversité des problématiques pouvant être rencontrées nécessite de prévoir au sein de notre modèle toutes les adaptations pouvant être nécessaires pour nettoyer convenablement les données malgré tout.

4.4 Identification des variables d'intérêt et des données correspondantes

Dans la présente section, on tâche d'identifier les variables d'intérêt, i.e. les variables à même d'influencer notablement les trajets et incidemment les TT des HT de la mine 1, ou qui contiennent des informations importantes, puis on recherche les données correspondantes dans la BDD de la mine 1.

4.4.1 Variables d'intérêt

On dresse ci-après une liste des variables relatives aux trajets des HT que nous considérons comme pertinentes et que nous utiliserons, en expliquant notre sélection.

- **TT** : c'est notre variable de sortie, et donc la variable cible de l'entraînement de notre modèle prédictif. Elle n'est pas directement disponible, nous allons donc devoir l'inférer bien que sa fiabilité soit critique pour chacun des trajets retenus dans notre étude. Des inexactitudes significatives et/ou nombreuses compromettraient sérieusement les performances de prédiction de notre modèle ;
- **Point de départ et point d'arrivée de l'itinéraire** : ils définissent l'itinéraire suivi et peuvent permettre d'entraîner les modèles de prédiction sur des itinéraires spéci-

fiques. Notons que, combinés au plan de la mine, ils permettent directement d'évaluer la distance de trajet au besoin. Cette dernière n'est pas nécessaire au bon fonctionnement de notre méthodologie ;

- **Sens vertical du trajet** (montée/descente) : cette variable est utile pour tous les trajets étudiés. En effet, dans le présent mémoire, on choisit de s'intéresser à des itinéraires de transport de minerai s'étalant sur une distance importante, durant lesquels les HT empruntent alors nécessairement la rampe (et se retrouvent donc prioritaires ou non selon qu'ils montent ou qu'ils descendent). La raison de ce choix réside d'une part dans le fait que les prédictions de TT sur des itinéraires de transport de minerai complets ou quasi-complets se prêtent bien mieux à être réutilisées lors de la planification des opérations minières ; et d'autre part dans le fait que l'extraction de TT réels à partir des données historiques de la mine 1 gagne en fiabilité et en précision relative lorsque l'on étudie des trajets plus étendus. On admettra que les trajets de HT en descente se font à vide tandis que leurs trajets en montée se font en charge ;
- **Type de HT** : étant donné qu'il peut exister plusieurs types de HT différents dans une mine, le type de HT effectuant un trajet donné est une variable d'intérêt qui peut fortement influencer sur le TT. Dans la mine 1, tous les HT sont de même modèle et ils sont tous pilotés par des conducteurs humains, cette variable n'est donc pas utilisée pour cette mine en particulier. Ce n'est pas le cas dans la mine 1 mais nous reviendrons sur ce point par la suite ;
- **Identifiant du HT** : bien que les HT parcourant la mine 1 soient tous issus du même modèle, la connaissance de l'identifiant du HT qui réalise un trajet donné peut permettre de prendre en compte les caractéristiques individuelles de ce HT. Ce dernier peut par exemple présenter des défauts mécaniques ou des détériorations, et certaines de ses pièces peuvent avoir été remplacées ou modifiées au fil de ses maintenances, affectant conséquemment ses performances. Pour un quart de travail donné, si chaque conducteur est invariablement attribué au même HT, l'identifiant de ce dernier est une variable encore plus pertinente ;
- **Détours ordinaires** : lorsque l'on s'intéresse à un itinéraire en particulier dans la mine, on aimerait savoir distinguer les cas où le HT a suivi très précisément la séquence de segments de route correspondante, des cas où il a effectué un détour qui relève de conditions opérationnelles ordinaires (les autres types de détours ne relèvent simplement pas du même itinéraire). Parmi ces détours, on peut prendre l'exemple d'un ravitaillement en carburant ou d'un passage du conducteur aux sanitaires. On associe ces événements à des actions facultatives internes à un trajet et on devrait donc considérer que le HT a respecté l'itinéraire étudié dans ces conditions. Une variable indiquant l'occurrence ou

non d'un détour ordinaire durant un trajet sur un itinéraire donné est donc importante pour notre étude et nous tâcherons de l'inférer ;

- **Pauses de longue durée** (booléen) : si le HT effectue une pause particulièrement longue au cours de son trajet, la validité du TT mesuré à la fin de son trajet peut être remise en question. De ce fait, l'existence ou non d'une pause anormale au cours d'un trajet est une variable d'intérêt. Ces pauses pourraient être déduites d'autres données déjà existantes dans la BDD de la mine 1, nous tâcherons donc de les identifier lorsque cela est possible ;
- **Nombre de HT actifs** : il s'agirait ici plus précisément d'une variable donnant, pour un quart de travail donné et pour un itinéraire donné, le nombre de HT se déplaçant notablement (i.e. ils ne sont ni en maintenance, ni occasionnellement actifs, ni limités à une zone de travail réduite, ni actifs essentiellement en surface) sur cet itinéraire. Cette variable permet de globalement rendre compte de l'état du trafic routier dû aux passages de HT dans la rampe et dans les galeries partagées. En effet, un plus grand nombre de HT actifs en souterrain favorise les conflits de trajectoires et est donc à même de causer un ralentissement global du trafic, plus particulièrement pour les trajets en descente car ceux-ci sont non prioritaires. Notons que cette variable serait connue avec une grande précision par les planificateurs en amont de chaque quart de travail ;
- **Quart de travail** : les conditions de fonctionnement peuvent varier selon le quart de travail (resp. de jour ou de nuit) et ce même en souterrain. On peut penser en particulier à la variation des équipes de travail et du nombre de superviseurs. On aimerait donc connaître le quart de travail durant lequel se déroule chaque trajet ;
- **Jour de la semaine** : on suspecte les conditions de fonctionnement de varier aussi selon le jour de la semaine, on prend donc cette variable en compte ;
- **Saison** : on considère que cette variable peut permettre à nos modèles de prédiction d'identifier des motifs saisonniers latents dans les conditions opérationnelles. Les HT transitent en effet par la surface et on peut supposer en particulier que les conditions hivernales peuvent rendre globalement plus difficiles les trajets, même une fois que les HT pénètrent à nouveau dans la mine via la rampe d'accès. On cherchera donc à prendre en compte l'alternance des saisons ;
- **Position temporelle du trajet relativement au quart de travail** : on cherche à évaluer le moment auquel le trajet a eu lieu durant le quart, car on considère que le début et la fin des quarts provoquent nécessairement des variations de conditions opérationnelles ;
- **Position temporelle du trajet relativement à la plage temporelle des données**

disponible : cette variable vise à donner la possibilité à nos modèles de ML d’inférer une réduction progressive globale des TT sur un itinéraire donné au fur et à mesure que l’itinéraire gagne en ancienneté. En effet, on peut imaginer un gain de fluidité au fil des années, liés à divers processus d’amélioration continue ; et

- **Écart temporel entre les trajets** : cette variable n’a d’intérêt que pour les RNN. Nous jugeons qu’il ne serait pas pertinent d’explicitier sa préparation avant le chapitre 5, car elle nécessite d’être contextualisée avec l’utilisation qu’en fera notre RNN. Aussi, nous ne la réévoquerons qu’au prochain chapitre.

On peut remarquer que cette liste de variables d’intérêt n’inclut pas l’identifiant du conducteur de chaque HT pour chaque trajet, car cette donnée ne nous a pas été rendue disponible pour des raisons de confidentialité. Lorsque les conducteurs n’ont pas chacun leur HT attribué, cette variable est pourtant intéressante puisqu’elle permettrait d’incorporer dans l’étude des notions sous-jacentes telles que le niveau de compétence des conducteurs et leur style de conduite. L’identifiant du conducteur n’est pas retrouvé indirectement dans le contexte de la mine 1, puisque les HT de la mine 1 changent fréquemment de conducteur sans nous permettre d’établir une relation fiable entre l’identifiant du HT et l’identifiant du conducteur. Par ailleurs, ces changements réguliers de conducteur peuvent induire une plus grande volatilité imprévisible des TT de chacun des HT comparativement à une mine dans laquelle un unique conducteur serait attribué à chaque HT.

Il manquerait aussi une variable quantifiant le kilométrage de chaque HT. Cette variable est intéressante et complémentaire avec l’identifiant du HT puisqu’elle permet de prendre en compte le fait que les performances des HT devraient globalement décliner avec leur ancienneté. L’identifiant du HT ne permet pas de rendre compte à lui seul de cette information, qui dépend de la période temporelle étudiée (i.e. seuls les RNN peuvent inférer cette relation). Notons que des HT déjà anciens côtoient des HT neufs progressivement introduits dans la mine au fil des années, induisant potentiellement des différences de performances plus importantes entre HT de même modèle, à une période temporelle donnée.

Le chargement (ou tonnage) des HT ne sera pas non plus utilisé. Il n’est pas censé avoir d’impact sur les TT en descente mais peut assurément faire varier la durée des trajets en montée. Il serait donc pertinent de l’intégrer dans de futurs travaux.

La détection des ravitaillements en carburant pourrait aussi être pertinente. Elle pourrait être déduite des variables de télémétrie indiquant le niveau de carburant des HT.

Notons enfin qu’il faudrait nécessairement ajouter à cette liste un éventail de conditions météorologiques si l’on cherchait à s’intéresser directement aux sections de trajet en surface.

4.4.2 Données correspondantes

Les variables d'intérêt qui permettent d'expliquer les TT de HT ayant maintenant été identifiées, il s'agirait de retrouver l'expression de chacune de ces variables via des données réellement disponibles. Naturellement, de très grandes quantités de données sont contenues dans la BDD de cette mine, et il est important de choisir précautionneusement dans celle-ci les attributs qui semblent les plus pertinents.

Notre exploration de cette BDD révèle en fait que cette tâche est notablement moins triviale que d'apparence puisque la majorité des variables énoncées dans la sous-section précédente ne correspondent directement à aucun attribut de la BDD. Elles doivent donc être retrouvées indirectement en utilisant les données d'attributs existants pour générer les données permettant d'exprimer les variables désirées. La complexité de ce processus d'« ingénierie des caractéristiques » (ou « *feature engineering* ») dépend à la fois des variables recherchées, de la fiabilité avec laquelle on veut exprimer celles-ci, des attributs disponibles dans la BDD et de la qualité des données historiques qui y sont associées. Ainsi, dans les cas où l'on cherche à exprimer une variable d'intérêt avec une fiabilité élevée, la complexité de ce processus peut s'accroître de façon spectaculaire si la relation entre la variable et les attributs est trop indirecte et/ou que les données exprimant ces attributs manquent de fiabilité ou sont trop lacunaires. L'apparition de problématiques liées aux données historiques entraîne en effet l'implication additionnelle de nouveaux attributs explicatifs souvent moins directement liés à la variable recherchée et eux-mêmes potentiellement exprimés par des données défectueuses, et ainsi de suite.

Par ailleurs, notre analyse de la BDD de la mine 1 nous mène au constat supplémentaire suivant : les variables d'intérêt citées plus haut pourront et devront **toutes** être inférées au moyen de données issues des balises de détection de véhicules ou bien de données externes à la BDD. La fiabilité de ces balises est donc critique pour limiter l'apparition des problématiques d'ingénierie des caractéristiques détaillées au paragraphe précédent.

Étant donnée leur criticité, on détaille dans le présent paragraphe la nature des données produites par les balises de la mine 1. Ces dernières ne produisent en fait pas une donnée unique à chacune de leurs détections, mais un lot de trois données : lorsqu'un HT donné passe dans le périmètre de détection d'une balise, celle-ci détecte sa puce RFID et relève l'identifiant de HT associé ; elle signe simultanément de son propre identifiant ce relevé, et l'horodate avant que ce lot de données de détection ne soit stocké dans la BDD de la mine 1. Pour un HT donné se déplaçant activement dans la mine, l'ensemble des données de détection de toutes les balises permet alors théoriquement de retrouver une liste horodatée de certaines positions occupées successivement par ce HT sur n'importe quelle période, rendant

ensuite possible la reconnaissance des trajets de ce HT et l'association de multiples variables contextuelles à ce trajet.

L'exploration de la BDD finit aussi de justifier que nous ne nous intéressions qu'aux trajets se déroulant en souterrain : bien que le GPS permette habituellement de localiser facilement les HT dès lors qu'ils sont en extérieur, la BDD de la mine 1 n'inclut pas de données permettant de retracer les trajets historiques des HT se déroulant à la surface de ce site minier. Étant donné que nous avons aussi mis en évidence à la sous-section 4.2.2 la quasi-absence de balises à la surface de la mine 1, nous ne disposons pas de moyens de retracer rigoureusement les trajets de HT en surface et nous les excluons donc de notre étude.

On présente successivement dans le tableau 4.1 les attributs finalement sélectionnés dans la BDD de la mine 1, tous issus de détections de balises. On précise aussi le format de chacune de ces variables, et l'inventaire impressionnant de variables d'intérêt que nous allons tâcher d'en extraire. Au vu de la complémentarité des attributs inclus dans les lots de données de détection que génèrent les balises, il est pertinent de rassembler les variables d'intérêt pouvant en être issues dans une unique case du tableau.

TABLEAU 4.1 Caractérisation des attributs sélectionnés dans la BDD de la mine 1

Source	Attribut	Format	Variables d'intérêt à extraire
Détections des balises	Horodatage des détections	<i>YYYY-MM-DD hh:mm:ss</i>	TT Points de départ et d'arrivée Sens vertical du trajet Type de HT Identifiant du HT
	Identifiant des balises rencontrées	Suites de caractères alphanumériques de six à 42 caractères	Détours ordinaires Pauses de longue durée Nombre de HT actifs Quart de travail Jour de la semaine Saison
	Identifiant du HT détecté	<i>XXXXXnn</i> avec <i>nn</i> allant de 01 à 11	Moment dans le quart de travail Moment dans le jeu de données Écart temporel entre les trajets

Concernant l'identifiant des balises, décrit comme une « suite de caractères alphanumériques allant de six à 42 caractères » dans le tableau 4.1, il s'agit en fait plus spécifiquement d'une succession très variable de mots-clés de localisation (rampe, type de galerie, portail de surface, points cardinaux, etc.), ainsi que d'un éventuel nombre donnant la profondeur du niveau

auquel est installée la balise et d'un éventuel second identifiant de la balise entre parenthèses, bien plus concis.

Dans le contexte de notre étude, nous disposons d'un plan des différents niveaux de la mine 1 sur lequel figure l'emplacement des balises de localisation qui y sont installées. Cette source de données est absolument nécessaire à la fine compréhension des séquences de détection d'un même HT par les balises. En effet, les données brutes de détection à elles seules peuvent être notablement difficiles à interpréter et ne permettent donc pas toujours de retracer et d'analyser les trajets. On pourrait aussi approximer via ce plan une variable secondaire, la distance entre balises. Toutefois, nous n'en aurons pas besoin.

4.5 Préparation des données

Une fois tous les attributs pertinents repérés, les données opérationnelles qui y sont associées doivent être extraites, nettoyées et préparées au cours des nombreuses étapes décrites dans l'ordre des sous-sections suivantes. L'extraction des trois attributs retenus de notre BDD se feront via des requêtes écrites en langage SQL. Le langage informatique utilisé dans toutes les autres tâches de programmation de ce mémoire sera le langage Python.

4.5.1 Extraction du jeu de données d'intérêt

On commence par se fixer l'objectif d'extraire de la BDD l'ensemble des données opérationnelles liées aux attributs pertinents.

Dans un premier temps, on cherche la plus ancienne occurrence d'un lot de données de détection de balise dans la BDD pour borner la période historique maximale de données que nous allons étudier. Elle semble manifestement correspondre au moment où les détections d'une majorité des balises qui sont actuellement installées ont commencé à être enregistrées dans la BDD de la mine 1. Ainsi, un réseau de balises maillait déjà le réseau minier au moment de ce premier enregistrement et, en supposant que les enregistrements ont été réguliers jusqu'à ce jour, nous devrions disposer de plus de trois ans de données opérationnelles à extraire.

Cette période temporelle identifiée, on analyse plus précisément la structure de la vaste BDD de la mine 1 pour comprendre entièrement les relations qui lient les attributs pertinents les uns aux autres à travers les différentes tables de données.

Une fois cette étape d'analyse consciencieusement réalisée, on formule en langage SQL la requête permettant d'extraire les données qui nous intéressent, en mentionnant les dénominations précises des attributs pertinents tels qu'ils sont enregistrés dans la BDD avec leurs

relations, et en limitant les véhicules étudiés aux seuls engins dont le chemin primaire est « `\Mine1\Trucks\` », car on a pu l’identifier comme renvoyant exclusivement les 11 HT de la mine 1. Les données historiques de tous les attributs mentionnés dans le tableau 4.1 sont alors effectivement extraits dans le jeu de données qui sera le théâtre des manipulations à suivre. Nous disposons alors de quelques trois millions de détections de HT par les balises.

4.5.2 Identification du niveau de profondeur de chaque balise

Cette sous-section vise à automatiser l’association de chaque balise à la profondeur du niveau auquel elle est installée, et ce, dans le but de faciliter certaines manipulations ultérieures. Aucun attribut de la BDD de la mine 1 ne correspond en effet à la variable de profondeur. Par ailleurs, bien qu’il soit théoriquement possible d’utiliser la lecture du plan de la mine 1 pour associer chaque identifiant de balise à une profondeur, le processus ne serait pas automatisé et il faudrait donc intervenir régulièrement pour prendre en compte de nouvelles balises. Il n’est pas non plus souhaitable de se risquer à effectuer plus de cent associations manuellement grâce au plan. Ce dernier n’est d’ailleurs potentiellement déjà plus à jour. Ajoutons aussi que les balises y ont clairement été positionnées à la main et que l’on retrouve quelques rares erreurs de frappe handicapantes dans les identifiants qui y sont inscrits. Toutes ces réflexions nous mènent à définitivement exclure cette alternative, pour plutôt nous intéresser aux informations contenues dans les identifiants des balises.

Tout d’abord, on programme un algorithme pour qu’il accède au jeu de données préalablement extrait, qu’il trouve la colonne correspondant à l’attribut « identifiant de balise », qu’il lise chacune des chaînes de caractères qui s’y trouvent (i.e. les suites de caractères correspondant auxdits identifiants de balises), qu’il repère toutes les chaînes distinctes et, enfin, qu’il les renvoie à l’utilisateur sous forme de liste. Grâce à cet algorithme, on dispose donc de la liste de tous les identifiants distincts de balises ayant été enregistrés dans la BDD de la mine 1 sur l’intégralité de notre période d’étude.

Dans cette liste, on observe l’ensemble de ces identifiants pour déterminer précisément les expressions régulières qui indiquent une profondeur **et** qui ne devraient jamais être confondues avec d’autres informations contenues dans les identifiants des balises. Ces derniers sont parfois étendus et riches en successions diverses de caractères variés, à l’image des exemples **fictifs** suivants, librement inspirés de balises existantes : « `RAMP1 575-550 (AN136)@Ramp 550-575 (AN128)` » ; « `525 FW1 West DP136 (AN49)` » ; ou encore « `DOVE PORTAL 3 (AN37)` ». Fait anecdotique, comme dans le tout premier identifiant évoqué, on trouve effectivement des double-espaces dans certains identifiants, ce qui manque cruellement de praticité et d’intuitivité lorsque l’on recherche un identifiant en particulier, de mémoire.

Pour déterminer les règles à imposer à la fin de cette manœuvre, voici le processus détaillé que nous avons suivi :

1. On observe que l'on doit repérer dans l'identifiant une séquence de trois chiffres désignant la profondeur approximative de la balise, en mètres. Au minimum, la profondeur indiquée est de 125 mètres, il y a donc bien toujours trois chiffres ;
2. Il est nécessaire de rechercher les séquences de trois chiffres isolées ou non par des espaces, car un nombre considérable d'identifiants de balises ne comportent pas ces espaces ;
3. Malheureusement, on remarque facilement des séquences de trois chiffres sans lien avec la profondeur des balises. On ajoute une condition supplémentaire : le nombre formé par les trois chiffres doit être un multiple de 25. C'est une caractéristique commune des profondeurs de niveaux due à l'écart inter-niveaux de 25 mètres ;
4. Toutefois, on remarque que certains identifiants de balises contiennent deux indicateurs de niveau de profondeur distincts. C'est en particulier le cas pour des sections de rampe, qui relient deux niveaux. On décidera de moyenner les deux valeurs pour obtenir une profondeur approximative de la balise ; et
5. Finalement, on voit que des balises n'incluent pas leur indicateur de profondeur sous la forme d'un nombre mais sous la forme d'une expression textuelle. On identifie en particulier les chaînes de caractères suivantes : « surface », « SURFACE », « portal », « Portal », « PORTAL » et « OUTSIDE ». Ces termes désignent la surface ou bien des portails qui la précèdent immédiatement. On associe les identifiants correspondants à une profondeur nulle.

On a ainsi dressé les principales règles à suivre pour réaliser les associations entre balises et profondeur. On peut maintenant s'atteler à la programmation d'un algorithme capable de reconnaître lui-même automatiquement les indicateurs de profondeur contenus dans les identifiants des balises, même pour les nouvelles balises régulièrement ajoutées au réseau minier. On crée d'abord un tableau associatif Python, appelé « dictionnaire » en programmation, qui visera à stocker la correspondance entre chaque identifiant de balise et le niveau de profondeur pouvant en être inféré. On crée ensuite notre algorithme de fouille de texte (plus précisément un algorithme d'extraction d'information textuelle) qui prend en entrée la liste des chaînes de caractères correspondant aux identifiants des balises. Il lira chacune de ces chaînes en ciblant les indicateurs de profondeur qui pourraient s'y trouver d'après les règles précédentes. Simultanément, il remplira le dictionnaire Python avec chaque paire d'éléments ayant été déduite. On demande aussi à l'algorithme de signaler les éventuels identifiants de balises qu'il serait incapable de gérer.

Après un premier fonctionnement de l'algorithme, quelques incohérences sont trouvées dans le dictionnaire qu'il renvoie, avec toutes la même source : les suites de chiffres incluses dans le sous-identifiant des balises (lui-même inclus dans l'identifiant) peuvent à la fois dépasser la centaine et être des multiples de 25. Ces sous-identifiants sont bien reconnaissables car mis entre parenthèses, on les exclut donc aisément des suites de caractères étudiées en modifiant le code de l'algorithme. Par ailleurs, concernant le cas où l'algorithme déclare ne trouver aucune indication de profondeur, on constate deux causes distinctes : le mot « DOME » n'a pas été associé à une profondeur nulle, ce à quoi on remédie immédiatement ; et quatre identifiants ne contiennent absolument pas d'indice permettant de se douter de la profondeur des balises correspondantes. Ces dernières sont situées sur des sections de rampe et le plan de la mine 1 nous renseigne sur le fait que ces balises sont proches de l'entrée de la mine. On estime la profondeur respective de ces trois balises d'après le plan puis on ajoute quelques lignes de code à l'algorithme pour qu'il insère lui-même ces valeurs dans le dictionnaire des balises lorsqu'on l'appellera.

Après une seconde itération de l'algorithme, on vérifie facilement la concordance de toutes les associations contenues dans le dictionnaire que retourne l'algorithme. Par ailleurs, toutes les balises ont bien été associées à une profondeur. Il est clair que si la difficulté de la tâche d'association des balises à leur profondeur respective avait été plus complexe encore ou que de nouvelles erreurs avaient été commises par l'algorithme, on aurait pu continuer de s'appuyer sur ce processus d'amélioration itérative de l'algorithme. Ainsi, après chaque tentative de ce dernier, on rechercherait dans les identifiants de balises des indicateurs de profondeur plus subtils non pris en compte par l'algorithme jusque-là. On améliorerait ensuite l'algorithme pour prendre en compte ces derniers, puis on réitérerait les opérations précédentes jusqu'à épuiser la liste de balises de profondeur inconnue.

Finalement, à l'issue des étapes présentées dans cette sous-section, nous disposons d'un algorithme qui renvoie un dictionnaire complet des couples « balise-profondeur » de la mine 1 en totale cohérence avec les profondeurs indiquées dans les identifiants et sur le plan. L'algorithme serait capable de s'adapter à de nouvelles balises dont les identifiants contiennent des indications de profondeur de format similaire, que ce soit pour la mine 1 ou, avec d'éventuelles modifications mineures du programme, pour d'autres mines. Il est clair que cet algorithme ne pourra en revanche pas gérer les nouvelles balises dont les identifiants ne comportent aucune indication de profondeur mais, en l'état actuel des choses, il n'est de toute façon pas possible de combler automatiquement cette lacune avec les informations contenues de la BDD de la mine 1. Les modifications qui pourraient être apportées pour régler ce type de problématique seront présentées dans la section 4.6.

Une problématique intrigante est toutefois soulevée : le dictionnaire complet retourné par l'algorithme contient une trentaine de couples distincts balise-profondeur de plus que le nombre de balises visibles sur le plan de la mine 1, pourtant récent. Nous nous pencherons sur cette problématique dans la sous-section suivante.

4.5.3 Correction d'incohérences entre les périodes de collecte de données

La durée de collecte des données étudiées s'étant étalée sur plusieurs années, diverses modifications sur les capteurs existants ou dans la collecte des données pourraient avoir amplement affecté les valeurs de certains attributs, les rendant visiblement incohérents d'un intervalle temporel à l'autre.

Pour en trouver un exemple probant, on commence par s'intéresser à la problématique soulevée dans la sous-section précédente. Pour cela, on commence par se fixer une profondeur quelconque correspondant à l'un des niveaux de la mine, puis on s'intéresse aux balises correspondantes d'après le dictionnaire qui avait été finalement obtenu en les comparant au plan de ce niveau censé recenser toutes les balises. On relève les identifiants des balises qui n'apparaissent pas sur le plan, mais qui apparaissent dans le dictionnaire. On remarque rapidement que chacun d'entre eux est notablement similaire à l'un des identifiants figurant sur le plan. On suspecte dès lors très fortement l'occurrence de nombreuses modifications des identifiants de balises dans les dernières années qui n'auraient pas été prises en compte dans la BDD de la mine 1, i.e. que des identifiants obsolètes de balises auraient été laissés tels quels dans les lots de données de détection anciennement émis par les balises et stockés dans la BDD. Des investigations plus poussées des historiques de détection des balises de la mine nous indiquent en effet que de nombreux identifiants cessent subitement d'apparaître dans la BDD à partir d'une certaine date (en existant bien sûr toujours dans la BDD puisque l'algorithme précédent les a rencontrés) au profit d'un nouvel identifiant qui apparaît quasi-simultanément et qui est très similaire à l'ancien. Par ailleurs, on ne trouve aucune indication dans la BDD permettant de rétablir les correspondances entre anciens et nouveaux identifiants. Une investigation similaire à celle décrite dans ce paragraphe est donc nécessaire. On ne trouve pas non plus de moyen externe de retracer l'historique des identifiants de chaque balise, ce qui s'avère problématique pour rattacher chaque ancien identifiant au nouvel identifiant correspondant avec certitude.

Concernant la tâche d'association des anciens identifiants avec leur nouvel identifiant respectif, on constate rapidement qu'il s'agit d'une tâche délicate puisque d'une part, une confusion entre balises pourrait avoir des conséquences critiques sur le fonctionnement de nos futurs algorithmes, et d'autre part, les modifications d'identifiants ayant pu survenir au cours des

dernières années sont considérablement diverses, ce qui limite l'utilité d'une automatisation du processus. Enfin, il n'y a qu'une trentaine d'associations à faire, ce qui représente un volume relativement réduit d'opérations.

Pour ces deux raisons, nous ne trouvons pas de meilleure solution que de lier les identifiants manuellement, niveau par niveau, en se fixant un protocole de vérification pour les cas les plus délicats. Ainsi, si l'on n'est pas absolument certain que deux identifiants donnés (que l'on présume désigner la même balise) désignent une même balise, on compare :

- D'une part l'horodatage de la première détection associée à l'identifiant qui figure actuellement sur le plan de la mine ; et
- D'autre part, l'horodatage de la dernière détection associée à l'identifiant qui est théoriquement le plus ancien.

Dans toutes les vérifications réalisées, ces deux horodatages étaient quasi-simultanés. Nous sommes ainsi assurés de la pertinence de l'association puisque le remplacement de l'identifiant est évident. Lorsqu'il n'y a aucun doute sur le fait que les deux identifiants désignent la même balise, on réalise immédiatement l'association.

Au final, après avoir réalisé toutes les associations, on retrouve le même nombre de balises que sur le plan, toutes associées à leurs identifiants les plus récents. On renomme donc dans notre jeu de données toutes les occurrences d'anciens identifiants par l'identifiant correspondant le plus récent. La connaissance préalable de la profondeur respective de chaque balise a facilité ce travail de correction puisque la recherche de fortes similarités textuelles entre les identifiants a pu se faire par groupes réduits, correspondant chacun à un niveau de la mine.

Après avoir été confrontés à la problématique précédente, on entreprend de vérifier la cohérence temporelle des autres variables dont nous disposons déjà via des représentations graphiques. Ainsi, on trace l'évolution des valeurs de notre seule variables quantitative pour le moment, la profondeur de chaque HT, en fonction du temps. Par ailleurs, un peu plus subtilement, on représente l'évolution du nombre d'occurrences des modalités des données qualitatives (i.e. l'identifiant de chaque HT et l'horodatage de chaque détection de chaque HT par les balises) en fonction du temps. Aucune incohérence ne transparaît.

4.5.4 Identification des itinéraires prédominants

Dans cette sous-section, on cherche à identifier les itinéraires de transport de minerai les plus empruntés historiquement, qui ne sont pas nécessairement les plus empruntés à ce jour. Ces itinéraires sont particulièrement intéressants pour tester la viabilité de notre modèle de prédiction, puisqu'ils correspondront à de très nombreux trajets réels, donc théoriquement à

un très grand nombre de détections par balises, ce qui pourrait permettre de disposer d’une quantité remarquable de données pour l’entraînement et la validation de notre modèle de prédiction sur ces itinéraires uniquement. L’obtention de bons résultats de notre modèle sur les itinéraires les plus empruntés est une condition nécessaire à son utilisation potentielle sur des itinéraires nettement moins fréquentés.

Avant d’identifier ces itinéraires prédominants, on commence par mieux les caractériser. En plus de devoir être remarquablement empruntés, ce qui reste notre critère principal, ils devront idéalement aussi couvrir une longueur importante. En effet, comme indiqué dans la sous-section 4.4.1, les itinéraires de transport de minerai quasi-complets sont mieux adaptés à être réutilisés dans le cadre de la planification des opérations minières, et l’étude de trajets plus étendus permet de retrouver des TT historiques plus fiables et précis. Ce critère est aussi utile pour que le nombre de balises rencontrées au cours du trajet soit suffisamment grand.

Afin d’identifier ces itinéraires, on trace l’évolution du niveau de profondeur des différents HT en fonction du temps d’après la liste des balises successivement croisées au cours de leurs trajets. On obtient alors des graphes permettant de mieux borner les niveaux entre lesquels sont détectés la plupart des trajets sur toute la période des données : on utilise les extrémités inférieure (i.e. la plus profonde) et supérieure (i.e. la moins profonde) les plus souvent observées sur le graphe. Il est important de noter que, même si l’on s’intéresse ici à une mine à accès par rampes dans laquelle les HT font essentiellement des allers-retours entre la surface et les fronts de taille, le niveau extrême supérieur que l’on identifie peut tout à fait se situer sous la surface si la dernière balise qui détecte habituellement le HT se situe dans un niveau souterrain. Ainsi, la prise en compte d’un niveau extrême supérieur a bien un sens pour identifier le niveau initial ou le niveau terminal des trajets les plus souvent observés.

Plusieurs observations peuvent être faites à partir des graphes obtenus. On représente l’un d’eux à la figure 4.3 et on fait apparaître par de fines lignes pointillées rouges les niveaux initiaux/terminaux profonds des trajets les plus populaires sur la période visée. Le cercle orange correspond quant à lui à des niveaux initiaux/terminaux de trajets peu classiques sur la période d’étude, pourtant fréquemment observés sur de courts laps de temps. Ils pourraient trahir des trajets non liés au transport de minerai mais nous ne pousserons pas cette dernière investigation. Tout comme le graphe de la figure 4.3, la plupart des graphes indiquent, quelle que soit la période, que les HT sont détectés entre la surface et un niveau dont la profondeur est au moins égale à 300 mètres dans la très grande majorité des cas. En particulier, de très nombreux trajets partent visiblement de la surface et semblent s’arrêter au niveau 300 puisqu’il y a visiblement nettement moins de détections plus en profondeur. Les niveaux inférieurs sont donc plus rarement visités, et le phénomène s’amplifie naturellement pour les

niveaux les plus profonds.

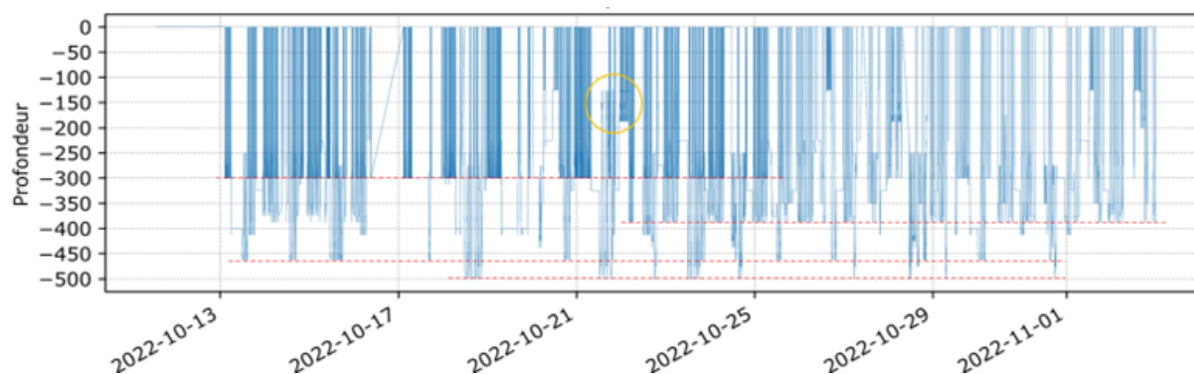


FIGURE 4.3 Graphe d'évolution de la profondeur d'un HT de la mine 1 sur 22 jours

Ces observations nous mènent à considérer que, dans la mine 1, l'itinéraire prédominant de grande longueur le plus pertinent relie la surface au niveau 300. Étant donné que tous les itinéraires complets passent par l'une des rampes, on commencera par accorder un intérêt tout particulier à la section de chacune des rampes qui relie la surface au niveau 300. On s'intéresse donc dans un premier temps à quatre itinéraires (en comptant les itinéraires allers et les itinéraires retours) délimités par les balises d'accès à la surface et par les balises d'accès au niveau 300. Si des dysfonctionnements particulièrement gênants de ces balises situées aux extrémités des itinéraires venaient à être remarqués par la suite, les critères exprimés plus tôt voudraient que l'on s'intéresse aux mêmes itinéraires que l'on raccourcirait autant que nécessaire. Par exemple pour la rampe 1, si la balise d'accès au niveau 300 présente un dysfonctionnement, on s'intéressera aux itinéraires (aller et retour) reliant la surface et la balise précédente, et ainsi de suite si d'autres dysfonctionnements étaient détectés pour les balises les plus profondes du trajet. On agirait similairement dans le cas où les balises les plus proches de la surface venaient à présenter un dysfonctionnement très considérable : on s'intéresserait à des itinéraires similaires que l'on aurait simplement raccourcis, en ne prenant plus en compte les balises problématiques situées aux extrémités de ces itinéraires. Rappelons bien qu'il s'agit uniquement ici de l'identification d'un itinéraire prédominant optimal ou quasi-optimal et que notre méthodologie pourra tout à fait s'appliquer à d'autres itinéraires moins empruntés.

4.5.5 Listage des balises de changement de niveau

Certaines étapes de notre modèle de préparation de données, présentées dans les sous-sections qui suivront, mais aussi de notre modèle de prédiction ne peuvent être rendues possibles

qu'en disposant de deux listes de balises bien particulières que l'on décrit dans le présent paragraphe. Ces deux listes seront respectivement dédiées aux trajets en descente et en montée et incluront chacune des balises que nous pourrions définir comme une « balise de changement de niveau », i.e. une balise située dans la rampe qui, dans un sens de trajet donné, détecte habituellement en dernière les HT avant qu'ils n'aient la possibilité d'accéder aux galeries de développement d'un nouveau niveau de la mine. On tâche d'illustrer cette notion par la figure 4.4, dans laquelle la balise entourée par un cercle bleu est une balise de changement de niveau en descente uniquement (si elle ne dysfonctionne pas), tandis que la balise entourée par un cercle orange est une balise de changement de niveau en montée uniquement (idem). Le fait de préciser « en dernière » dans la définition précédente a pour but de minimiser le risque d'attribuer à tort un changement de niveau à un HT (puisque'il est ainsi obligé de parcourir une majeure partie du dénivelé séparant deux niveaux avant d'être considéré comme changeant effectivement de niveau), et vise par ailleurs à mieux homogénéiser les propriétés des balises identifiées comme des balises de changement de niveau entre deux sites miniers différents. Ajoutons enfin que la balise de changement de niveau n'est pas nécessairement la dernière balise **installée** avant le niveau suivant, car cette dernière peut dysfonctionner régulièrement (et par conséquent ne pas détecter habituellement le HT en dernière avant son changement de niveau). Dans un tel cas, on s'intéressera successivement aux balises précédentes pour trouver la meilleure candidate.

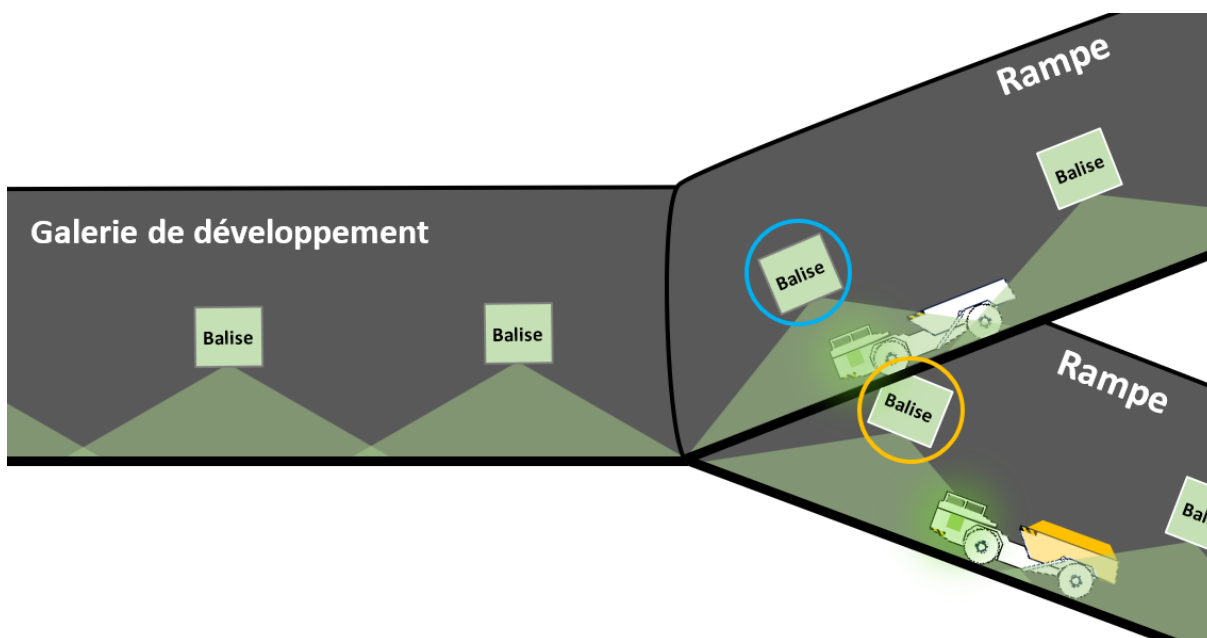


FIGURE 4.4 Emplacement des balises de changement de niveau

Passons à l'identification de ces balises. Pour la rampe 2, la tâche n'a pas vraiment de

sens, puisqu'elle ne dessert que le niveau 300. On passe donc directement à la rampe 1. La figure 4.4, évoquée dans le paragraphe précédent, ne permet hélas pas immédiatement de déduire comment automatiser la création des listes de balises désirées quelle que soit la mine étudiée. Le nom que l'on donne à ces balises (« de changement de niveau ») doit être interprété avec précaution puisqu'il peut mener à accorder du crédit à la stratégie naïve illustrée fictivement par le tableau 4.2. Cette dernière semble pourtant rarement donner des résultats satisfaisants. Elle consiste à suivre l'historique des données de détection des balises ayant identifié un certain HT se déplaçant uniquement le long de la rampe 1, tout en observant l'évolution de sa profondeur (donnée par notre dictionnaire des couples balise-profondeur). On répertorie alors chaque identifiant de balise qui correspondrait théoriquement à l'arrivée du HT au niveau suivant (i.e. à une nouvelle profondeur multiple de 25 et incluse entre 125 et 500) en considérant qu'il s'agirait là d'une balise de changement de niveau. On opère de même dans l'autre sens de circulation du HT. Dans la figure 4.2, on colore en rouge les nouvelles profondeurs rencontrées multiples de 25 et incluses entre 125 et 500 lors d'un trajet en descente, ainsi que les balises de changement de niveau correspondantes d'après cette stratégie. Toutefois, en mettant ici de côté d'éventuels dysfonctionnements des balises rencontrées dans la rampe, cette stratégie n'est malheureusement pas suffisamment fiable. En effet, ses résultats dépendent de l'architecture du réseau de balises et de la manière dont a été préalablement établie la valeur de la profondeur des balises situées dans la rampe d'accès. Deux cas posent en effet problème :

1. Le cas le plus simple est rencontré lorsque le niveau de profondeur d'une balise a mal été identifié auparavant. Cette identification incorrecte pourrait être causée par notre algorithme d'identification de profondeur, mais pas seulement. En effet, ce dernier accorde une confiance aveugle à la profondeur indiquée dans l'identifiant de chaque balise ; des erreurs pourraient donc exister parmi la ou les centaine(s) d'identifiants recensés pour un site minier donné. La stratégie précédente pourrait alors nous mener à détecter un changement de niveau qui n'en est pas un, ce qui serait une erreur critique pour le fonctionnement de notre algorithme ; et
2. Le second cas est plus complexe et, bien que peu dérangeant dans le contexte de la mine 1, il risquerait d'occasionner de nombreux faux changements de niveau difficiles à détecter dans certains autres sites miniers, selon la manière dont y ont été paramétrées les balises de détection. Considérons un itinéraire en descente $A \rightarrow B$ et un HT en train de descendre la rampe, situé entre ces balises A et B et en provenance de A. Ce HT se rend au niveau B', qui précède immédiatement le niveau de la balise B. Il n'est donc pas en train de faire un trajet qui nous intéresse, et ce dernier ne devrait pas être inclus dans notre étude. Pourtant, au moment où il s'apprête à pénétrer dans ce

niveau, l'une des balises situées dans la rampe quelques mètres en aval le détecte à son tour. S'il ne s'agit pas d'une balise de changement de niveau et que cette balise n'a pas été identifiée comme une balise « interdite » au cours de ce trajet, il ne se passe rien. Toutefois, sur le site minier fictif étudié, la profondeur qui a été associée à chaque balise correspond à la profondeur du niveau qui la suit immédiatement. Ladite balise est donc associée à la profondeur du niveau B. Le HT précédent vient donc de faire un trajet $A \rightarrow B$ sans jamais avoir descendu la dernière section de rampe. Notre algorithme considérerait automatiquement que son TT est valable. Par ailleurs, ce paramétrage des couples balise-profondeur est tout à fait crédible, et pourrait avoir été réellement mis en place volontairement. Il permettrait en effet de savoir dès qu'un HT commence à descendre la rampe vers le niveau suivant. En montée, il permettrait de savoir dès qu'un HT arrive à un nouveau niveau. Il permettrait aussi d'éviter de devoir mesurer la profondeur de chaque balise de la rampe si la profondeur associée à chaque balise doit être une valeur numérique (par exemple si la profondeur de balise est une véritable variable numérique de la BDD). Bref, dans une telle mine, un HT se rendant à n'importe quel niveau en descente risquerait d'être détecté au niveau suivant, à cause de notre choix malencontreux des balises de changement de niveau en descente. Inversement, en montée, il est possible que la balise de changement de niveau (détectée avec la stratégie naïve précédente) se trouve étonnamment plus haut dans la rampe que ne se situe l'entrée du niveau correspondant. En effet, lors d'un trajet $A \rightarrow B$ en montée dans cette mine, un HT ne rencontrera assurément pas de balise dont la profondeur est celle du niveau B avant d'atteindre l'entrée de ce même niveau. De plus, si aucune balise ne se trouve immédiatement à l'entrée de ce niveau, alors le HT doit continuer à remonter pour être détecté par la balise de changement de niveau correspondante. Il aurait donc pu rentrer dans le niveau B sans être détecté par la balise de changement de niveau correspondante, qui perd tout son intérêt puisqu'elle serait pourtant supposée nous permettre de prendre connaissance de l'arrivée du HT en B. Ces situations aberrantes, non causées par de véritables défauts du réseau de balises, nous amènent à chercher une autre stratégie.

La tâche peut donc théoriquement s'avérer difficile à automatiser, mais ce n'est pas le cas dans la mine 1. En effet, on a conscience d'une part que toutes les balises de changement de niveau recherchées se situent nécessairement soit sur la rampe soit à proximité immédiate (de manière à ce que chaque HT qui emprunterait la rampe se trouve à un instant donné de son trajet dans le périmètre de détection de cette balise). D'autre part, on remarque que chaque niveau de la mine 1 est invariablement doté d'une balise située sur l'accès via la rampe à la galerie de développement de ce niveau, et que les identifiants de ces balises possèdent tous un

TABLEAU 4.2 Stratégie naïve d'identification de balises de changement de niveau en descente

Identifiant de la balise	Profondeur
Balise 1	125
Balise 2	125
Balise 3	150
Balise 4	162,5
Balise 5	175
Balise 6	175
Balise 7	200
Balise 8	200

format identique. Ainsi, la liste des balises de changement de niveau de la mine 1 est censée être globalement la même en descente qu'en montée. En effet, en théorie, quel que soit le sens de circulation d'un HT se déplaçant dans la rampe et arrivant à un niveau donné, ce sera toujours la même de ces balises qui le détectera au plus proche de l'entrée dudit niveau. On mentionne tout de même un détail : pour assurer le bon fonctionnement des algorithmes ultérieurs, les balises de changement de niveau ne seront pas parfaitement identiques dans les deux sens de circulation. La balise du niveau 500 (resp. la balise d'accès à la surface) figure en effet uniquement dans la liste dédiée aux trajets en descente (resp. en montée). Sinon, on détecterait un changement de niveau lorsque le HT part du niveau 500 (resp. emprunte la rampe depuis la surface) pour rejoindre la rampe, ce qui ne correspond pas à la définition que nous nous en sommes fait puisque l'on veut que le HT ne change de niveau que lorsqu'il approche du niveau suivant. Pour le reste, les deux listes sont identiques. En prenant en compte ces considérations, nous créons facilement les listes de balises de changement de niveau de la mine 1.

Dans l'éventualité où ces balises n'existeraient pas, on aurait tout avantage à chercher une autre forme de récurrence visible des identifiants des balises de changement de niveau, en montée et en descente. Si cette approche ne donne pas les résultats escomptés, on testera la stratégie naïve mentionnée précédemment et on s'intéressa à la qualité des balises résultantes par rapport aux observations pouvant être faites sur le plan de la mine. Si ces balises sont suffisamment pertinentes, on pourra s'appuyer sur ces résultats préliminaires. Par la suite, on devra vérifier l'intégralité des listes obtenues en ajoutant manuellement toutes les balises de changement de niveau non repérées par la stratégie naïve et identifiables sur les plans de la mine. Il s'agit là d'un travail potentiellement laborieux et pouvant mener à diverses erreurs,

mais qui reste difficilement évitable si aucune mesure n'a été mise en place en amont pour le faciliter.

Bien que le choix des balises de changement de niveau soit critique dans notre méthodologie, il est possible de gérer un manque de fiabilité partiel de ces balises en rendant plus flexible les critères de filtrage du reste de notre méthodologie globale. Nous n'en aurons pas besoin pour la mine 1 mais tâcherons de préciser par la suite la marche à suivre dans une telle situation.

4.5.6 Création d'un algorithme de reconnaissance de trajets sans détours

Dans cette sous-section, on s'appuie dans un premier temps sur notre dictionnaire de couples balise-profondeur pour développer un algorithme de reconnaissance de trajets et de calcul des TT associés. Par la suite, on tâche d'améliorer les performances de cet algorithme en identifiant les détours ayant pu avoir lieu durant chaque trajet grâce à notre connaissance des listes de balises de changement de niveau.

4.5.6.1 Création d'un algorithme primitif de reconnaissance de trajets

Pour expliquer le fonctionnement de notre algorithme, on se donne un itinéraire orienté qui relie deux balises distinctes, que nous appellerons respectivement « balise A » et « balise B ». L'itinéraire qui nous intéresse sera appelé « itinéraire $A \rightarrow B$ », de même que les trajets parcourant cet itinéraire seront appelés « trajets $A \rightarrow B$ ». On pourra donc, par extension, évoquer des « TT $A \rightarrow B$ ».

Les données d'entrée de notre algorithme seront alors d'une part l'identifiant de ces deux balises (A et B) et d'autre part les données de détection produites par ces dernières, soit les données des trois attributs explicités dans la sous-section 4.4.2. La variable obtenue en sortie correspondra aux TT de tous les trajets ayant pu être reconnus entre ces deux balises au cours de la période étudiée. Ces TT reconnus, qui pourraient être détectés par milliers, seront représentés graphiquement sous la forme d'un unique histogramme pour rendre possible l'interprétation des résultats.

Le codage de l'algorithme, quant à lui, combine quelques boucles simples, ici en pseudo-code : Successivement pour chacun des 11 HT, on parcourt la liste des détections historiques à la recherche d'une détection issue de la balise initiale A uniquement. Si l'identifiant de la balise A est observé dans la liste :

- On relève l'horodatage correspondant à cette détection ; et
- Tant qu'on ne lit pas l'identifiant de la balise B dans la suite de la liste :

- On veille à borner la zone de la mine dans laquelle le HT est censé se déplacer pour limiter le nombre de déplacements qui seraient détectés comme des trajets $A \rightarrow B$ par l'algorithme mais que nous n'interpréterions pas comme tels, c'est-à-dire en s'assurant des critères suivants pour un trajet en descente (resp. en montée) :
 - **Orientation cohérente du trajet** : la profondeur de chacune des balises qui détecte le HT doit nécessairement être supérieure (resp. inférieure) ou égale à celle la profondeur de la balise initiale A ; et
 - **Non-dépassement** : la profondeur de chacune des balises qui détecte le HT doit être inférieure (resp. supérieure) ou égale à celle de la balise B.

Si l'identifiant de la balise A est observé :

- On retourne en arrière dans le code jusqu'à l'endroit où l'on a détecté la balise A pour la première fois, et on recommence les étapes précédentes puisque le HT est revenu à son point de départ sans passer par la balise B.

Si l'identifiant de la balise B est observé :

1. On relève l'horodatage correspondant à cet événement ;
2. On calcule l'écart temporel entre les deux horodatages correspondant respectivement à la détection par la balise A et par la balise B pour obtenir la valeur d'un TT $A \rightarrow B$; et
3. On enregistre cette valeur dans une nouvelle liste.

Enfin, l'algorithme trace un histogramme sur la base des durées rassemblées dans cette dernière liste.

Sur un itinéraire donné, les TT de tous les HT seront donc confondus sur le même histogramme. En effet, nous n'affinerons pas notre analyse préliminaire des TT au point de distinguer les TT obtenus par chacun des HT avant d'atteindre la sous-section 4.5.11. D'ici-là, nous ne cherchons qu'à disposer d'un portrait global de l'allure de l'histogramme des TT sur un itinéraire donné.

Tâchons maintenant d'améliorer la fiabilité des TT obtenus par cet algorithme. Ce dernier est pour l'instant rudimentaire et ne permet aucunement d'identifier les détours.

4.5.6.2 Amélioration de l'algorithme par détection de boucles inter-niveaux

On s'intéresse ici à la détection de « boucles inter-niveaux », i.e. des détours durant lesquels un HT effectuant un trajet en descente (resp. en montée) a remonté (resp. redescendu) une section de rampe sur un dénivelé qui équivaldrait minimalement à un écart inter-niveaux complet avant de se diriger de nouveau vers la balise B.

Avant d'expliquer comment détecter de tels phénomènes, on cherche dans un premier temps à justifier le fait qu'une très faible minorité de ces observations de boucles inter-niveaux devraient selon nous provenir d'un trajet $A \rightarrow B$ qu'on qualifierait d'« ordinaire », i.e. d'un trajet durant lequel le HT a vraiment cherché à rejoindre la balise B ou une balise située en aval en partant de la balise A ou d'une balise située en amont. En effet, il n'est pas totalement évident d'identifier lesquels des demi-tours virtuellement observables dans les données de détection historiques des balises font partie de trajets ordinaires : des conflits de trajectoire peuvent régulièrement forcer un HT à reculer pour trouver une chambre creusée dans la paroi afin de laisser passer un HT prioritaire arrivant en sens inverse. Il s'agit ici d'une manœuvre de croisement parfaitement ordinaire qui serait virtuellement assimilable à un demi-tour si le HT qui recule est détecté par une balise qu'il a déjà franchie auparavant. Pour autant, nous n'excluons pas de notre étude un trajet incluant une telle marche arrière puisqu'elle ne constituerait pas une boucle inter-niveaux d'après la définition que nous avons donné à ce terme dans le paragraphe précédent. En effet, pour qu'une marche arrière puisse être identifiée comme une boucle inter-niveaux, il faudrait déjà que le HT parcourt un dénivelé correspondant minimalement à un écart inter-niveaux, soit 25 mètres de dénivelé dans la mine 1 i.e. environ 200 mètres linéaires pour une pente à 12%. De plus, étant donné que cette marche arrière supposée ordinaire devrait faire partie d'une manœuvre de croisement, il faudrait aussi supposer que, durant ces 200 mètres de marche arrière, le HT ait continuellement cherché à se stationner dans l'une des chambres dédiées à ce type de situation sans en trouver une seule. Enfin, cela signifie qu'il n'a pas non plus eu l'occasion de se dégager de la rampe en empruntant l'entrée du niveau qu'il a nécessairement rencontré avant de parcourir minimalement un écart inter-niveaux. En bref, on considère cette accumulation de conditions comme étant suffisamment improbable pour justifier le fait de considérer chaque boucle inter-niveaux comme une preuve que le trajet observé n'est pas un trajet $A \rightarrow B$ ordinaire, et qu'il doit donc être exclu de l'étude des TT $A \rightarrow B$. On peut considérer que ces observations correspondent soit à d'autres activités ordinaires de HT se déroulant entre les bornes A et B, soit éventuellement à des cas de force majeure qui sortent donc d'un cadre d'étude portant sur les trajets ordinaires.

Passons maintenant à la stratégie de détection des boucles inter-niveaux que nous allons ajouter au sein de notre algorithme. Grâce à notre connaissance des listes de balises de changement de niveau, le principe de détection est simple. Une fois qu'une détection de la balise A est remarquée dans la liste des détections de balises et que l'on commence à lire la suite de cette liste pour trouver une détection de la balise B, alors, dès que le HT croise l'une des balises de changement de niveau, on garde en mémoire la profondeur de cette dernière. Pour un trajet en descente (resp. en montée), s'il y a ensuite détection d'une

autre balise de changement de niveau dont la profondeur est inférieure (resp. supérieure) à celle de la précédente balise de changement de niveau détectée, alors une boucle inter-niveaux vient d'être observée. On invalide donc ce trajet en reprenant la lecture de la suite de la liste à la recherche d'une balise A qui pourrait marquer le début d'un nouveau trajet $A \rightarrow B$ (tout début de l'algorithme). Si toutes les balises de changement de niveau ne sont pas complètement fiables pour un site minier donné, notre méthode reste alors identique mais des trajets non ordinaires supplémentaires sont à prévoir dans l'histogramme final des TT puisqu'on ne remarquera pas nécessairement toutes les boucles inter-niveaux.

En bref, on vient d'ajouter un critère additionnel aux exigences d'« orientation cohérente du trajet » et de « non-dépassement » qui apparaissaient déjà dans l'algorithme primitif : il s'agit d'un critère d'« absence de boucle inter-niveaux ».

4.5.6.3 Amélioration de l'algorithme par détection de boucles dans les sections de rampe simple

Le terme « section de rampe simple » pourra désigner ici à la fois la section de la rampe 1 qui relie la surface à la première balise de changement de niveau (au niveau 125), et la rampe 2 au complet. Ces longues sections de rampe ont la particularité de ne donner accès à aucun niveau d'exploitation. Ainsi, elles ne disposent naturellement d'aucune balise de changement de niveau, ce qui nous empêche pour le moment de détecter aisément les activités anormales qui s'y dérouleraient. En théorie, de rares cas de demi-tours pourraient pourtant y avoir lieu, et tâcher d'exclure les potentiels trajets incluant ce type de détournement est pertinent.

Comme nous maîtrisons déjà le principe de détection de boucles inter-niveaux et que nous avons adapté notre algorithme à cet effet, la conception d'une méthode de détection de boucles dans les sections de rampe simple est en fait directe et rapide :

- On observe que les balises situées dans ces portions de rampe simple sont peu nombreuses et sont très espacées entre elles. Elles caractérisent donc chacune une augmentation (resp. une diminution) importante de la profondeur du HT ;
- On décide de considérer qu'elles feront office de balises de changement de niveau puisque, comme ces dernières, elles permettront d'identifier sans ambiguïté les demi-tours ayant lieu dans la rampe. Elles sont en effet trop distantes pour qu'un HT fasse une marche arrière de l'une à l'autre dans une situation ordinaire sans pouvoir se stationner. Elles auront bien sûr la particularité de ne pas donner accès à un niveau d'exploitation, mais elles jalonnent effectivement le parcours du HT dans la rampe.
- On les classe dans une liste ordonnée en fonction de leur profondeur respective.

- On applique le principe de détection de boucles inter-niveaux, présenté précédemment, aux séquences de détections de ces balises pour chacun des trajets $A \rightarrow B$ en cours d'étude. On vient ainsi d'ajouter à notre algorithme un critère d'« absence de boucle dans la rampe simple » à la suite du critère précédent.

4.5.6.4 Amélioration de l'algorithme par détection de boucles intra-niveaux

On s'intéresse maintenant aux « boucles intra-niveaux », i.e. des détours d'un autre type, durant lesquels le HT étudié est entré dans l'un des niveaux de la mine, puis en est ressorti avant de se diriger vers la balise B.

On rencontre alors une nouvelle problématique : ces détours peuvent faire partie de trajets ordinaires et de trajets non ordinaires. En effet, certaines conditions ordinaires peuvent mener le conducteur du HT à faire des boucles intra-niveaux et à y faire des pauses ponctuelles durant son trajet, lui-même ordinaire ; p. ex. les ravitaillements en carburant et les passages aux sanitaires puisqu'ils sont nécessaires, à court terme, au bon fonctionnement du HT et à celui de son conducteur. Ainsi, il est prévisible qu'un nombre considérable de trajets ordinaires puissent inclure naturellement des détours limités qui ne remettent pas en question le fait que le conducteur soit effectivement en train de se rendre jusqu'à la balise B ; et que, de surcroît, ce conducteur s'y rende efficacement puisque comme mentionné à la section 4.2, il est dans l'intérêt des conducteurs d'éviter les pauses et les détours superflus. Les autres boucles intra-niveaux ne seront généralement pas liées à un trajet ordinaire et devraient être considérées comme des anomalies. Les TT correspondants ne devraient donc finalement pas être retenus par notre algorithme. En résumé, nous pourrions être confrontés, sur un même itinéraire, à des détours ordinaires et non ordinaires ayant potentiellement lieu aux mêmes niveaux, dans le périmètre de détection des mêmes balises (en particulier à certaines balises du niveau 300 proches du garage) et dont les durées pourraient vraisemblablement coïncider.

Il semblerait que nous ne disposions pas de données historiques fiables et actualisées suffisamment régulièrement pour arriver à discerner correctement les boucles intra-niveaux ordinaires des boucles intra-niveaux non ordinaires ayant eu lieu dans la mine 1 durant toute la période étudiée. En théorie, si l'une des variables de télémétrie de la BDD de la mine 1 exprimant à intervalles réguliers le niveau de carburant de chaque HT était très régulièrement actualisée, il serait possible d'identifier les détours effectués pour un ravitaillement en carburant. Concernant les maintenances réactives (événement non-ordinaire), si les données de l'intégralité des bons de travail correspondant aux maintenances qui se sont déroulées sur la période étudiée avaient été transférés dans la BDD de la mine 1 et incluaient tous les renseignements concernant le type de maintenance, l'identifiant du HT en maintenance, la date et heure de

l'intervention ainsi que sa durée, il serait possible d'éliminer tous les trajets non ordinaires liés à des maintenances réactives. Il semble en revanche difficile de classer automatiquement certaines boucles intra-niveaux comme étant ordinaires ou non ordinaires, et encore plus délicat d'identifier clairement la raison, ordinaire ou non, de certaines boucles intra-niveaux (même si nous avons à notre disposition d'innombrables attributs supplémentaires dans la BDD de la mine 1) : passage aux installations sanitaires, maintenances non répertoriées, détour du conducteur qui aimerait aller chercher un objet quelconque ou parler à quelqu'un, etc. En somme, cette problématique est très vaste et complexe. Elle ne sera pas traitée dans le présent mémoire.

Finalement, la solution qui semble la plus cohérente consiste à reconnaître les boucles intra-niveaux et à mettre les trajets correspondants de côté, puisqu'ils forment un mélange de trajets ordinaires et de trajets non ordinaires à distinguer via des méthodes qui restent à développer, et ne devraient donc pas être oubliés. On s'intéressera donc aux trajets qui en sont dépourvus, i.e. les trajets que l'on pourrait qualifier de « directs », car sans véritable détour détectable. Il faut bien noter que le fait de mettre de côté ces trajets ne devrait pas réduire drastiquement le nombre de trajets observés sur les itinéraires prédominants identifiés dans la sous-section 4.5.4 pour la mine 1 puisque ces derniers avaient pour borne inférieure le niveau 300. Toutes les problématiques liées à des boucles intra-niveaux passant par les zones proches du garage de ce niveau ne nous concernent donc pas pour le moment puisque les HT arrivent toujours à destination avant d'avoir l'opportunité de se rendre au garage. Pour autant, nous avons livré toutes les explications permettant de savoir comment s'y prendre.

La reconnaissance des boucles intra-niveaux, quant à elle, est très simple. En effet, une fois que la balise initiale A a été détectée et qu'un nouveau trajet $A \rightarrow B$ est donc en cours d'étude, il suffit de vérifier, au fur et à mesure de la lecture des identifiants de balises ayant détecté le HT sur son trajet, qu'aucune balise interne à un niveau de la mine (i.e. dont le périmètre de détection ne chevauche théoriquement pas de la rampe) n'ait détecté le HT auquel on s'intéresse. Les identifiants des balises internes aux niveaux de la mine 1 sont facilement distinguables des identifiants des balises de la rampe (et de celles présentes à l'entrée de chaque niveau) puisque ces derniers possèdent un format particulier ou des mots-clés de localisation discriminants. L'automatisation de cette tâche ne constitue donc aucunement un obstacle pour la mine 1 et on l'intègre facilement à notre algorithme.

Au terme de cette sous-section, nous disposons donc d'un algorithme de reconnaissance de trajets capable d'isoler uniquement les trajets directs (en ignorant les trajets non ordinaires incluant une boucle inter-niveaux et les trajets indirects, ordinaires ou non, incluant une boucle intra-niveau) et de dresser l'histogramme des TT correspondants.

4.5.7 Analyse d'histogrammes de TT sur les trajets prédominants et améliorations résultantes

La présente sous-section vise à analyser les histogrammes de TT produits par l'algorithme de reconnaissance de trajets développé précédemment lorsqu'il est appliqué aux trajets prédominants identifiés à la sous-section 4.5.4, et à en déduire diverses adaptations intéressantes, p. ex. concernant notre algorithme ou les itinéraires à étudier.

Dès lors que l'on voit les histogrammes issus d'itinéraires passant par la rampe 1 uniquement, on remarque des distributions tout à fait atypiques des TT correspondant aux trajets en descente. Ces observations sont particulièrement évidentes sur certaines portions de l'itinéraire prédominant en descente sur la rampe 1 qui relie la balise d'accès en surface à la balise de changement de niveau qui correspond au niveau 300 ; on abrégera maintenant cet itinéraire en « itinéraire surface→niveau300 », et on appliquera ce même format d'abréviation aux autres itinéraires et trajets de la rampe 1. En particulier, la distribution des TT sur l'itinéraire surface→niveau200 met parfaitement en évidence l'existence d'un phénomène qui doit être étudié. On la représente à la figure 4.5.

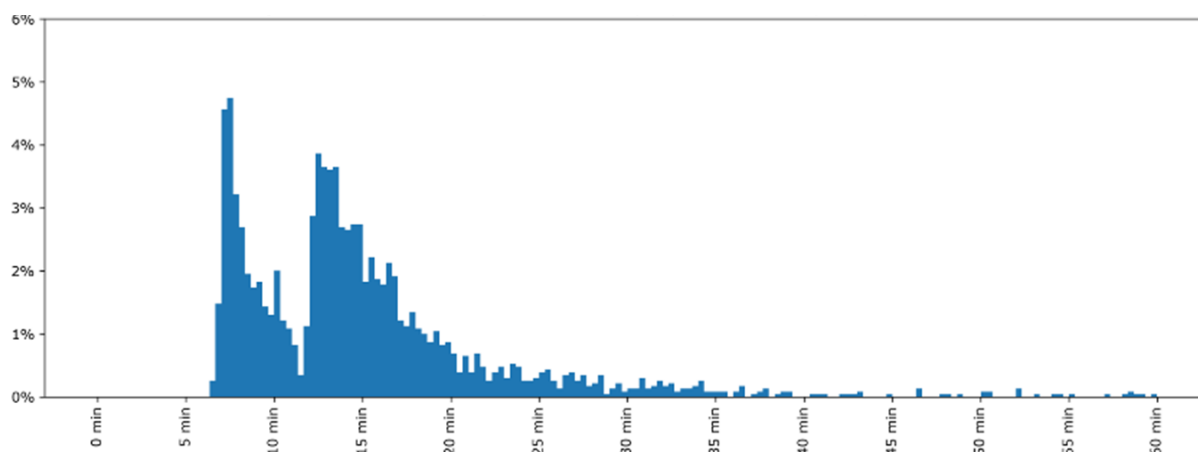


FIGURE 4.5 Histogramme des TT observés sur l'itinéraire surface→niveau200 dans la mine 1

Cette distribution, formée de plusieurs milliers de TT distincts calculés, a une allure qui s'éloigne fortement de celle d'une simple loi normale asymétrique. On distingue en effet deux modes dans la distribution, ce qui semble pointer vers l'existence de deux types de trajets bien différents au sein d'un même itinéraire. Pourtant, notre algorithme est supposé éliminer les types principaux de détours, i.e. théoriquement les boucles inter-niveaux et les boucles intra-niveaux. L'interprétation de ce phénomène intense ne coulant pas de source, une analyse approfondie est requise.

4.5.7.1 Étude approfondie de la distribution à deux modes

Pour essayer de comprendre le sens physique de cette distribution, visiblement peu exploitable en l'état pour prédire avec confiance les TT ordinaires sur cet itinéraire, on procédera par élimination. Dans un premier temps, on modifie séparément la balise A et la balise B de ce trajet dans les paramètres de l'algorithme, en demandant à ce dernier de représenter, pour chacune des combinaisons testées, l'histogramme des TT correspondant. Dans un deuxième temps, on rassemble d'une part toutes les combinaisons de balises A et B ayant visiblement provoqué une distribution à deux modes et d'autre part celles qui n'ont pas recréé le phénomène. Enfin, on analyse les regroupements obtenus afin d'étudier les caractéristiques communes des combinaisons qui ont provoqué ces distributions. Ainsi, il est très vite limpide que le fait de modifier la balise B n'a pas d'incidence notable sur l'existence ou non des deux modes, tandis que le fait de modifier la balise A fait quasiment toujours disparaître entièrement le phénomène, alors que le nombre total de trajets détectés varie très peu. Ces observations nous permettent donc de fortement soupçonner notre balise A initialement choisie d'être entièrement responsable de l'occurrence de ces distributions à deux modes. Cette balise A n'avait pourtant pas été choisie au hasard puisqu'elle est particulièrement intéressante du point de vue de l'architecture de la mine : elle correspond à l'accès en surface (ou « portail ») via la rampe 1.

Reste à investiguer sur le phénomène à l'origine de ces distributions à deux modes au niveau de la balise A pour savoir comment le filtrer et se ramener à une distribution se rapprochant de celle d'une loi normale asymétrique, ou bien pour établir les bonnes pratiques permettant simplement d'éviter l'apparition de telles distributions.

Pour ce faire, on utilise l'un des histogrammes faisant apparaître très clairement ce phénomène (p. ex. l'histogramme des TT de l'itinéraire surface→niveau200) et on relève le TT qui permet graphiquement de séparer au mieux les deux modes (i.e. qui correspond au creux qui sépare les modes). Dans le cas de l'itinéraire surface→niveau200, cette valeur se situe aux alentours de 11 minutes. Par la suite, on crée une version modifiée de notre algorithme purement destinée à analyser la problématique. Ainsi, on copie puis on modifie l'algorithme existant pour qu'il garde en mémoire, à chaque détection de la balise A, une séquence de détection correspondant aux trois balises qui précèdent chronologiquement la détection initiale de la balise A ainsi que les trois balises suivantes. Il s'agit uniquement là d'un choix arbitraire, qui vise à prélever des données de détection de quelques balises proches de A pour nous informer sur le contexte de la détection. Une fois que l'algorithme confirme que la détection de la balise A correspondait effectivement au début d'un trajet valide surface→niveau200, il calcule le TT correspondant puis classe la séquence de détections dans l'une des deux listes

prévues à cet effet, selon ce TT finalement calculé. S'il est inférieur ou égal à 11 minutes, on classe la séquence de détections dans la première liste ; sinon, on la classe dans la seconde. Finalement, on dispose de toutes les séquences de détections qui occasionnent l'apparition respective du premier mode et du second sur notre graphique. Ainsi, les deux types de trajets qui font apparaître chacun des deux modes ne devraient pas être confondus entre eux lors de l'analyse des listes de séquences de détections obtenues. Ceci est bien sûr rendu possible grâce au fait que les deux modes sont notablement différents.

On compare maintenant les listes obtenues, avant de passer à l'interprétation. La différence entre les deux listes saute rapidement aux yeux : dans la première liste, qui correspond au premier mode, une détection du HT par l'une des balises situées à la surface est observée avant que la balise A ne détecte le HT à son tour, ce qui débute le trajet surface→niveau 200 ; tandis que dans la seconde liste, les détections de balises qui précèdent celle de la balise A sont issues de balises situées dans la rampe, i.e. exactement les mêmes balises qui détectent à nouveau le HT lorsqu'il s'enfonce à nouveau dans la mine, sans qu'il ait jamais été détecté en surface dans ce cas de figure. Bien que l'on sache maintenant cela, la raison qui explique cette distribution déconcertante n'est toujours pas évidente à saisir. C'est parce qu'il manque une dernière information pour bien comprendre la source de ce phénomène. En fait, l'enregistrement des détections des balises dans la BDD de la mine 1 est programmé de sorte à ce que seules les détections qui apportent de nouvelles informations y soient sauvegardées. Plus explicitement, l'entrée d'un HT dans le périmètre de détection d'une balise donnée ne donnera lieu à l'enregistrement du lot de données de détection correspondant dans la BDD de la mine 1 **que si** la toute dernière détection de ce HT a été effectuée par une autre balise. Lorsque l'on exploite les lots de données de détection de la BDD de la mine 1, on doit donc considérer que les balises deviennent virtuellement aveugles à un HT donné au moment où elles le détectent, et restent aveugles tant que le HT n'a pas été détecté ailleurs (elles restent en revanche parfaitement attentives à la présence d'autres HT). Ce phénomène peut donner lieu à des activités/trajets non retraçables aux alentours des zones les moins équipées en balises : un HT, qui passerait successivement dans des zones non couvertes par les balises et dans le périmètre de détection de la balise à l'origine de sa toute dernière détection, peut circuler autant qu'il veut sans qu'il y ait de trace de l'une de ses localisations dans notre BDD, comme illustré à la figure 4.6. Appliqué à notre exemple, ce corollaire permet de modéliser la suite d'évènements à l'origine du second mode :

1. Antérieurement à un trajet surface→niveau200 donné, alors que le HT remonte sur la rampe 1 durant son trajet précédent avec sa cargaison de minerai à décharger en surface, il est détecté par la balise A (supposée incontournable lorsqu'un HT remonte par la rampe 1) dont le lot de données de détection qu'elle génère sera enregistré dans

la BDD ;

2. Il poursuit son trajet en surface jusqu'à la zone de déchargement, suivant un certain itinéraire (qui semble plutôt habituel au vu du volume de TT qui forment le second mode) qui a la particularité de ne le faire passer par aucun des périmètres de détection des rares balises de la surface ;
3. Il décharge sa cargaison ;
4. Il repart, prenant un itinéraire aux propriétés identiques (ce peut être le même itinéraire en sens inverse) pour se rendre de nouveau à la rampe 1 ;
5. Il passe sous la balise A, mais bien qu'elle détecte effectivement la présence de sa puce RFID, aucune donnée relative à cet évènement n'est enregistrée dans la BDD ; et
6. Finalement, il descend le long de la rampe, générant enfin de nouvelles détections.

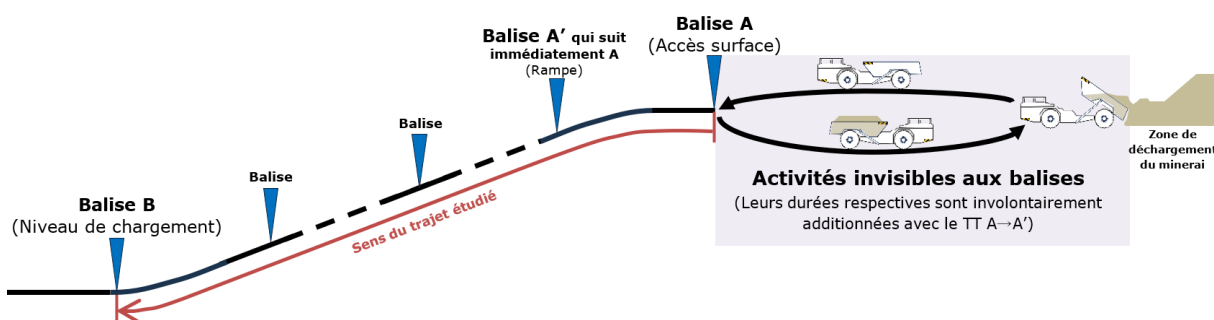


FIGURE 4.6 Modélisation d'une configuration de trajet qui rallonge le TT $A \rightarrow A'$, pouvant provoquer un second mode

Cette modélisation implique que le décalage temporel entre les deux modes observés devrait correspondre presque exactement à la différence entre, d'une part, la durée habituelle d'un trajet direct reliant la balise A à la balise qui la suit immédiatement dans la rampe (on la note « A' »), et, d'autre part, la durée médiane des segments de trajets indirects $A \rightarrow A'$ i.e. ceux qui correspondent à notre modélisation. On retrouve effectivement une durée similaire au décalage temporel des deux modes lorsque l'on calcule la différence entre la médiane des TT $A \rightarrow A'$ issus des trajets $A \rightarrow B$ du premier mode et la médiane de ceux issus des trajets $A \rightarrow B$ du second, ce qui achève de nous convaincre que notre modélisation explique conformément le phénomène observé.

Enfin, on représente l'histogramme des TT entre A et A' pour encore mieux visualiser l'aberrance des distributions de TT provoquées par ce phénomène (voir fig. 4.7).

Par ailleurs, notre modélisation implique aussi que le phénomène peut être observé n'importe où, y compris en souterrain et en montée, du moment qu'il existe des zones non couvertes

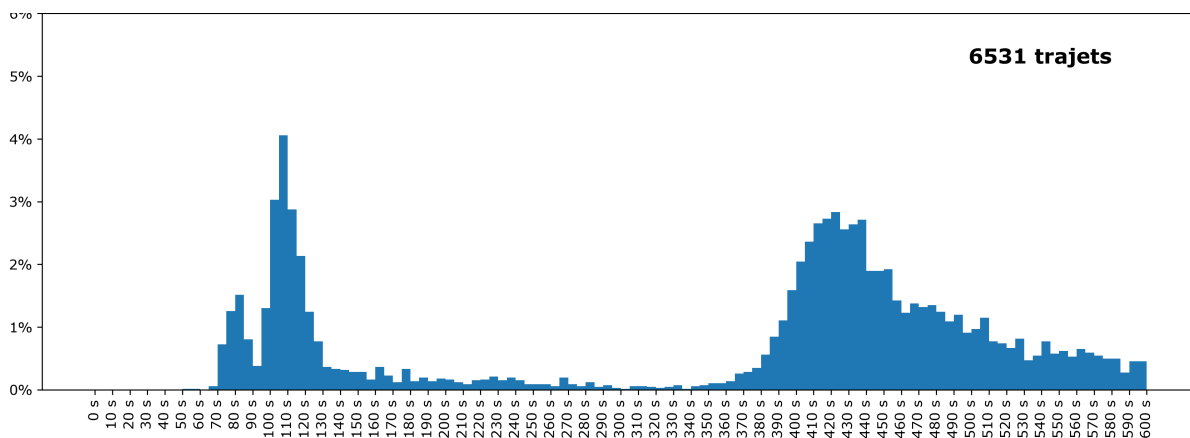


FIGURE 4.7 Histogramme des TT $A \rightarrow A'$ repérés par notre algorithme

par les balises, permettant au HT de rester intraçable.

Remarque importante, il se pourrait parfaitement que le second mode ne soit tout simplement pas observable si la quasi-totalité des trajets d'un itinéraire donné reproduisent ce phénomène d'activités non retraçables. On risquerait dans ce cas d'être confiant que l'histogramme des TT obtenu est fiable malgré le fait qu'il puisse être décalé d'une durée considérable, ce qui mettrait à mal toutes les prédictions de TT réels sans laisser présager de l'existence d'un tel phénomène. En d'autres termes, on doit toujours s'assurer que l'allure de nos distribution de TT n'est pas influencée par ce dernier pour éviter des erreurs de prédiction de TT réels considérables.

Notre compréhension fine du phénomène précédent ne nous donne malheureusement pas directement de méthode miraculeuse pour régler la situation qui a mené à une telle distribution. Toutefois, on peut réfléchir à plusieurs possibilités pour contourner le problème (à défaut de le supprimer complètement), mais chacune d'elles entraîne une concession. En effet, pour tous les trajets $A \rightarrow B$ pouvant théoriquement inclure un détour invisible, on peut :

- Calculer la durée du sous-trajet $A \rightarrow A'$ et la comparer à la distribution des TT de ce sous-trajet pour estimer, avec une confiance raisonnable, si l'on est en présence ou non d'activités invisibles, puis :
 - exclure de l'étude tous les trajets $A \rightarrow B$ que l'on a supposés contenir des activités invisibles, i.e. l'équivalent des trajets du second mode, mais cette solution n'est viable que s'ils sont peu nombreux puisqu'on veut maximiser le nombre de trajets retenus sur chacun des itinéraires étudiés (on perd un grand nombre de trajets réels dans l'opération) ; ou
 - corriger les TT supposés contenir des activités invisibles en leur retirant la du-

rée moyenne de ces dernières, mais cette solution est très peu recommandable pour étudier les TT réels puisqu'elle génère des approximations non quantifiables que l'on cherche vraiment à éviter (elle permet néanmoins de conserver le même nombre de TT).

- Choisir d'abandonner la balise A, géographiquement intéressante, mais peu commode à cause de ce phénomène, et s'intéresser plutôt aux trajets $A' \rightarrow B$, qui sont justement protégés du phénomène par la présence en amont de la balise A. Cette solution réduit la longueur du trajet étudié et pourrait réduire incidemment la valeur apportée par nos futures prédictions aux planificateurs miniers, mais elle permettrait de se débarrasser du phénomène indésirable et de garder un nombre de trajets étudiés très similaires, dont les TT n'auraient pas été approximatés. On choisit finalement cette option.

Concernant la solution que nous avons choisie, une subtilité additionnelle doit être évoquée : la balise A doit fonctionner continuellement sur la période étudiée pour être le foyer de tous les phénomènes indésirables d'activités invisibles, sous peine de laisser ces dernières contaminer temporairement les TT $A' \rightarrow B$. Pour parer à cette éventualité, on peut ajouter à notre algorithme une condition supplémentaire : vérifier que la balise A ait bien été détectée avant de s'intéresser au trajet $A' \rightarrow B$ qui la suit immédiatement.

Précisons enfin que la prévention de la problématique explicitée au paragraphe précédent n'est malheureusement pas une précaution exagérée : d'après les données historiques de la mine 1, la balise A n'a pas toujours été victime du décalage temporel lié à des activités invisibles. En effet, au début de la période étudiée, une autre balise située en surface détectait visiblement toujours le HT après qu'il soit passé par la balise A en montée, ce qui décalait la problématique d'activités invisibles à cette autre balise et permettait d'obtenir uniquement des durées de trajets directs $A \rightarrow B$, sans détours liés aux activités invisibles. Cette balise a visiblement été désactivée depuis et le volume de trajets détectés sur cette période reste assez faible, on n'envisage donc pas de se servir de ces trajets uniquement.

4.5.7.2 Investigation sur les TT irréalistiquement courts

En s'intéressant aux trajets $A' \rightarrow B$, une nouvelle anomalie toute autre est rapidement identifiée. Comme le montre la figure 4.8, certains histogrammes présentent une faible proportion de trajets exceptionnellement courts. Cette proportion réduite peut sembler anecdotique à première vue mais il n'en est rien : il n'est de toute évidence pas question ici de quelques trajets effectués en excès de vitesse mais bien d'observations virtuelles de déplacements physiquement impossibles, car quasi-instantanés (abrégié « Q-I »). En particulier, on relève un trajet de 500 mètres de dénivelé identifié comme ayant duré moins de 10 secondes par l'algo-

rithme ; à lui seul, il pourrait mériter de chercher la cause de ce phénomène pour éviter qu'il ne se reproduise.

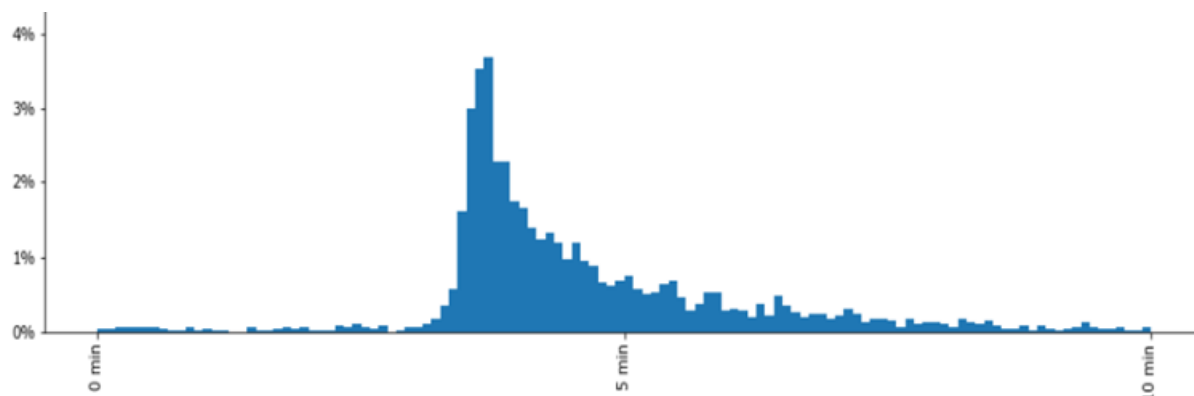


FIGURE 4.8 Histogramme des TT observés sur l'itinéraire A'→Niveau150, faisant figurer des TT Q-I

On décide alors d'intégrer à l'algorithme une nouvelle fonction qui vise à détecter et à renvoyer à l'utilisateur des données relatives aux trajets dont le TT calculé semble physiquement incohérent. Dans le cas précis de l'itinéraire de la figure 4.8, on peut considérer graphiquement que cette incohérence physique correspondra p. ex. à des TT inférieurs à deux minutes, pour garder une marge de sécurité. On applique ensuite l'algorithme à ce dernier itinéraire. Il est alors aisé de remarquer une corrélation significative entre l'apparition de ce phénomène et certaines périodes temporelles, lesquelles englobent toutes les occurrences observées. En lisant les séquences de détection correspondantes, directement issues de la BDD, on remarque que le phénomène est fréquent sur ces périodes puisque de nombreuses incohérences sont visibles et que des balises sans lien direct avec l'itinéraire étudié ici sont visiblement aussi concernées.

On décide donc de s'intéresser au cas général en créant un algorithme dédié à l'étude de ces phénomènes de déplacements apparemment Q-I via l'utilisation de toutes les données de détection de la BDD de la mine 1, pour ensuite mieux identifier le contexte d'apparition de ces déplacements virtuels et ainsi mieux les gérer. Cet algorithme nécessite d'établir un critère d'identification de déplacements Q-I, ce qui nous demande de préalablement expliciter ce qui constitue un TT Q-I dans notre contexte. Ainsi, étant donné que l'on aimerait maximiser la précision de détection de ces TT par notre algorithme (i.e. que l'on aimerait que chaque TT réputé Q-I corresponde effectivement à une erreur des données de détection et non à un trajet réel), on propose une définition notablement réductrice de cette notion, que l'on pourrait modifier si elle ne nous suffisait pas à globalement identifier les périodes d'occurrences du phénomène : arbitrairement, on décide qu'un TT Q-I correspondra nécessairement à un trajet de plus de 100 mètres de dénivelé parcouru en moins d'une minute. On précise

que parcourir 100 mètres de dénivelé en une minute semble impossible dans notre contexte puisque c'est moins de la moitié de la durée des trajets physiquement cohérents les plus courts ayant été observés en descente. Sur la base de cette définition, on peut programmer notre algorithme. Dans un premier temps, on liste toutes les dates incluses dans la période sur laquelle s'étendent nos données. On demande ensuite à l'algorithme de lire toutes les données historiques de détection de la mine 1, en calculant très simplement le TT qui sépare deux détections d'un même HT immédiatement successives et en déterminant si le TT obtenu est Q-I ou non au vu de la différence de dénivelé entre les deux balises immédiatement successives. Si une telle observation survient, l'algorithme incrémente le nombre d'observations de TT Q-I de la date étudiée. Si l'on a choisi de repérer uniquement les TT Q-I provenant de trajets entre balises immédiatement successives, c'est parce qu'il semble à première vue que des identifiants aléatoires de balises soient générés parmi les identifiants attendus. Les balises immédiatement successives devraient donc ici suffire à analyser le phénomène, mais même dans le cas contraire il n'aurait naturellement jamais été nécessaire d'étudier tous les itinéraires possibles de la mine, seulement quelques-uns parmi les plus fréquentés.

La liste du nombre de détections journalières de TT Q-I renvoyées par notre algorithme nous permet d'établir le constat suivant : au total, seuls 59 jours de détections sont victimes de ce phénomène, répartis en 3 périodes temporelles, mais ils peuvent être extrêmement affectés : notre reconnaissance partielle des TT Q-I permet de détecter de 70 à 2271 de ces TT incohérents durant chacun de ces jours, pour tous les HT confondus ; et aucun hors de ces jours. On découvre qu'ils contiennent 26% des données dont nous disposons, ce qui représente un volume disproportionné pour moins de deux mois de données de détection.

L'étude approfondie des détections que nous avons isolées nous fait découvrir que les détections « spontanées » ne sont pas complètement aléatoires comme on pourrait le croire de prime abord : lorsque les détections d'un HT donné sont victimes de TT Q-I, on peut distinguer dans ces détections plusieurs trajets se déroulant en parallèle dont les balises correspondantes s'alternent. Par ailleurs, et dans le même temps, certains HT cessent d'être détectés pendant des périodes inhabituellement longues.

Ces deux observations nous mènent à adhérer au scénario explicatif suivant : durant trois périodes distinctes, plusieurs véhicules ont été régulièrement et temporairement reconnus sous un même identifiant de HT, et leurs détections respectives au cours de leurs trajets respectifs se sont confondues. De plus, étant donné que le protocole d'enregistrement des données de détection dans la BDD de la mine 1 n'enregistre que les nouvelles balises rencontrées (comme explicité plus tôt dans la présente sous-section), le mélange des trajets sous un même identifiant de HT persuade le protocole que de nouvelles balises sont constamment

rencontrées (p. ex. on peut observer des cycles de détection d'un même HT à seulement deux balises différentes plusieurs dizaines de fois de suite en quelques secondes), ce qui justifie l'enregistrement de quantités importantes de détections redondantes durant ces 59 jours particuliers. Toutefois, la source du dysfonctionnement initial (plusieurs véhicules enregistrés sous un même identifiant de HT) reste inconnue.

Nous n'avons trouvé aucune méthode de nettoyage de ces données historiques fortement polluées, étant donné qu'il n'est plus possible de démêler avec assurance les vrais détections (celles du HT initial) des fausses détections (celles des autres véhicules), en particulier, car les véhicules se rapprochent et se croisent dans la mine, mêlant leurs détections et brouillant le retraçage des itinéraires empruntés. On décide donc d'éliminer ces trois périodes temporelles de notre étude. Pour ce faire, il suffit de modifier notre requête SQL d'extraction de données de la BDD de la mine 1 pour exclure de notre extraction les données de ces périodes, puis de ne travailler qu'avec ce nouveau jeu de données. Après cette suppression, il nous reste finalement plus d'un million et demi de détections de balises.

Par ailleurs, les successions de détections de balises pouvant être considérées comme partiellement aléatoires, on ne pense pas avoir manqué d'autres périodes de TT Q-I. En effet, on considère que si notre critère d'identification actuel n'a permis de détecter aucun TT Q-I sur une date donnée, il est très peu probable que d'autres formes de TT Q-I dus aux mêmes dysfonctionnements s'y cachent (p. ex. des déplacements trop rapides entre des balises espacées de moins de 100 mètres de dénivelés, ce que nous ne surveillons pas) étant donné que nous avons relevé un nombre élevé de TT Q-I sur chacune des dates problématiques identifiées.

4.5.7.3 Amélioration de l'algorithme par filtrage des trajets anormalement longs

Hormis les TT Q-I, d'autres TT anormaux sont observés sur la plupart des histogrammes obtenus : il s'agit de trajets anormalement longs, certains atteignant plusieurs heures. On admet qu'une majeure part de ces trajets provient de trajets non ordinaires. Ils ne sont pas extrêmement nombreux mais ils sont gênants puisque leurs durées exceptionnellement longues seraient jugées très significatives par nos futurs modèles prédictifs, ce qui nuirait à la phase d'apprentissage de ces derniers. On cherche donc idéalement à identifier leur cause pour ensuite déduire une solution efficace permettant d'exclure de notre étude tous les TT longs qui ne correspondent pas à des trajets ordinaires. À défaut, on exclura les valeurs aberrantes les plus élevées de TT de notre étude.

Notons d'abord que ces trajets anormalement longs sont relativement étonnants.

- D'une part, on a déjà éliminé de nombreux détours auparavant et, au vu de leur durée extraordinairement longues, il semble improbable que ces TT correspondent à des tra-

jets ordinaires directs. Ce sont pourtant théoriquement les seuls trajets que conserve maintenant notre algorithme.

- D'autre part, bien qu'une explication simplement basée sur l'occurrence de pauses d'une durée importante (i.e. plusieurs dizaines de minutes) au milieu de certains des trajets étudiés soit très séduisante pour expliquer le phénomène, il y a apparemment une incohérence puisque notre algorithme est conçu pour reconnaître les détours intra-niveau et exclure les trajets correspondants. Les longues pauses doivent donc avoir lieu directement dans la rampe. Cette conclusion semble assez incongrue puisqu'il s'agit d'un lieu de passage déjà étroit, et les chambres creusées dans les parois des rampes sont théoriquement dédiées aux croisements.

Après de plus amples analyses, on propose l'hypothèse complexe suivante.

- Une part importante des TT anormalement longs peuvent provenir d'arrêts volontaires de la part des conducteurs dans les quelques mètres ou dizaines de mètres qui suivent généralement l'entrée d'un niveau, puisqu'il y est plus facile de faire une longue pause que directement dans la rampe. Le HT serait alors uniquement détecté par la balise d'entrée du niveau, et aucune balise ne signalerait donc de détour. Le conducteur reprendrait ensuite normalement son trajet dans la rampe sans s'aventurer plus loin dans le niveau, ce qui empêcherait notre algorithme de détecter une boucle intra-niveau ;
- Une part non négligeable de nos observations pourrait bel et bien s'expliquer par des arrêts **volontaires** de conducteurs dans les chambres creusées dans la rampe, surtout dans les sections de rampe simple. En effet, ces dernières mesurent respectivement près d'un kilomètre et près de trois kilomètres ; de rares longues pauses seraient donc compréhensibles. Cette situation est particulièrement envisageable lors de quarts de travail durant lesquels peu de HT sont actifs dans une rampe donnée, ce qui laisse plus de flexibilité sur l'utilisation des chambres ; et
- Certaines de nos observations correspondraient enfin sûrement à de longues pauses involontaires ayant directement lieu dans la rampe ou dans les chambres qui y sont creusées. Elles seraient provoquées par un incident quelconque (panne du HT ou problème technique incapacitant, accident, rampe obstruée, etc.) qui générerait une attente potentiellement longue d'un ou plusieurs HT(s) entre deux balises de la rampe.

Précisons bien qu'une longue pause volontaire peut faire partie d'un trajet ordinaire ou non ordinaire selon la situation. Ainsi, bien qu'une pause correspondant p. ex. à une longue discussion inopinée entre le conducteur du HT et l'un des usagers croisés en chemin puisse être considérée comme faisant partie intégrante d'un trajet ordinaire, une pause qui correspondrait en revanche à un HT stationné durant une pause-repas fera partie d'un trajet non ordinaire.

Effectivement, le conducteur n'est alors pas en train d'essayer de se rendre activement à la balise B de l'itinéraire $A \rightarrow B$ étudié ou en aval, et il n'est pas non plus supposé être en train de le faire durant cette période du point de vue des planificateurs. Notons bien que ces considérations peuvent dépendre de la définition d'« ordinaire » que l'on adopte. De plus, selon le site minier étudié, la fréquence des différents types de pauses peut changer le caractère ordinaire ou non ordinaire de ces événements. D'après notre conception des conditions ordinaires dans lesquelles se déroulent les activités de transport de minerai dans la mine 1 et bien que l'on ne puisse pas le démontrer à l'aide de données historiques, on estime finalement que la part des trajets ordinaires provoquant des pauses particulièrement longues devrait être significativement faible par comparaison à celle des trajets non ordinaires, que l'on désire exclure de l'étude des TT.

Les circonstances à même de causer de longues pauses étant diverses, nous ne disposons finalement pas d'une méthode d'ingénierie des caractéristiques permettant d'utiliser les données disponibles pour identifier avec assurance la raison de chacune de ces pauses, pour ensuite distinguer les trajets ordinaires et non ordinaires. Aussi, nous finissons par nous résoudre à employer une méthode moins subtile : filtrer directement les valeurs aberrantes élevées de TT.

Notre méthode de filtrage ne sera toutefois pas totalement conventionnelle, car toutes les remarques précédentes nous permettent de l'améliorer en l'implémentant de manière plus réfléchie. Voici les diverses considérations qui guident notre choix de méthode :

- il faut maximiser le nombre de trajets non ordinaires exclus de notre étude et le nombre de trajets ordinaires inclus dans notre étude ;
- les longues pauses génèrent essentiellement des trajets non ordinaires ;
- ce sont les longues pauses qui génèrent la plupart des TT anormalement longs ; et
- une accumulation de très nombreuses courtes pauses ou ralentissements dans la rampe aurait plutôt tendance à témoigner d'un trajet ordinaire, même si le TT qui en découle pourrait être très long.

Tout d'abord, les deux dernières considérations impliquent qu'il est tout à fait possible que le même TT anormalement long puisse correspondre à deux trajets bien différents : l'un ayant été victime d'une longue pause et l'autre ayant inclus de très nombreuses courtes pauses ou ralentissements. Ces trajets risqueraient alors d'être collectivement inclus ou exclus de notre étude à la suite du filtrage des TT obtenus en sortie d'algorithme. Or, la combinaison des considérations précédentes nous indique que nous devons chercher à exclure les longues pauses et à inclure les courtes pauses. Par conséquent, nous proposons de ne pas filtrer les valeurs aberrantes élevées des TT inclus dans la liste de TT obtenue en sortie de l'algorithme. Nous

nous concentrerons plutôt sur l’exploitation de chacune des listes des TT entre deux balises de changement de niveau immédiatement successives (on qualifiera ces trajets d’« inter-niveaux », un terme qui était jusqu’alors dédié aux boucles). Ces dernières mettront en exergue les longues pauses et éviteront aux courtes pauses et aux ralentissements d’être perçus comme aberrants par notre filtrage. Les TT inter-niveaux incluront aussi les TT qui séparent deux balises immédiatement successives dans les sections de rampe simple, pour favoriser la concision des paragraphes suivants.

Ces TT inter-niveaux ne sont pour l’instant malheureusement pas collectés par notre algorithme. De ce fait, on décide de créer une nouvelle fonction qui sera spécialisée dans l’exécution de cette tâche et dans la création de listes correspondant chacune à l’un des itinéraires inter-niveaux et visant à stocker les TT observés sur chacun d’entre eux. Notre fonction s’intéressera successivement à tous les itinéraires inter-niveaux de la mine 1 et, pour chacun d’eux, analysera tout l’historique des détections de balises contenues dans notre jeu de données. Lors de la lecture de ces dernières, elle cherchera à reconnaître les trajets qui correspondent à l’itinéraire inter-niveaux en cours d’étude. Elle calculera les TT correspondants et les stockera dans la liste dédiée à l’itinéraire inter-niveaux en question. Enfin, elle déduira de tous ces TT le seuil de filtrage à partir duquel on devrait considérer un TT comme une valeur aberrante élevée sur l’itinéraire en question.

Reste à trouver l’approche la plus adaptée pour fixer ce seuil pour chacun des itinéraires inter-niveaux. Plusieurs approches différentes peuvent en effet être envisagées. Toutefois, il est particulièrement difficile de trancher dès maintenant en faveur d’une formule invariable étant donné la diversité des distributions de TT inter-balises. Nous chercherons à vérifier la pertinence des différentes approches envisageables à la sous-section 4.5.11, et nous utiliserons d’ici-là un critère de filtrage classique basé sur la suppression des 5% des données les plus élevées. Quelle que soit l’approche adoptée, la valeur de seuil obtenue sera placée dans une autre liste, avec toutes les autres valeurs de seuils de filtrage.

Les étapes précédentes devront être effectuées pour tous les itinéraires inter-niveaux, à la fois en descente et en montée. Chaque seuil de filtrage obtenu devra être comparé manuellement avec l’histogramme des TT inter-niveaux correspondant, de sorte à ce que l’on évite de fixer des seuils de filtrage problématiques. En effet, certains de ces histogrammes pourraient comporter par exemple une distribution de TT à deux modes marqués si l’on n’a pas vérifié suffisamment sérieusement auparavant que les balises qui avoisinent les balises de changement de niveau empêchent ce phénomène indésirable de se produire. Dans ce cas, la formule précédente ne permettrait pas de filtrer le second mode puisque l’écart interquartile pourrait correspondre à la distance séparant les deux modes, ou puisque $Q1$ et $Q3$ pourraient tous

deux se situer dans le second mode.

L'utilisation de la liste précédente se fera ensuite directement dans notre algorithme principal. Au cours d'un trajet sur un itinéraire $A \rightarrow B$ donné, il vérifiera à chaque nouvelle balise de changement de niveau rencontrée (et chaque balise située dans la rampe simple rencontrée) que le TT inter-niveaux qui vient de s'écouler depuis la détection de la balise de changement de niveau précédente, ne dépasse pas le seuil de filtrage qui correspond à ce trajet inter-niveaux dans la liste générée par la fonction précédente. Si les balises de changement de niveau ne sont pas fiables, il faudra :

- D'une part, conserver ce seuil de filtrage inter-niveaux entre les détections effectives de balises de changement de niveau immédiatement successives ; et
- D'autre part, lorsque l'une ou plusieurs des balises de changement de niveau ne fonctionne(nt) pas, ajouter un seuil de filtrage nettement plus flexible entre les deux balises de changement de niveau non immédiatement successives afin d'absorber les variances cumulées des multiples trajets inter-niveaux inclus entre ces deux balises.

Les histogrammes de TT finalement tracés par notre algorithme seront en grande partie nettoyés des TT non-ordinaires qui s'y trouvaient. Quelques valeurs de TT anormalement longues pourraient quand même y subsister, mais, malgré leur aberrance, il devrait vraisemblablement s'agir d'observations légitimes de TT ordinaires qui doivent être conservées pour prendre en compte les conditions les plus difficiles dans lesquelles peuvent se dérouler ces trajets.

4.5.7.4 Comparaison d'histogrammes de TT en montée et en descente

La comparaison des histogrammes de TT ordinaires directs en montée et en descente ne se fait qu'à ce stade, car tous les développements précédents ont permis d'améliorer la fiabilité des histogrammes pouvant être obtenus au moyen de notre algorithme.

Attention toutefois, pour rappel, les distributions affichées ont toutes subies un filtrage très significatif (seuil des 5% sur **chaque** segment inter-niveaux) des TT les plus longs, pour que les TT non ordinaires ne polluent pas nos observations des distributions.

On représente les histogrammes respectifs de deux itinéraires inverses i.e. $A' \rightarrow \text{Niveau200}$ (en descente) et $\text{Niveau200} \rightarrow A'$ (en montée) sur les figures 4.9 et 4.10.

Sur ces deux histogrammes correspondant à deux itinéraires de longueur identique, plusieurs points de comparaison sont notables. On les mentionne et les justifie ci-dessous :

- **Allures générales distinctes** : les différences fondamentales entre les trajets en descente et en montée sont ici mises en relief. Alors que les HT qui descendent sont aidés

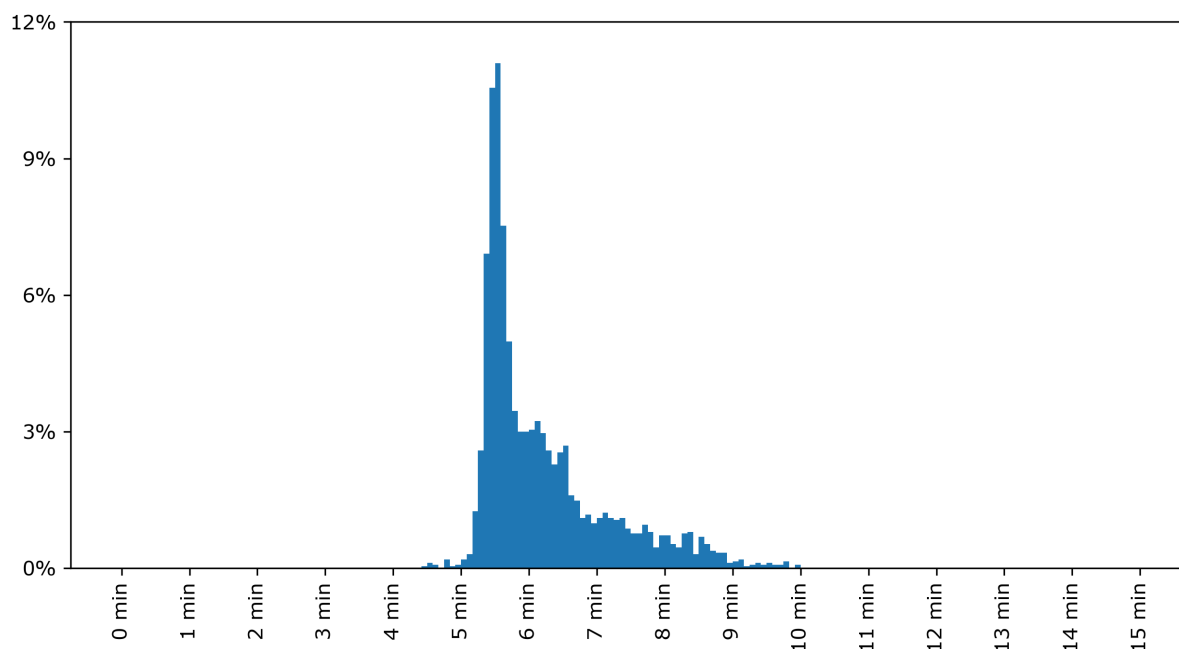


FIGURE 4.9 Histogramme des TT observés sur l'itinéraire A'→Niveau200 dans la mine 1

par la pente et ne sont pas chargés (donc plus maniables) les HT qui remontent sont handicapés à la fois par leur lourde cargaison et par la pente. Ces derniers sont prioritaires et subissent peu de ralentissements aléatoires, tandis que les autres doivent céder la priorité tant bien que mal dans la rampe à voie unique. Ces éléments suffisent à expliquer la différence d'allure générale entre les graphes précédents, qui témoignent de dynamiques éloignées malgré des moyennes respectives de TT étonnamment moins contrastées que l'on ne pourrait l'imaginer.

- **Grand nombre de trajets particulièrement rapides en descente** : contrairement aux trajets en montée, dont les conditions de chargement et de pente limitent inéluctablement l'accélération et la vitesse maximale des HT, les trajets en descente peuvent offrir des conditions bien plus clémentes à la prise de vitesse, éventuellement sur toute la longueur de la rampe à supposer que le HT concerné ne croise pas d'autres véhicules. Ainsi, l'important pic visible dans la distribution des trajets en descente serait en fait dû à un effet de seuil, puisqu'il correspondrait aux trajets durant lesquels aucun croisement n'a eu lieu. En effet, hormis évidemment lors de croisements avec d'autres véhicules, les HT qui descendent peuvent théoriquement rouler à une vitesse proche de la vitesse maximale autorisée avec très peu d'aléas, d'où un grand nombre de TT rapides et quasiment identiques qui correspondent au pic imposant tout à gauche de la distribution.

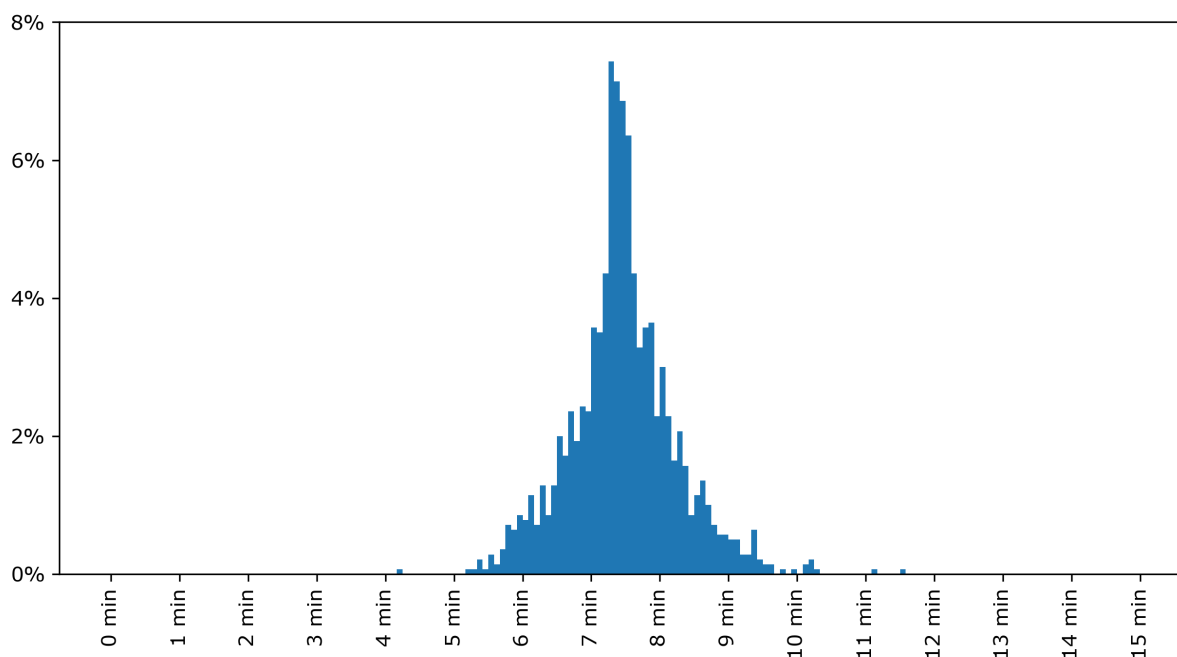


FIGURE 4.10 Histogramme des TT observés sur l'itinéraire Niveau200→A' dans la mine 1

- **Distribution ressemblant considérablement à celle d'une loi normale en montée** : cette allure indique en particulier que, contrairement à l'itinéraire inverse, il n'y a pas ici d'effet de seuil. Les TT en montée semblent donc varier moins brutalement que les TT en descente. Ils sont sûrement bien moins influencés que ces derniers par la vitesse maximale autorisée dans la rampe et sont en revanche plus affectés par plusieurs variables continues comme la masse de minerai chargée, le kilométrage du HT ou encore le niveau d'expérience et le style de conduite du conducteur. Les TT varient ainsi plutôt symétriquement autour d'un TT moyen fréquent, résultant en une distribution notablement proche de celle d'une loi normale.
- **Différence importante du nombre de trajets** : bien que non visible directement sur les histogrammes, c'est une constatation déstabilisante puisque l'on détecte en effet deux fois plus de trajets ordinaires directs (après filtrage) en descente qu'en montée sur cette section (resp. environ 3400 et environ 1700 sur notre période d'étude). Pour expliquer ce phénomène, notons d'une part que l'on détecte près de 30% de trajets quelconques (ordinaires ou non) supplémentaires en descente qu'en montée sur cette section (resp. environ 6500 contre environ 5000), ce qui semble indiquer que cette rampe n'est pas utilisée aussi intensivement dans les deux sens de circulation (l'existence de la rampe 2 permet aux HT qui descendent via la rampe 1 vers des niveaux profonds d'ensuite remonter via la rampe 2). D'autre part, toutes nos données concernant les

trajets ordinaires directs en montée indiquent que les camions font effectivement des détours variés dans les niveaux lors de leur remontée, et que ces détours ne sont pas attribuables à des erreurs de détection puisque les HT se déplacent considérablement dans ces niveaux à chacun de leurs détours. Il faut toutefois bien noter que nous ne détectons pas très exactement la totalité des trajets ordinaires directs ayant lieu sur ces itinéraires. En effet, lors d'un trajet donné dans la mine 1, nous imposons au HT d'être détecté par chacune des balises de changement de niveau situées sur l'itinéraire suivi car celles-ci sont visiblement fiables. On peut ainsi vérifier si le trajet du HT respecte toutes les propriétés que nous avons associées aux trajets ordinaires directs, et en déduire si nous le retenons ou non dans notre étude. En revanche, cela induit que chacune des « non-détections » de balises de changement de niveau peut supprimer le trajet correspondant de notre étude même s'il est ordinaire direct. Finalement, il nous manque sûrement quelques-uns des trajets ordinaires directs de chaque itinéraire, mais nous ne pouvons pas détecter plus de trajets ordinaires directs dans la mine 1 en gardant la fiabilité actuelle de nos détections. Par ailleurs, la suppression de ces quelques trajets est globalement aléatoire. Notre échantillon des TT est donc très légèrement réduit, ce qui ne cause aucune problématique importante.

On représente aussi les histogrammes de TT correspondant aux itinéraires de la rampe 2 reliant la surface au niveau 300 et inversement sur les figures 4.11 et 4.12. Ces trajets sont plus longs et se déroulent sur une longue section de rampe simple, ce qui modifie la dynamique des trajets et le nombre potentiel de conflits de trajectoires des HT. On observe une différence considérable de la médiane des TT de chacun de ces histogrammes, mais leur allure générale correspond effectivement aux observations précédentes. On arrive à détecter un très grand nombre de trajets ordinaires directs sur cette rampe avec resp. environ 18000 trajets ordinaires directs en descente et environ 22000 trajets ordinaires directs en montée sur notre période d'étude de 11 mois.

4.5.7.5 Comparaison d'autres histogrammes

On termine cette sous-section en réalisant les comparaisons d'histogrammes qui semblent pertinentes dans le contexte du site minier étudié. Bien que cette étape de notre méthodologie soit facultative, nous n'aurons pas d'autres occasions de faire ce travail alors qu'il permet d'accroître la compréhension des enjeux liés aux particularités de la mine en question.

Dans le cas de la mine 1, nous jugeons pertinent de comparer les TT observés dans chacune des deux rampes sur la section la plus longue possible et de même longueur, i.e. entre la surface et le niveau 300. Cette double rampe d'accès est en effet une particularité notable de

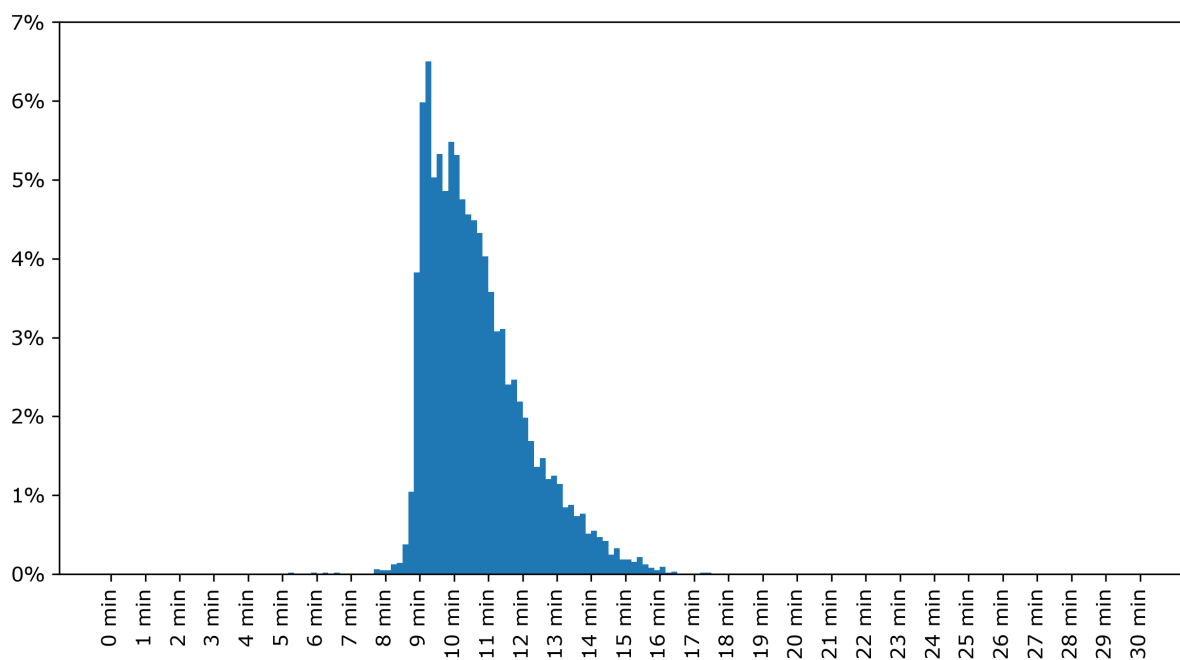


FIGURE 4.11 Histogramme des TT observés sur l'itinéraire Surface→Niveau300 via la rampe 2 dans la mine 1

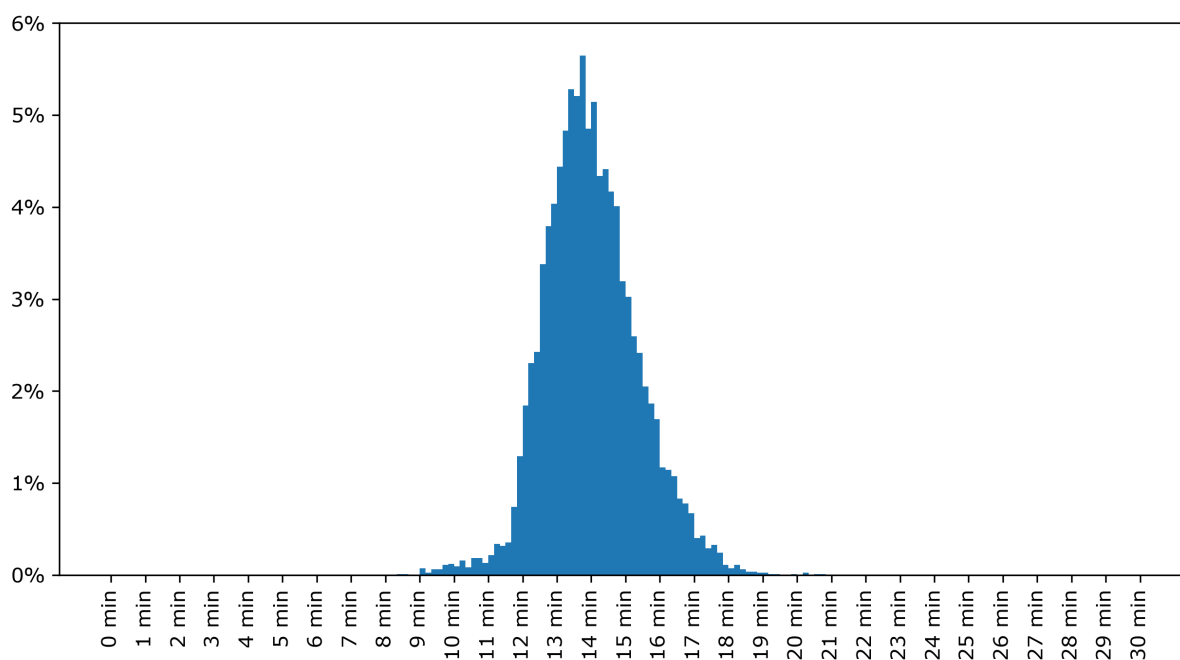


FIGURE 4.12 Histogramme des TT observés sur l'itinéraire Niveau300→Surface via la rampe 2 dans la mine 1

cette mine et il peut être intéressant de comparer les différences entre les trajets se déroulant sur chacune des deux rampes.

La problématique rencontrée avec la balise A de la rampe 1 nous empêche malheureusement d'exploiter les itinéraires en descente partant de la balise A, et il n'y a pas d'équivalent exact de la balise A' dans la rampe 2 pour comparer des trajets démarrant à un niveau de profondeur identique. Aussi, on comparera les TT observés en montée sur ces longs itinéraires puisque rien ne nous en empêche.

On adapte notre algorithme à ces nouveaux itinéraires, puis on récupère les histogrammes de TT correspondants.

Les figures 4.13 et 4.12 sont respectivement les histogrammes des deux itinéraires menant directement les HT du niveau 300 de la mine 1 vers la surface, l'un via la rampe 1 et l'autre via la rampe 2.

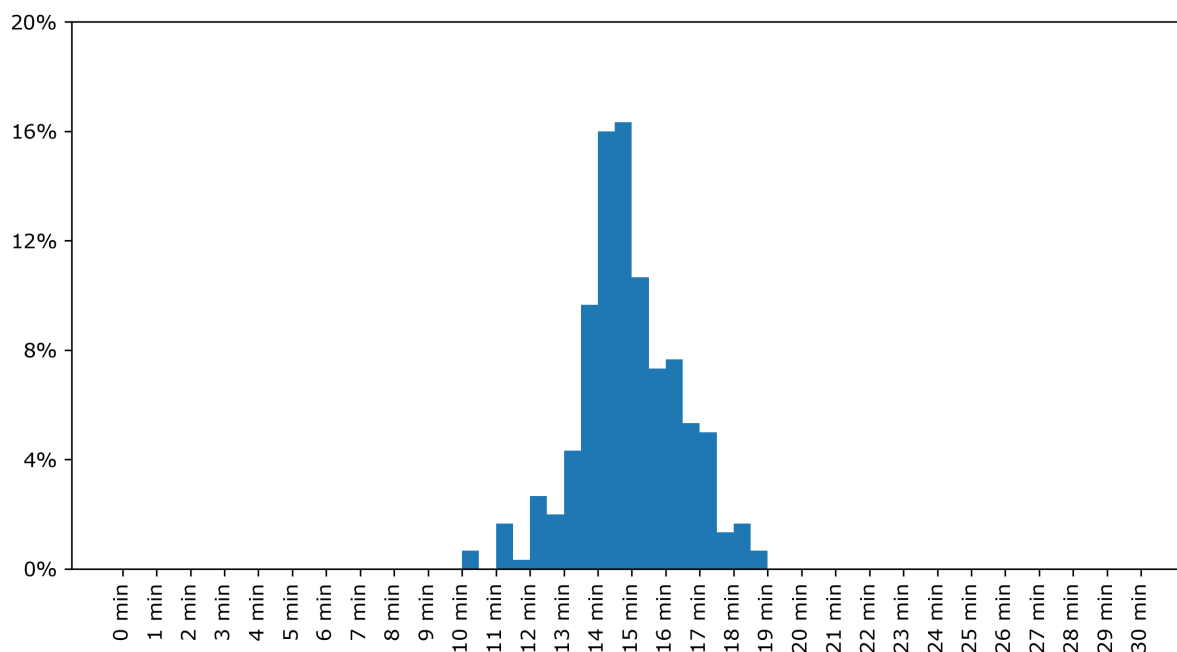


FIGURE 4.13 Histogramme des TT observés sur l'itinéraire Niveau300→Surface via la rampe 1 dans la mine 1

Une seule propriété relative à la détection de trajets par notre algorithme est particulièrement frappante : il s'agit de la différence entre le nombre de trajets ordinaires directs sur ces deux itinéraires puisqu'on ne dispose que d'environ 500 trajets ordinaires directs pour la rampe 1, contre près de 6000 pour la rampe 2. Pourtant, comme nous l'avons fait auparavant, on a ici vérifié à quel point les détours observés dans la rampe 1 étaient indubitablement de

vrais détours dans les niveaux rencontrés, ce qui confirme le fait que très peu de HT utilisent couramment la rampe 1 pour remonter vers la surface une fois qu'ils se sont chargés dans les niveaux les plus profonds de la mine et que la rampe 2 est clairement privilégiée dans ce cas. En revanche, ils remontent régulièrement la rampe 1 pour se diriger vers le niveau 275, 250, 225 ou 200 et s'y déplacer considérablement. On pourrait imaginer qu'ils sont alors à vide et qu'ils proviennent de la rampe 2, avec pour objectif de se charger dans les niveaux qui précèdent le niveau 300 dans la rampe 1, ce qui offre une grande panoplie de trajets possibles.

Les histogrammes indiquent que le temps médian de remontée via la rampe 2 est plus court que via la rampe 1 mais ce n'est pas un raccourci particulièrement impressionnant. On pourrait donc supposer que la fréquentation très variable de chacune de ces rampes en montée ordinaire directe s'explique par diverses raisons, p.ex. une volonté de fluidification des trajets dans chacune des deux rampes en limitant les croisements de HT via plus de trajets en montée dans la rampe 2, une plus grande difficulté des longs trajets en montée dans la rampe 1 (constamment en virage, avec de nombreuses intersections d'entrées de niveau et une faible visibilité) par rapport à la rampe 2 (constituée de longs morceaux rectilignes, sans intersections avec une meilleure visibilité), ou encore un lieu plus favorable de débouchage en montée de la rampe 2 à la surface.

4.5.8 Étude de déplacements erratiques de HT sur de faibles distances

Étant donnée la grande variété d'opérations de nettoyage de données effectuées dans les sections précédentes sur les données de la BDD de la mine 1, et étant donnée de surcroît la cohérence appréciable des histogrammes des figures 4.9 et 4.10, on pourrait s'attendre à ce que les données isolées soient globalement exemptes de bogues. Pour autant, en portant une attention redoublée aux séquences de détections de balises dans notre jeu de données, il est encore possible de discerner des anomalies, visiblement récurrentes dans les données la mine 1. En effet, en mettant en relation les séquences de détections de HT avec le plan minier sur lequel figure les balises, on observe des déplacements virtuels qui paraissent erratiques sur des distances très faibles (allant d'une dizaine de mètres de galerie à un peu plus d'une centaine de mètres). Ces distances sont bien plus faibles que pour les trajets Q-I tels que nous les avons définis avec notre critère de plus de 100 mètres de dénivelé entre deux balises en une minute. Ces anomalies ne devraient pas être issues du même phénomène puisque nous avons éliminé les périodes durant lesquels les trajets Q-I étaient très certainement les plus fréquents. De plus, nos observations sont notablement différentes : virtuellement, les HT semblent constamment avancer et reculer ou bien faire de multiples demi-tours entre les balises d'un même niveau, et semblent par ailleurs faire de nombreux détours difficilement justifiables dans certaines

galeries à côté desquelles les itinéraires les plus directs et les plus prédominants sont établis. Étant donné que l'on s'intéresse par ailleurs à des véhicules lourds et peu manœuvrables dont les conducteurs ont généralement une intention de haute productivité (voir sous-section 4.2.1), l'allure de ces déplacements semble encore moins naturelle. Tout ceci nous mène à considérer que le traitement des détections brutes de HT par le système de balises de la mine 1 n'est hélas sûrement pas aussi fluide qu'espéré.

Notre modèle explicatif du phénomène se base sur deux problématiques entrant simultanément en jeu. Tâchons de les décomposer.

- D'une part, le périmètre de détection des balises de la mine 1 est particulièrement étendu. En effet, on observe des séquences de détections qui indiquent qu'un HT donné peut régulièrement être détecté quasi-simultanément par deux balises pourtant situées à plus d'une centaine de mètres l'une de l'autre, sans que l'on puisse donc savoir quelle est la position réelle du HT. Cette situation n'est visiblement rencontrée qu'en l'absence de paroi capable de gêner le champ de détection des balises, ce qui appuie la thèse de périmètres de détection trop étendus plutôt que de bogues fréquents. Il est possible que le périmètre de détection des balises ne soit simplement pas réglable, mais il est aussi tout à fait vraisemblable d'imaginer que l'angle du cône de détection de chaque balise a volontairement été réglé très large pour s'assurer de la détection de tous les HT qui passeraient à proximité, et même plus loin aux alentours, sans forcément prévoir que ce réglage réduirait la fiabilité des détections de HT : il est en effet encore moins valide dans ces conditions de supposer naïvement que la position réelle d'un HT correspond à la position de la balise qui l'a détecté puisqu'il pourrait éventuellement se trouver à près d'une centaine de mètres de cette dernière au moment de la détection. Cette première problématique ne suffit pas à expliquer certaines observations bien qu'elle soit apparemment nécessaire pour expliquer les séquences de détection observées ; et
- D'autre part, lorsqu'un HT chevauche les périmètres de détection de deux balises, celles-ci entrent visiblement en conflit plutôt que de céder la priorité. En effet, les séquences de détections décrivant virtuellement un trajet durant lequel un HT avance et recule successivement ou fait de petits allers-retours correspondent en fait vraisemblablement au cas de figure concret suivant :
 - Un HT arrive dans une zone couverte par deux balises à la fois, qui détectent sa présence à intervalles réguliers ;
 - La balise dont il est nouvellement à portée le détecte pour la première fois ;
 - Cette détection inédite est enregistrée dans la BDD de la mine 1 selon le protocole habituel d'enregistrement des détections (voir sous-section 4.5.7) ;

- La balise précédente, qui continue à détecter le HT, émet une détection qui est maintenant différente de celle enregistrée en toute dernière dans la BDD ;
- Elle paraît inédite au protocole d’enregistrement des détections, qui l’enregistre donc naïvement dans la BDD ; et
- Les quatre étapes précédentes se répètent jusqu’à ce que le HT quitte la zone de chevauchement des deux périmètres de détection.

Bien qu’il puisse sembler curieux qu’un tel cas de figure ait été laissé au hasard lors de la conception du protocole d’enregistrement des détections des balises, le phénomène qui en résulterait est entièrement cohérent avec les observations que l’on peut faire sur les détections dont on dispose maintenant, lesquelles ont été notablement nettoyées et épurées, ce qui laisse finalement peu de place à d’autres interprétations.

Ces deux problématiques se renforcent puisque l’étendue considérable des périmètres de détection favorise l’apparition de zones de chevauchement plus vastes de ces derniers, lesquelles brouillent davantage les enregistrements des détections dans la BDD. Hormis le fait que la qualité des détections historiques finalement disponibles en soit nécessairement affectée, d’autres conséquences mineures peuvent être évoquées. En effet, cette combinaison de problématiques semble augmenter considérablement le nombre de détections finalement enregistrées, ce qui accroît inutilement le volume de données stockées dans la BDD de la mine 1. Aussi, tous les algorithmes de préparation de données que l’on applique ensuite à ces dernières sont ralentis par ce volume supplémentaire de données redondantes.

Avec notre compréhension plus approfondie de ces détections indésirables, il s’agirait d’expliquer pourquoi nos histogrammes semblent visuellement peu affectés par celles-ci afin de savoir comment prévenir la situation inverse. Fondamentalement, cette explication réside dans le fait que toutes les sections d’itinéraires auxquelles nous nous sommes intéressés directement (i.e. seulement les itinéraires complets en eux-mêmes) et indirectement (i.e. les portions de rampe entre deux balises de changement de niveau) s’étendent sur des distances considérables et se situent de surcroît essentiellement dans la rampe 1. Celle-ci est hélicoïdale entre les niveaux successifs, ses courbes et ses parois limitent donc efficacement le périmètre de détection des balises. De plus, la très grande distance qui sépare les balises de la rampe 2 élimine aussi efficacement le phénomène. Ainsi, il ne peut pas y avoir de conflits de détection entre les balises particulières auxquelles nous nous intéressons (i.e. la balise A’ et les balises de changement de niveau) malgré le périmètre de détection visiblement très étendu des balises, le suivi des trajets reste donc plutôt fiable. Une seule subtilité aurait pu persister : lorsque le HT descend la rampe et passe devant l’entrée d’un niveau sans y entrer, les périmètres de détections étendus des autres balises implique qu’il pourrait théoriquement être détecté par

celles-ci, en plus de celle qui se trouve précisément à l'entrée de ce niveau. Dans ce cas, notre algorithme comptabiliserait ce trajet dans les trajets avec détour intra-niveau et l'éliminerait des trajets étudiés. Toutefois, mis à part au niveau 300 et 350, la disposition des balises rend cette situation improbable puisque la rampe est un angle mort pour toutes les balises du niveau non situées à l'entrée de ce dernier d'après notre plan de la mine 1. Au niveau 300 et 350 en revanche, on doit empêcher notre algorithme de considérer les détections d'une autre balise comme étant la preuve d'un détour intra-niveau. Pour que l'algorithme ignore leurs détections, il suffit d'ajouter ces balises à la liste de celles qui sont censées pouvoir détecter les HT lors de leurs trajets **directs** se déroulant uniquement dans la rampe. Les détections des autres balises continueront à être assimilées à des détours intra-niveau par l'algorithme, qui ignorera donc le trajet correspondant.

Si nous avons adopté une autre approche concernant le sectionnement de nos itinéraires, notre algorithme aurait pu être en revanche complètement désorienté par les problématiques précédentes. En particulier, si l'on s'était intéressé à un itinéraire très balisé se déroulant le long d'une galerie rectiligne avec de nombreuses galeries perpendiculaires, et si l'on avait décidé de surcroît de segmenter l'itinéraire complet en utilisant chacune des balises supposées détecter le HT sur son trajet ; alors on cumulerait de l'une des manières les plus défavorables qui soient les effets des deux problématiques, avec des périmètres de détection très peu restreints par les parois (comparativement aux trajets dans la rampe) et donc de multiples zones de chevauchement de périmètres de détection successives à même de recréer en boucle le phénomène virtuel des déplacements erratiques. C'est d'ailleurs un cas similaire qui nous a amené à observer les déplacements erratiques remarqués initialement. Dans cette configuration, l'algorithme qui aurait été initialement développé serait potentiellement incapable de reconnaître un seul des trajets correspondant à cet itinéraire en raison de la confusion totale des détections enregistrées successivement dans la BDD : pour chaque trajet historique étudié, un algorithme naïf identifierait un détour dès la première détection qui ne suit pas parfaitement la succession prévue des balises qui devraient être exclusivement rencontrées d'après le plan de la mine 1. Il faut pourtant bien noter que le type d'itinéraire décrit dans ce paragraphe est classique et semble à première vue particulièrement simple à étudier. Si nos deux problématiques n'existaient pas, la méthode de sectionnement utilisant toutes les balises successives semble bien pertinente puisque chacune d'entre elles génère une information contextuelle supplémentaire qu'il pourrait être pertinent d'exploiter. L'étude de tels itinéraires via une telle méthode n'est donc pertinente que si les données de détection disponibles sont d'une fiabilité particulièrement élevée, ce qui n'est pas le cas ici. De plus, l'ajustement de cet algorithme demanderait d'accorder moins d'importance aux détections imprévues. Or, de tels itinéraires sont plutôt courts dans la mine 1 et les balises étant de

surcroît capables de détecter à une distance considérable les HT circulant dans les galeries rectilignes, on démultiplierait inévitablement les chances de l'algorithme de comptabiliser des trajets virtuels, qui n'ont donc jamais eu lieu en réalité. On pourrait p. ex. penser au cas d'un HT qui ne se déplacerait concrètement pas sur l'itinéraire étudié, mais qui opérerait dans la même zone en cumulant les détections de nombreuses balises, qu'elles soient avoisinantes ou situées à bonne distance de lui. Finalement, les trajets qui suivent un tel itinéraire sont exceptionnellement difficiles à étudier rigoureusement dans notre contexte comparativement à ceux qui suivent un itinéraire le long de la rampe. L'étude des trajets dans la rampe semble véritablement plus faisable ici, d'autant plus qu'ils correspondaient aux trajets prédominants. Notre approche par balises de changement de niveau semble par ailleurs bien plus robuste que l'approche visant à exploiter chacune des balises rencontrées durant le trajet.

4.5.9 Génération de variables d'intérêt temporelles

On génère dans un premier temps une variable permettant d'identifier le quart de travail durant lequel a eu lieu chaque trajet. On croise pour cela l'horodatage de la détection du HT pendant le trajet et la plage horaire de chacun des quarts. Dans la mine 1, ils sont au nombre de deux, respectivement diurne et nocturne, et on connaît la plage horaire sur laquelle ils s'étendent précisément. Si un HT a débuté son trajet plus d'une demi-heure avant le début du quart de travail (resp. terminé son trajet plus d'une demi-heure après la fin), on préfère supprimer le trajet correspondant pour limiter le nombre de trajets qui ne se sont pas déroulés dans les conditions opérationnelles d'un quart. Sinon, on place le trajet dans le quart correspondant.

Via ces mêmes plages horaires, on peut aussi créer la variable que nous avons dénommée « position temporelle du trajet relativement au quart de travail ». Elle sera égale au nombre de minutes écoulées depuis le début du quart (plus précisément depuis la demi-heure précédente) lorsque le HT observé débute chaque trajet, quel que soit le quart.

Nous décidons aussi de prélever l'indice du jour de la semaine durant lequel se déroule chaque trajet. À partir d'un horodatage passé au format `datetime` de la bibliothèque `pandas`, la conversion est instantanée via la méthode `isoweekday` du module `Timestamp` de cette même bibliothèque.

Ajoutons la dernière variable temporelle, l'indice du jour de l'année en cours au moment du trajet. De même qu'au point précédent, il suffira d'utiliser la bonne méthode issue de `pandas`, i.e. `dateofyear`.

Les trois dernières variables que nous avons décrites, qui peuvent être assimilées à des va-

riables continues et qui décrivent un phénomène cyclique, devront être pertinemment pré-traitées pour pouvoir être exploitées correctement par nos modèles de ML. Nous nous y emploierons dans la section 4.7.

4.5.10 Estimation du nombre de HT se déplaçant activement sur un itinéraire donné durant un quart de travail donné

Afin d'évaluer l'influence du trafic dans la rampe sur la durée des trajets des HT, la présente sous-section aura pour objectif d'estimer le nombre de HT réellement actifs par quart de travail sur un itinéraire donné.

Il est important de bien préciser ici que nous ne tenterons pas d'utiliser le potentiel complet des données historiques disponibles. En effet, il nous serait théoriquement possible de déterminer le nombre de conflits de trajectoires exact ayant opposé un HT donné avec les autres HT de la mine 1 à chaque quart de travail, voire p. ex. à chaque demi-heure. Pour autant, ce nombre de conflits de trajectoires est évidemment très difficilement prévisible en amont des opérations, et nous cherchons donc plutôt à exprimer un élément de contexte connu des planificateurs miniers avant chaque quart de travail, i.e. le nombre de HT supposés se déplacer activement sur un itinéraire donné. On disposera ainsi d'une variable supplémentaire pour améliorer les prédictions de TT de notre modèle prédictif, qui sera aussi accessible aux planificateurs miniers en amont des activités. À des fins de meilleure compréhension de l'amplitude de l'impact des croisements de HT sur leurs TT, nous comptabiliserons tout de même par ailleurs le nombre de détections de chaque HT effectuées par chaque balise de l'itinéraire étudié à chaque demi-heure. Nous comparerons au prochain chapitre les performances de notre modèle de prédiction de TT lorsqu'on lui fournit une telle masse de données, qu'il pourrait potentiellement réussir à exploiter avec succès, par rapport à une simple estimation du nombre de HT se déplaçant activement sur une section de la mine sur un quart de travail au complet.

Pour obtenir cette dernière estimation, on commence par se donner un itinéraire (dans notre cas un itinéraire prédominant), puis on liste les balises qui jalonnent cet itinéraire. On découpe ensuite notre jeu de données par portions temporelles correspondant aux quarts de travail. On s'intéresse alors à chacune de ces portions en comptabilisant sur chacune d'entre elles le nombre de détections respectif de chaque HT sur l'ensemble des balises de cet itinéraire. Notons bien que toutes ces détections ne correspondent pas obligatoirement à un véritable déplacement du HT en raison de la problématique de déplacements erratiques des HT et plus précisément des conflits de périmètres de balises évoqués dans la sous-section 4.5.8; ils ne donnent donc pas une estimation totalement fiable des déplacements de chaque HT, mais

cela n'est pas nécessaire. Il ne reste alors plus qu'à fixer un seuil arbitraire du nombre de détections nécessaires pour qu'un HT puisse vraiment être considéré comme « se déplaçant activement sur cette section ». Pour cela, on compare le nombre de détections totalisées par chaque HT sur différents quarts de travail. On constate généralement un écart important entre les HT les plus actifs (dont le nombre de détections se chiffre souvent en centaines) et les moins actifs (qui comptabilisent souvent moins de 20 détections sur la rampe 1 et moins de 40 sur la rampe 2). À défaut de pouvoir trouver analytiquement un seuil assurément pertinent, on finit par établir un seuil arbitraire de 25 détections pour la rampe 1 et de 50 détections pour la rampe 2 sur un quart de travail donné avant de considérer un HT comme étant actif sur l'itinéraire prédominant étudié. Ces seuils ne sont valables que pour la mine 1, et un travail similaire devra donc être réalisé pour les autres sites miniers.

4.5.11 Évaluation de l'influence de variables d'intérêt sur les TT

La présente sous-section vise à permettre au lecteur de visualiser graphiquement l'impact de quelques variables d'intérêt sur les TT. Elle donne ainsi un aperçu de l'intérêt et de l'influence attendue de certaines variables d'entrée facilement interprétables sur le fonctionnement du modèle prédictif que nous décrirons au chapitre suivant.

Comme évoqué précédemment, nous allons tester dans le même temps différentes approches pour calculer le seuil de filtrage des trajets anormalement longs. Avant de présenter les approches envisagées, précisons que l'on note « $Q1$ » le premier quartile, « $Q3$ » le troisième quartile et « EI » l'étendue interquartile ($EI = Q3 - Q1$) d'une distribution donnée. Voici maintenant les approches retenues pour tâcher de filtrer de manière pertinente les trajets anormalement longs de nos distributions :

- Le critère de la limite supérieure de la moustache de Tukey [36], que l'on peut écrire $Limite = Q3 + (1,5 \times EI)$ et qui est un critère couramment utilisé pour le filtrage de données aberrantes ;
- Le critère classique basé sur le filtrage des 5% des données les plus élevées ; et
- Un seuil arbitraire extrêmement élevé, qui consistera à exclure uniquement les TT de plus d'une demi-heure sur chaque segment de trajet, avant tout dans le but d'effectuer une comparaison avec les autres méthodes de filtrage pour évaluer leur performance ; ce seuil arbitraire aura tout de même le mérite de filtrer les pauses les plus longues des HT, qui trahissent des trajets évidemment non ordinaires.

Nous nous intéresserons à la fois à la médiane et à la moyenne des TT pour avoir deux perspectives simultanées et éviter de fausser nos conclusions. La médiane est en effet habituellement bien plus pertinente que la moyenne pour étudier les distributions fortement

asymétriques comme celles correspondant à nos itinéraires en descente, mais son utilisation concrète en planification est globalement limitée. La moyenne des TT sur chaque itinéraire peut quant à elle permettre d'estimer la productivité horaire de la flotte de HT, elle est abondamment utilisée par les planificateurs comparativement à la médiane.

Par ailleurs, nous inférerons l'intervalle de confiance sur chacune de nos valeurs de moyenne/médiane, à un niveau de confiance de 95% ou de 99,7% selon les cas. Cette inférence se fera au moyen d'une technique de *bootstrap*, i.e. en simulant un millier de rééchantillonnage à partir de l'échantillon de TT dont nous disposons [37].

Enfin, nous représenterons les moyennes/médianes calculées via un diagramme à barres, en faisant figurer chacun des intervalles de confiance inférés. On ne représentera que les barres formées de plus de 30 échantillons (on considère généralement qu'un échantillon devient représentatif statistiquement pour $n \geq 30$, avec n le nombre d'échantillons), et la taille maximale autorisée pour les barres d'erreur sera fixée à 45 secondes pour limiter l'occurrence d'immenses intervalles de confiance, qui ne favorisent aucunement la compréhension.

4.5.11.1 Influence du quart de travail

Pour un itinéraire donné, on calcule tous les TT directs et on les trie dans deux listes en fonction du quart de travail durant lequel a eu lieu le trajet correspondant. On calcule ensuite les grandeurs statistiques évoquées précédemment pour chacune des listes de TT. On procède ainsi pour quelques itinéraires différents, afin de vérifier si nos observations restent valables à l'échelle de la mine 1 au complet. C'est effectivement le cas, et on représentera donc ici seulement des diagrammes à barres obtenus sur l'un des itinéraires prédominants de la mine.

Premièrement, quel que soit le seuil de filtrage sélectionné, nous pouvons affirmer à un niveau de confiance de 99,7% que la médiane des TT de jour est inférieure à la médiane des TT de nuit sur cet itinéraire, et dans le reste de la mine 1. Cette affirmation repose sur la lecture graphique de chacun des diagrammes à barres, comme celui représenté à la figure 4.14. Cette figure représente en particulier le diagramme à barres observé pour la valeur médiane des TT avec le seuil de filtrage basé sur le critère de la limite supérieure de la moustache de Tukey ; les autres diagrammes représentant la médiane sont extrêmement similaires.

En revanche, concernant les TT moyens, le filtrage des TT supérieurs à 30 minutes ne permet pas de conclure à un niveau de confiance de 99,7% que la moyenne des TT de jour est inférieure à celle des TT de nuit comme en témoigne la figure 4.15. L'utilisation d'un niveau de confiance plus faible n'aurait pas changé ce constat étant donné la très grande proximité des deux moyennes calculées. On aboutit pourtant à cette conclusion avec les deux autres méthodes

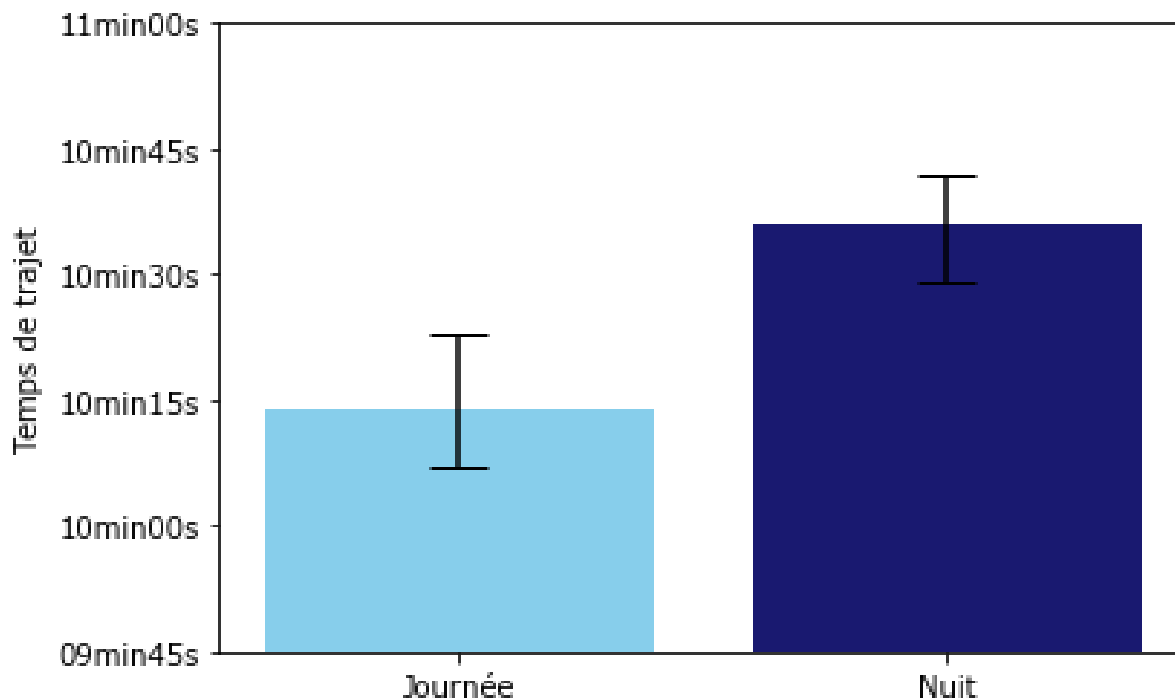


FIGURE 4.14 Diagramme à barres du TT médian selon le quart de travail sur l'itinéraire Surface→Niveau300 de la rampe 2 et barres d'erreur pour un niveau de confiance de 99,7%, seuil de filtrage établi via Tukey

de filtrage à un niveau de confiance de 99,7%. Elles sembleraient donc filtrer efficacement les TT les plus longs puisque ceux-ci sont parfaitement capables de biaiser la moyenne calculée.

Quoi qu'il en soit, l'écart entre la durée des TT de jour et de nuit semble réduit, même lorsqu'on prouve sa significativité statistique. Précisons enfin que les TT par quarts sont dépendants d'autres variables que nous allons examiner, comme le nombre de HT actifs en fonction du quart.

4.5.11.2 Influence du nombre de HT actifs

On s'appuiera bien sûr ici sur notre travail précédent ayant permis d'identifier le nombre de HT considérés comme actifs à chacun des quarts de travail depuis le début de la période temporelle étudiée. Similairement à notre évaluation des quarts de travail, on calcule tous les TT directs pour un itinéraire donné, et on les trie dans le même temps dans l'une des 11 listes correspondant au nombre de HT actifs sur l'itinéraire étudié durant le quart de travail correspondant (il y a au minimum 1 HT actif, i.e. celui ayant permis de détecter le trajet, et la flotte ne dépasse jamais les 11 HT toutes périodes confondues). On calcule les

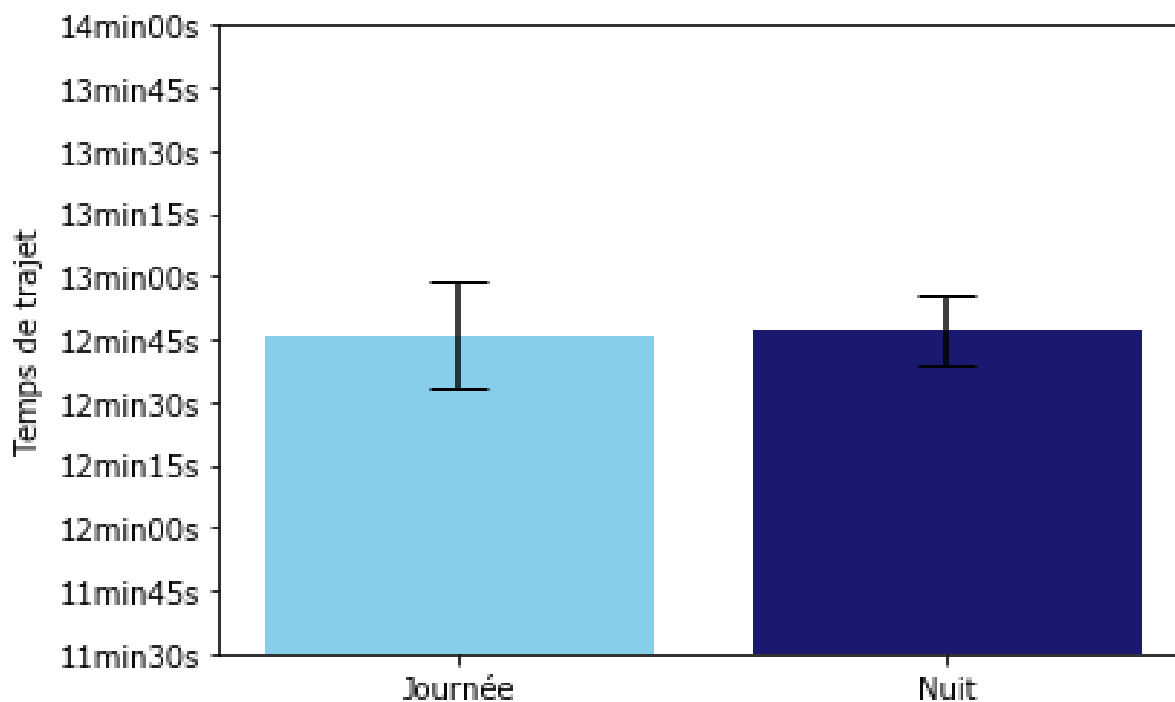


FIGURE 4.15 Diagramme à barres du TT moyen selon le quart de travail sur l’itinéraire Surface→Niveau300 de la rampe 2 et barres d’erreur pour un niveau de confiance de 99,7%, filtrage des TT supérieurs 30 minutes

grandeurs statistiques correspondantes. Nos observations seront différentes en montée et en descente pour toutes les raisons évoquées auparavant, et, étant donné qu’il nous faut satisfaire le critère $n \geq 30$ pour un grand nombre de barres différentes, nous ne nous intéresserons qu’aux trajets sur la rampe 2, qui est bien plus fréquentée sur notre période d’étude. Nous représenterons enfin les diagrammes à barres correspondants avec les différentes méthodes de filtrage. Rappelons que les croisements peuvent avoir lieu entre un HT et un autre véhicule de la mine, mais que nous n’en tenons pas compte et que nous perdons donc nécessairement une partie de l’impact des croisements sur les TT observables.

Concernant la moyenne des TT, notre constat est surprenant puisqu’une seule des trois méthodes de filtrage permet de conclure à la significativité statistique de l’évolution du TT moyen en fonction du nombre de HT considérés comme actifs sur cet itinéraire en descente, à un niveau de confiance de 95% (plus exactement, on peut conclure qu’un grand nombre de HT considérés comme actifs implique un TT moyen significativement supérieur au TT moyen observé lorsque ce même nombre est faible, à un niveau de confiance de 95%). Comme le montre les figures 4.16, 4.17 et 4.18, il s’agit aussi de la seule méthode de filtrage qui permet de clairement visualiser la corrélation positive qui existe pourtant assurément entre ces deux

grandeurs en descente d'après les professionnels du secteur et d'après notre revue de la littérature. La rampe 2 correspond pourtant globalement au cas de figure classique. De surcroît, ladite méthode de filtrage n'est autre que la suppression des TT inter-niveaux supérieurs à 30 minutes, qui devait pourtant servir de modèle de référence pour prouver la pertinence des autres méthodes de filtrage classiques. Autrement dit, les méthodes de filtrage classiques, supposées améliorer ici la significativité des TT en excluant essentiellement des TT particulièrement longs issus de trajets non ordinaires, sont potentiellement contre-productives sur ce point. En effet, elles auraient tendance à supprimer partiellement ces données tout à fait pertinentes pour l'étude de TT de HT en descente. Ce sont effectivement les TT particulièrement longs mais pourtant ordinaires (longues pauses ordinaires supposément dues à de la congestion) qui ont le plus d'effet sur la moyenne des TT en fonction du nombre de HT, et le filtrage semble trop sévère avec les critères classiques. On pourrait chercher à modifier arbitrairement ces critères pour retrouver la significativité désirée tout en filtrant plus de trajets non ordinaires que le critère des 30 minutes, mais il semblerait qu'il soit nécessaire de faire appel à des experts du site minier étudié pour trouver le seuil pertinent à fixer pour chaque segment inter-niveaux. On recommande une telle approche quel que soit le site minier étudié puisqu'il ne semble pas possible de généraliser une méthode satisfaisante. En modifiant les seuils, on risquerait d'ailleurs de perdre au passage une partie de la significativité des quarts de travail sur la durée moyenne des TT avec les deux méthodes de filtrage classiques.

Les diagrammes à barres représentant la médiane des TT en fonction du nombre de HT indiquent logiquement une corrélation moins forte de ces deux grandeurs qu'avec la moyenne. On ne peut pas affirmer à un niveau de confiance de 95% que la médiane des TT augmente avec le nombre de HT actifs, mais c'est le filtrage des TT supérieurs à 30 minutes qui permet là aussi de clairement visualiser la corrélation positive entre ces deux grandeurs.

Enfin, on constate que la corrélation est plus faible en montée qu'en descente pour la moyenne et pour la médiane, sur ce même itinéraire. C'est l'observation attendue. Une corrélation positive est quand même bien visible avec le filtrage des TT supérieurs à 30 minutes, signe qu'il est plus difficile de monter rapidement la rampe lorsque de nombreux HT sont actifs, bien que le HT qui monte ait la priorité.

Rappelons que l'existence de deux rampes dans la mine 1 peut biaiser en partie nos résultats puisque les HT qui descendent via la rampe 2 n'y remontent pas obligatoirement et ne gênent donc pas forcément les autres véhicules en descente. Parfois, un grand nombre de HT actifs sur la rampe 2 pourrait donc seulement signifier que de nombreux HT l'ont empruntée en descente (resp. en montée), et non qu'ils l'ont emprunté massivement pour remonter (resp. redescendre). En étudiant indépendamment chaque quart de travail de la période étudiée, on

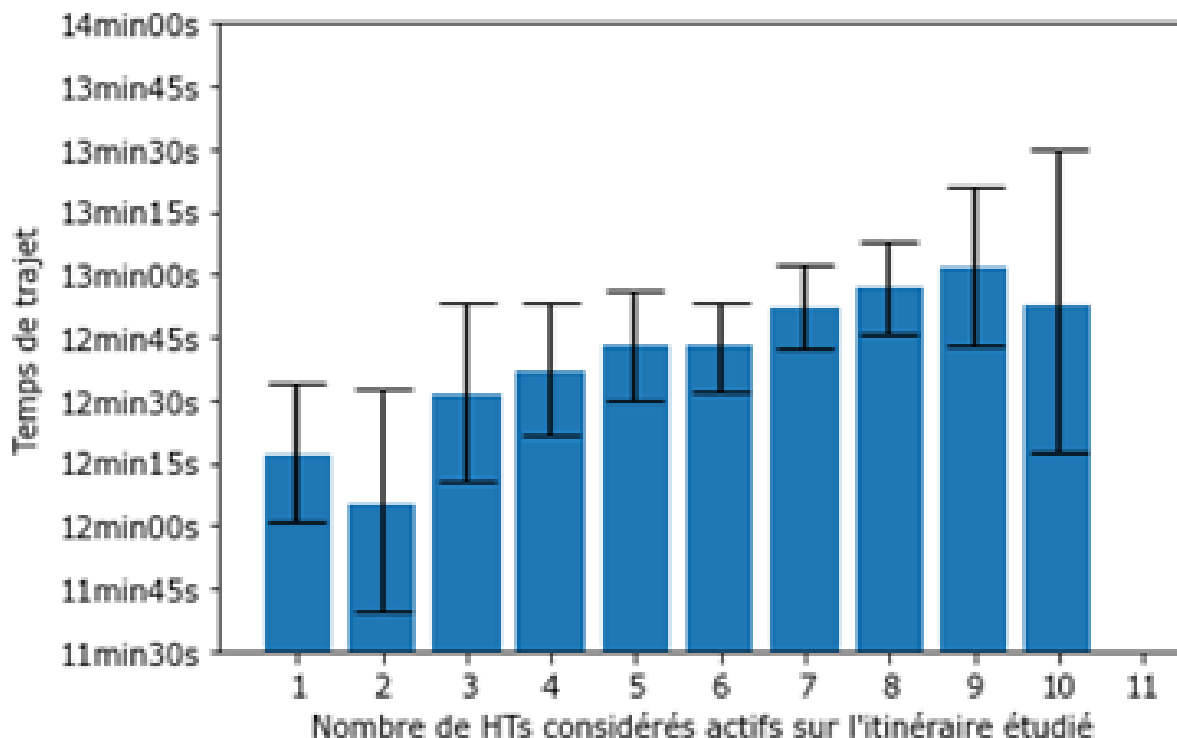


FIGURE 4.16 Diagramme à barres du TT moyen selon le nombre de HT actifs sur l'itinéraire Surface→Niveau300 de la rampe 2 et barres d'erreur pour un niveau de confiance de 95%, filtrage des TT supérieurs 30 minutes

peut tout de même inférer l'existence d'un coefficient de proportionnalité globalement stable entre le nombre de trajets sur la rampe 2 resp. en montée et en descente. Des variations sont néanmoins visibles, ce qui peut là aussi réduire la significativité de la corrélation positive étudiée.

4.5.11.3 Identifiant du HT

On sait que les HT de la mine 1 n'ont pas leur conducteur attribué, mais on ne connaît toutefois pas les modalités exactes selon lesquelles se font le changement de conducteur. On suppose que le choix du HT que va conduire chaque opérateur à chaque quart de travail n'est pas assimilable à un tirage aléatoire, et on suppose donc qu'un identifiant de HT donné réfère plutôt à un nombre limité de conducteurs fréquents différents. Conséquemment, on établit que le TT moyen d'un HT donné peut non seulement dépendre des propriétés du HT en question, mais aussi de l'aptitude, du style de conduite et des habitudes desdits conducteurs fréquents. En plus de garder en tête ce point lors de l'analyse des diagrammes à barres, on décide d'anonymiser l'identifiant de chaque HT pour anonymiser dans le même temps leurs

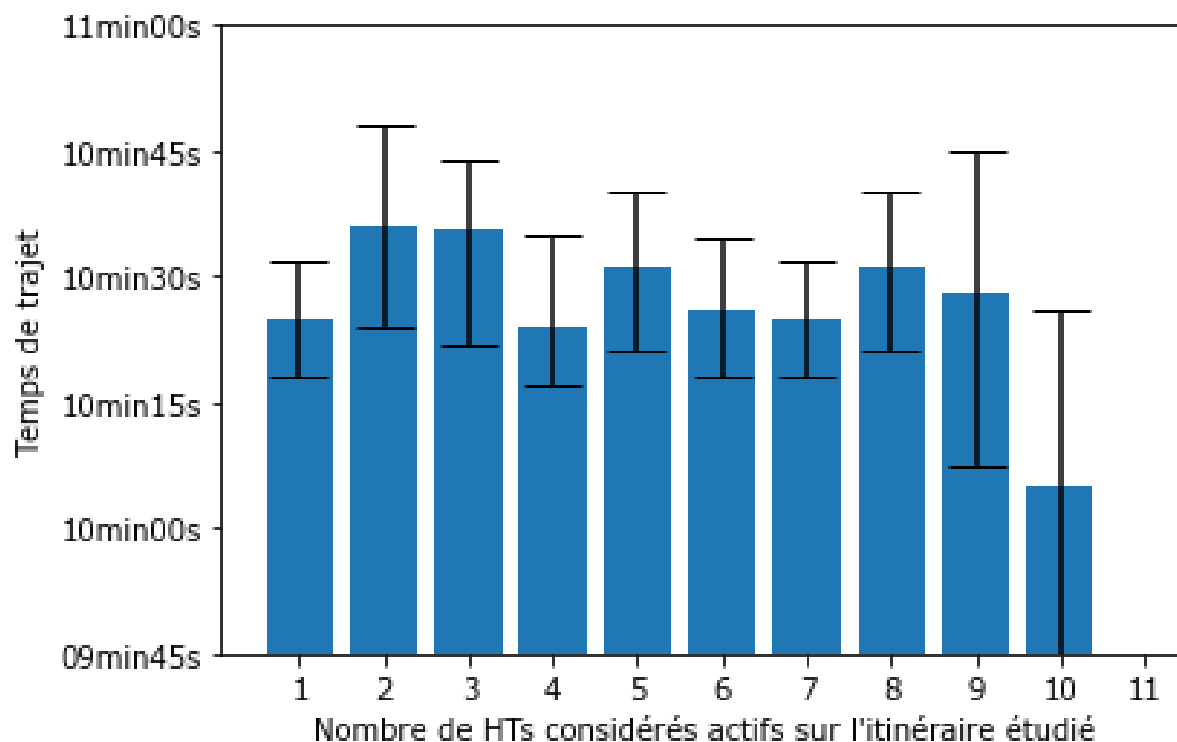


FIGURE 4.17 Diagramme à barres du TT moyen selon le nombre de HT actifs sur l'itinéraire Surface→Niveau300 de la rampe 2 et barres d'erreur pour un niveau de confiance de 95%, seuil de filtrage établi via Tukey

conducteurs les plus fréquents.

À l'image de la démarche suivie pour les autres variables étudiées, on classe dans une liste chaque TT observé en fonction de l'identifiant du HT correspondant, avant de calculer les grandeurs statistiques désirées et de représenter les diagrammes à barres. Là encore, il y a 11 modalités possibles et on décide donc de plutôt s'intéresser à un itinéraire extrêmement emprunté de la rampe 2, en l'occurrence ici l'itinéraire en descente.

On observe que les trois critères de filtrage permettent de conclure que le TT moyen varie significativement en fonction de l'identifiant du HT correspondant, à un niveau de confiance de 95% voire de 99,7% lorsque l'on compare certains HT. En effet, comme on peut l'observer à la figure 4.19, trois HT se démarquent par des TT significativement plus longs que ceux de la plupart des autres HT. On constate que la différence de TT moyen peut être considérablement élevée entre deux HT donnés.

Concernant la médiane, la significativité est un peu plus faible mais les mêmes HT se démarquent encore significativement, à un niveau de confiance de 95%.

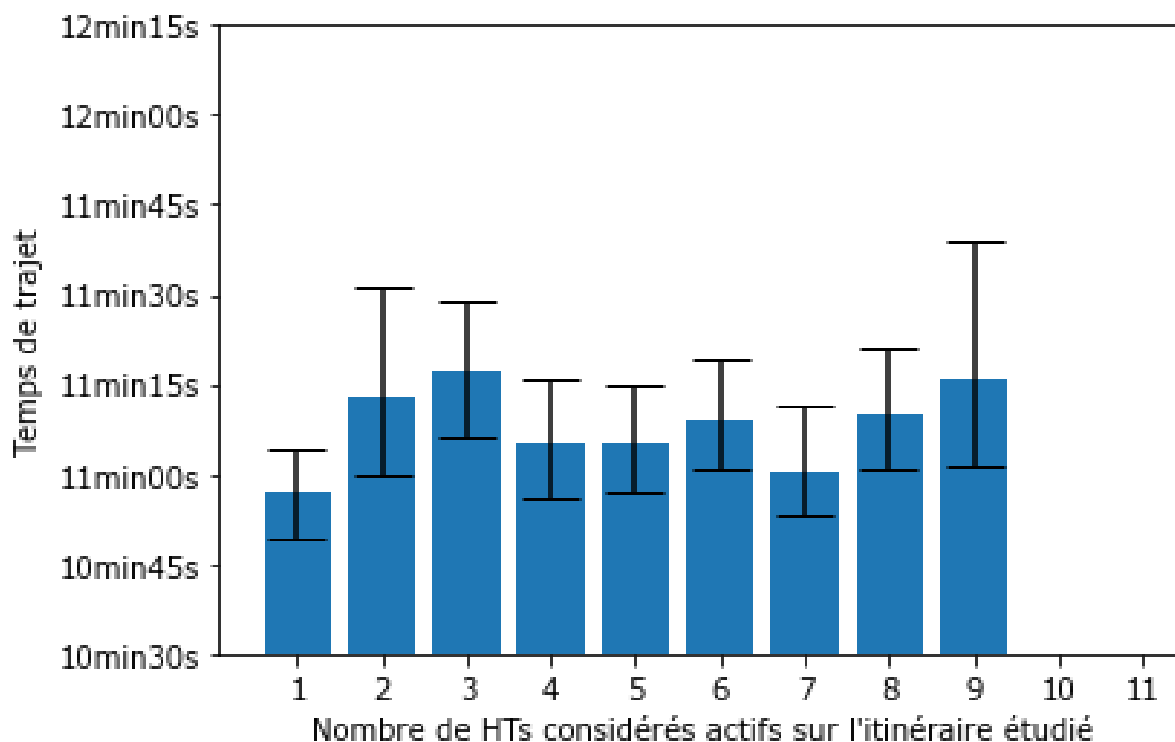


FIGURE 4.18 Diagramme à barres du TT moyen selon le nombre de HT actifs sur l'itinéraire Surface→Niveau300 de la rampe 2 et barres d'erreur pour un niveau de confiance de 95%, filtrage des 5% des TT les plus longs

4.6 Inventaire des pistes d'amélioration relatives aux capteurs, à l'acquisition des données et à leur gestion

À la lumière des observations réalisées au cours de la préparation des données, la présente section vise à inventorier les diverses évolutions relatives aux données que nous conseillons d'envisager dans la mine 1 afin de permettre un jour de disposer de données de meilleures qualités, ou inédites, pour la prédiction de TT de HT. On dresse successivement l'inventaire, dans trois sous-sections dédiées, des pistes d'amélioration relatives aux capteurs de la mine 1, puis celles relatives à l'acquisition des données, pour terminer avec celles qui concernent la gestion des données.

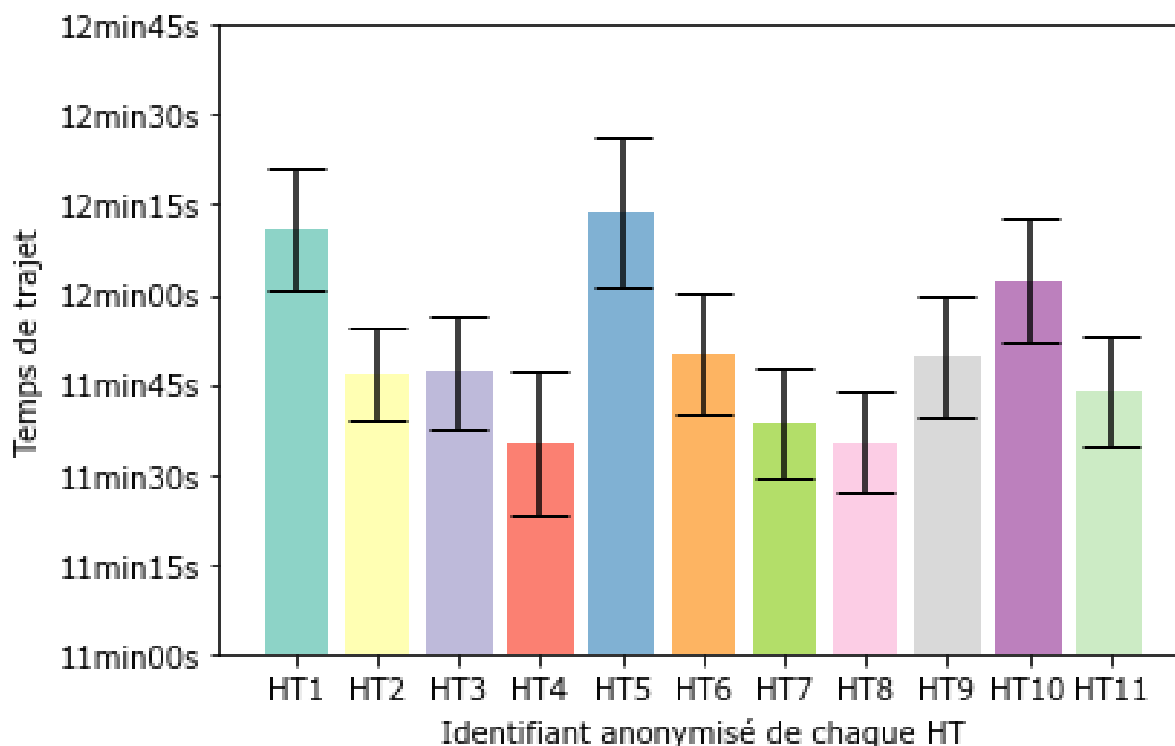


FIGURE 4.19 Diagramme à barres du TT moyen selon l'identifiant du HT correspondant sur l'itinéraire Surface→Niveau300 de la rampe 2 et barres d'erreur pour un niveau de confiance de 95%, filtrage des 5% des TT les plus longs

4.6.1 Pistes d'amélioration relatives aux capteurs

4.6.1.1 Implantation idéale des balises de changement de niveau

Sur ce point, la mine 1 a été plutôt performante mais une amélioration est possible : on remarque, dans des niveaux très profonds (donc récents), qu'on ne trouve parfois pas de balise située précisément à l'intersection entre la rampe et un nouveau niveau. Il serait intéressant d'installer, si cela est possible, une balise aux quelques intersections où ces balises précieuses sont manquantes. Ce sont en effet les balises de changement de niveau idéales, comme nous avons pu l'expliquer. Il serait donc aussi dommageable de ne plus installer de balise à ces endroits précis lors de l'implantation du réseau de balises.

4.6.1.2 Prévention de l'apparition de distributions à deux modes

La compréhension fine du phénomène menant à l'apparition de distributions à deux modes permet d'entrevoir des solutions, à la fois concernant le système de capteurs et l'acquisition

des données, qui permettraient d'éviter définitivement qu'un tel phénomène se reproduise. Ici, on présente simultanément deux solutions consistant à modifier le système de capteurs existant.

Ces solutions se basent sur un principe simple : il faut éviter qu'une balise pouvant être jugée pertinente en tant que balise initiale d'un trajet ne soit précédée d'une zone entièrement aveugle aux balises. En particulier, deux solutions peuvent être imaginées.

- Il serait tout à fait possible d'installer dans le sens du trajet des balises « fusibles » juste avant les balises importantes à même de générer des distributions à deux modes. Ces quelques balises ne seraient pas utilisées pour la prédiction de TT puisqu'elles feraient sinon apparaître des distributions à deux modes. Elles auraient ainsi pour unique but d'empêcher les balises importantes d'être victimes de cette problématique ; et
- Idéalement, il faudrait installer de nouvelles balises pour couvrir toutes les zones qui sont actuellement accessibles aux HT mais hors du périmètre de détection de toutes les balises. On s'assurerait ainsi d'être débarrassé de toutes les situations dans lesquelles une distribution à deux modes pourrait apparaître. Un nombre conséquent de balises pourrait être nécessaire, en particulier en surface, ce qui peut limiter la viabilité de cette solution.

4.6.1.3 Réglage du rayon de détection des balises

On recommande tout simplement de réduire le périmètre de détection des balises souterraines à quelques mètres seulement, en s'assurant au passage que les HT sont quand même détectés à chacun de leurs passages.

4.6.2 Pistes d'amélioration relatives à l'acquisition des données

4.6.2.1 Nettoyage automatisé des identifiants périmés de balises et normalisation du format des identifiants de balises

Il semble essentiel de remplacer dans un premier temps les identifiants de balises périmés apparaissant dans la BDD de la mine 1 par les identifiants actuels de ces mêmes balises. De plus, il serait pertinent de programmer un algorithme capable d'automatiser ce remplacement des identifiants périmés lors de la modification d'un identifiant de balises. Enfin, on propose une codification idéale des identifiants de balises, basée sur les identifiants déjà existants. Pour une balise donnée, cette codification serait composée :

- Du niveau de profondeur approximatif de la balise en question, même si elle se situe en surface (profondeur 000) ;

- D'un mot-clé désignant la zone dans laquelle elle se situe (rampe 1, rampe 2, galerie, surface) ;
- D'un mot-clé spécifiant si cette balise est une balise de changement de niveau (« CN ») ou non (balise régulière « R ») ;
- De tous les autres mots-clés quelconques déjà existants dans l'identifiant de la balise ; et
- De l'identifiant secondaire qui apparaît déjà à la fin des identifiants des balises.

Par ailleurs, pour contourner complètement la problématique qui nous a amené à coder notre algorithme d'identification de profondeur, il suffirait d'ajouter une simple information. En effet, si on ajoute dans la BDD un nouvel attribut qui donnerait le niveau de profondeur associé à chaque balise (y compris pour celles nouvellement installées), il n'y aurait plus besoin d'appliquer cet algorithme.

4.6.2.2 Résolution de dysfonctionnements de balises de changement de niveau

D'après notre lecture de certaines séquences de détection de HT, quelques balises de changement de niveau ne détectent pas toujours tous les HT qui passent en-dessous d'elles, malgré leur très grand périmètre de détection, en particulier durant quelques périodes temporelles restreintes. Il pourrait être intéressant de programmer un algorithme simple vérifiant régulièrement si toutes ces balises sont parfaitement fonctionnelles (p. ex. d'après les détections d'autres balises voisines). En cas de dysfonctionnement, une investigation rapide serait recommandée puisque, s'il manque une seule détection de la part de l'une des balises de changement de niveau, nous avons décidé d'éliminer le trajet correspondant de l'étude avec notre stratégie actuelle de reconnaissance de trajets pour la mine 1.

4.6.2.3 Acquisition de données permettant d'identifier la raison de certaines boucles intra-niveaux

En théorie, si l'une des variables de télémétrie de la BDD de la mine 1 exprimant à intervalles réguliers le niveau de carburant de chaque HT était très régulièrement actualisée, il serait possible d'identifier les détours intra-niveau effectués pour un ravitaillement en carburant. Les TT correspondants pourraient alors être ajoutés aux autres TT déjà retenus pour permettre de prendre en compte cette source de variabilité de TT supplémentaire. Concernant les maintenances réactives (événement non-ordinaire), il faudrait disposer des données de l'intégralité des bons de travail correspondant aux maintenances qui se sont déroulées sur la période étudiée étaient transférées dans la BDD de la mine 1 ; en incluant les renseignements

qui nous intéressent tout particulièrement tels que l'identifiant du HT en maintenance, la date et heure de l'intervention ainsi que sa durée, il serait possible d'éliminer tous les trajets non ordinaires liés à des maintenances réactives.

4.6.2.4 Prévention de l'apparition de distributions à deux modes

Cette problématique handicapante, à laquelle nous avons déjà proposé des solutions qui concernaient le système de capteurs, devrait pouvoir aussi être réglée via une solution d'acquisition des données, qui semble simple et que l'on juge donc réaliste.

En effet, si l'on modifiait la manière dont les données de détection sont enregistrées dans la BDD, il serait possible de s'en débarrasser. La méthode ajustée consisterait à toujours générer un lot de données de détection dans la BDD lorsqu'un HT pénètre dans le périmètre de détection d'une balise (en vérifiant par exemple que cette balise ne détectait pas sa présence durant les 10 secondes précédentes), même si celle-ci est déjà la dernière à l'avoir détecté. Ce réglage permettrait, dans la BDD de la mine 1, d'obtenir deux détections successives de la balise d'accès à la surface A, là où était généré le phénomène indésirable. Ces deux détections seraient distinguables par leur horodatage, ce qui permettrait de savoir, au niveau de la balise A, quand un HT donné s'est dirigé vers la surface et quand il en est revenu. Ainsi, dans le cas d'un trajet $A \rightarrow B$ en descente, notre algorithme de reconnaissance de trajets ne se soucierait que de la seconde détection de la balise A, et ne commencerait à compter la durée du trajet qu'à partir de là. On aurait de plus une excellente manière d'obtenir l'histogramme des temps totaux dédiés aux activités en surface, qui peuvent être trouvés entre deux détections de la balise A, et qui pourraient s'avérer intéressants pour d'autres études de durées d'activités. Enfin, on juge que cette méthode ne devrait générer qu'une faible proportion de lots de données de détection supplémentaires dans la BDD, ce qui ajoute à sa crédibilité.

4.6.3 Piste d'amélioration relative à la gestion des données : élimination des périodes de confusion d'identifiants de véhicules

Bien qu'il soit difficile de déterminer la cause exacte ayant mené des identifiants de véhicules à se confondre temporairement, il n'est pas difficile de mettre en place un protocole permettant d'identifier au jour le jour si des TT Q-I ont eu lieu ou non. Il suffit en fait d'analyser les données de détection quotidiennes d'après un critère similaire à celui que nous avons fixé arbitrairement dans la sous-section 4.5.7. Cela permettrait de s'attaquer à la racine du problème dès qu'un TT Q-I survient, en investiguant les modifications récentes relatives aux véhicules concernés pour déterminer la cause du dysfonctionnement, et éviter de polluer définitivement les périodes de détection ultérieures.

4.7 Prétraitement des données préparées

Nous allons maintenant nous atteler au prétraitement des données issues des étapes de préparation précédentes. Dans la présente section, notre prétraitement sera présenté quasi-intégralement, mais on gagnera en fluidité à expliciter quelques étapes mineures (qui ne concernent que les RNN) au chapitre suivant.

Pour une variable donnée, certaines valeurs peuvent favoriser une baisse du TT (resp. une hausse), on transformera alors généralement ces valeurs pour qu'elles deviennent les nouvelles valeurs les plus négatives (resp. les plus positives) de cette variable, favorisant ainsi une bonne interprétation de la variable correspondante par les réseaux de neurones. Précisons que de prime abord, cette transformation est contre-intuitive puisque l'on a tendance à considérer ces valeurs comme des conditions opérationnelles « positives » (resp. « négatives ») pour les HT et la productivité des opérations de transport de minerai.

Concernant les données opérationnelles, nous listerons ci-après tous les encodages et les normalisations de variables que l'on effectue.

Identifiant du HT : Nous décidons d'encoder l'identifiant du HT par variables indicatrices booléennes via l'encodage « One-Hot », sans tenter de grouper auparavant les HT qui semblaient donner des TT moyens similaires. En effet, nous avons peu de variables explicatives, et nous cherchons ici à aider notre modèle de prédiction de TT en gardant le maximum d'informations disponibles. On supprime au hasard l'une des colonnes résultantes car la combinaison des 10 autres donne déjà la totalité des informations aux modèles.

Nombre de HT actifs sur l'itinéraire étudié : Cette variable est positivement corrélée avec le TT. Nous la normalisons donc pour que le nombre minimal de HT actif (i.e. 1, uniquement le HT en cours d'observation), soit égal à -1 et que le nombre maximal de HT actifs (11 dans le cas de la mine 1) soit égal à 1. Les autres valeurs seront uniformément réparties entre ces deux bornes.

Quart de travail (jour ou nuit) : Cette variable est déjà presque prête puisqu'elle ne peut prendre que deux valeurs, 0 et 1. Pour autant, elle est positivement corrélée aux TT (lorsque l'on passe d'un quart de jour à un quart de nuit). On préfère donc modifier les valeurs associées aux quarts de jour en leur affectant plutôt la constante -1 .

Date d'occurrence du trajet (par rapport à l'intervalle complet de notre jeu de données) : On normalise cette variable sur $[-1; 1]$ puis on calcule l'opposé de chacune de ses valeurs. En effet, on peut supposer que la tendance des TT est légèrement à la baisse au fil des années, sur un itinéraire identique dans des conditions opérationnelles identiques.

Autres variables temporelles : Le **jour de la semaine**, le **jour de l'année** et la **position temporelle relative dans le quart de travail** subiront un encodage cyclique. On décomposera ainsi chacune de ces variables en deux nouvelles variables permettant conjointement de capturer totalement leur aspect cyclique : les composantes sinusoïdale et cosinusoidale, qui seront calculées via les fonctions éponymes. Cas particulier, pour le jour de l'année, on prend d'abord en compte le fait que les conditions hivernales les plus difficiles n'ont habituellement pas lieu le premier janvier de chaque année. Afin que cette variable soit la plus interprétable possible pour les modèles de ML, on implémente ainsi nous-mêmes un décalage temporel adapté à la situation géographique du site minier. Pour l'emplacement de la mine 1, les moyennes saisonnières les plus froides sont atteintes à la fois en janvier et février, on implémente donc un décalage d'un mois en soustrayant 31 aux indices des jours de l'année de chacun des trajets. On constate aussi que les deux mois les plus chauds y sont ex aequo les mois de juillet et août, la cohérence est donc maintenant totale.

Concernant les **TT**, ils devront d'abord être utilisés sous leur forme brute par le modèle de partitionnement inclus dans notre modèle de prédiction de TT. Par la suite, si le trajet est en montée, il suffira de normaliser la distribution de TT brute, car la distribution observée sera significativement similaire à une gaussienne (mis à part les longs TT qui subsisteront). En revanche, en descente, la distribution observée peut être profondément asymétrique. Nous pourrions réduire cette asymétrie et donc l'impact des TT les plus élevés, riches en trajets non ordinaires, en appliquant une transformation logarithmique naturelle à la distribution des TT. On supprime alors les quelques TT qui se détachent clairement de la distribution transformée de par leur faible durée. Dans notre cas, il y en a une vingtaine sur plus d'une dizaine de milliers d'observations. On normalise la distribution obtenue. Naturellement, lorsque nous utiliserons nos modèles de ML, chacune de leur prédiction de TT devra être dénormalisée puis retransformée via la fonction exponentielle. Par souci de concision, nous ne repréciserons globalement pas ces notions dans le prochain chapitre.

À l'issue de cette section, les conditions opérationnelles de chacun des trajets ont été pré-traitées pour être immédiatement utilisées par la quasi-totalité des modèles de ML que nous évaluerons dans le prochain chapitre. Elles forment l'ensemble noté « X ». L'un de nos modèles de ML étant un RNN, quelques transformations supplémentaires de X lui seront spécifiquement nécessaires et nous les décrirons alors. Les TT seront notés « y », qu'ils soient « *log-normalisés* » ou non puisque nous avons déjà précisé que seul le modèle de partitionnement utilisera la forme brute.

4.8 Conclusion

Commençons par souligner le fait que notre méthodologie de préparation des données n'est pas parfaitement généralisable. Nous avons en effet dû régulièrement nous résigner à admettre qu'il ne semble pas possible d'automatiser certaines étapes de la préparation de données. Au lieu de cela, nous avons finalement abouti au constat qu'il fallait fixer ou choisir arbitrairement certaines valeurs, critères et seuils en fonction des connaissances des industriels œuvrant sur un site minier donné. Ainsi, malgré ces ajustements notables nécessaires pour assurer le bon fonctionnement de notre méthodologie, la précision avec laquelle nous détaillons chacune de ses étapes devrait permettre de l'appliquer à toute autre mine souterraine étudiée.

Par ailleurs, la préparation des données issues de la mine 1 s'est révélée remarquablement complexe. En effet, nous avons dû exploiter des détections de balises dont la qualité dépend fortement de leur localisation et de la période temporelle étudiée, et dont l'interprétabilité peut s'avérer bien moindre qu'espérée. Nous avons tout de même pu créer ici une stratégie de reconnaissance de trajets ordinaires qui semble tout à fait fonctionnelle. Combinée à notre prétraitement des données, elle nous permet de bâtir de larges ensembles d'observations de trajets aux propriétés intéressantes. Nous avons en effet pu extraire de nombreuses variables potentiellement pertinentes pour évaluer les conditions opérationnelles futures, purement et uniquement sur la base des détections de balises. Comparativement à ce qui s'est fait dans la littérature, il est par ailleurs remarquable que nous n'ayons extrait aucune variable qui sera inconnue des planificateurs avant chaque quart de travail.

Au travers de ce chapitre, il faut aussi mentionner que le potentiel massif de la technologie de balises de détection de véhicules dans les mines souterraines a clairement transparu. Les défis que nous avons rencontré doivent donc être progressivement relevés pour exploiter son plein potentiel. Nous avons formulé de nombreuses pistes d'amélioration pour la mine 1, et des améliorations impressionnantes de notre préparation de données y sont clairement attendues si l'on entreprend d'appliquer nos prescriptions.

Pour l'heure, si nous ne voulons supprimer aucun trajet ordinaire de notre échantillon, nos données seront entachées d'une quantité substantielle de trajets non ordinaires. Il reste en effet impossible de filtrer ces derniers sans sacrifier au passage une large part des TT ordinaires les plus longs. Ces trajets non ordinaires peuvent défavoriser considérablement le modèle de prédiction de TT que nous allons implémenter au prochain chapitre. Pour autant, nous présumons que la large quantité de données relatives aux détections de trajets ordinaires, que nous avons méticuleusement extraites, et notre prétraitement de données mûrement réfléchi, devraient permettre conjointement d'entraîner convenablement notre modèle de prédiction.

CHAPITRE 5 DÉVELOPPEMENT D’UN MODÈLE INTÉGRÉ DE PRÉDICTION DE TEMPS DE TRAJETS

Ce chapitre présente le développement de notre modèle de prédiction de TT qui reposera sur une combinaison de sous-modèles. Au travers des pages suivantes, nous détaillerons les éléments constitutifs de ce modèle intégré afin de mettre en lumière sa conception et son fonctionnement, puis nous évaluerons ses performances.

Nous commencerons par présenter notre support informatique à la section 5.1 puis nous spécifierons les requis de notre modèle dans la section 5.2. Nous nous intéresserons ensuite à l’architecture adoptée pour notre modèle ainsi qu’aux sous-modèles sélectionnés au cours de la section 5.3. Nous poursuivrons avec la section 5.4, dédiée à la sélection d’un itinéraire pertinent pour l’application de notre modèle et à l’identification d’un seuil de filtrage pertinent des TT. La section 5.5 permettra quant à elle d’expliquer l’approche retenue concernant le partitionnement des données. La section 5.6, centrale, est divisée en deux sous-sections qui viseront chacune à exposer des méthodes visant à prédire les TT : la sous-section 5.6.1 présentera la démarche adoptée pour obtenir plusieurs prédictions initiales d’un même TT, tandis que la sous-section 5.6.2 présentera un modèle d’empilement se basant sur ces dernières pour proposer sa propre prédiction de TT. Nous évaluerons alors à la section 5.7 la qualité des prédictions obtenues avant de formuler nos remarques à la section 5.8. Enfin, nous conclurons.

5.1 Spécifications du support informatique utilisé

Toutes les manipulations du présent chapitre ont été réalisés sur un ordinateur portable dont les spécifications sont résumées dans le tableau 5.1.

5.2 Spécification des requis du modèle

Nous considérons que le modèle développé est destiné à être utilisé par des planificateurs miniers en amont des opérations. On établit donc les requis de notre modèle au regard de ce contexte d’utilisation.

- **Généralisation dynamique** : l’environnement minier souterrain se développant progressivement, notre modèle devra être capable de prédire des TT sur n’importe quel jeu de données obtenu en sortie de notre modèle de préparation de données. Si ce jeu de données est vide (aucun trajet historique complet détecté), notre modèle devra faire

TABLEAU 5.1 Spécifications techniques de l'appareil utilisé dans le présent mémoire

Composant	Spécification
Type d'appareil	Ordinateur portable
Processeur	Intel Core i5-10300H (2.50 GHz, 4 cœurs, 8 processeurs logiques)
Mémoire vive	16 Go
Stockage	SSD de 512 Go
Carte graphique	NVIDIA GeForce RTX 3060
Système d'exploitation	Windows 11

preuve d'adaptabilité pour générer des prédictions de TT grâce aux itinéraires connus. La marge d'erreur de telles prédictions pourra en revanche être considérablement élevée ;

- **Indicateurs de performance** : notre modèle devra retourner à l'utilisateur la RMSE et la MAE associées aux prédictions de TT sur l'itinéraire désigné.
- **Durée d'entraînement de notre modèle sur un itinéraire inconnu** : lors de l'utilisation de notre modèle prédictif pour prédire un TT de HT sur un itinéraire sur lequel il n'a jamais été entraîné, un délai d'entraînement conséquent pourrait être toléré. En effet, certaines problématiques rencontrées dans le chapitre 4 ne permettent pas de préparer des données pour n'importe quel itinéraire. Ainsi, la variété des itinéraires pour lesquels des TT pourraient être prédits sera limitée. L'entraînement de notre modèle sur un itinéraire inconnu restera donc un évènement relativement ponctuel puisque les itinéraires prédominants empruntés par les HT varient lentement, au fur et à mesure de l'expansion du réseau minier. Les planificateurs miniers auraient la possibilité d'anticiper l'entraînement du modèle sur le nouvel itinéraire prédominant qui les intéresse, pour pouvoir par la suite obtenir très rapidement des prédictions de TT sur cet itinéraire, tout au long de son utilisation et quelles que soient les conditions opérationnelles. On considère finalement qu'un temps de compilation supérieur à trois heures serait déraisonnable, puisqu'il ne permettrait pas nécessairement d'ajuster le modèle sur un nouvel itinéraire entre deux quarts de travail ;
- **Durée de prédiction sur les itinéraires connus** : une fois qu'il est entraîné, il paraît raisonnable de demander au modèle de réaliser chacune de ses prédictions de TT en moins d'une seconde pour limiter les temps morts durant son utilisation en planification.

5.3 Architecture du modèle, sélection des sous-modèles et fonctionnement théorique

Notre revue de littérature nous a permis d'identifier les approches et les techniques de prédiction de TT de HT les plus prometteuses. Aussi, nous décidons de retenir certaines d'entre elles et tâchons de les combiner judicieusement pour en tirer des résultats de prédiction plus fiables. Il convient de préciser que, au vu de notre revue de littérature systématique, le modèle intégré proposé dans cette section est inédit pour la prédiction de TT de HT dans un environnement minier.

Les principaux éléments de l'architecture adoptée sont représentés à la figure 5.1.

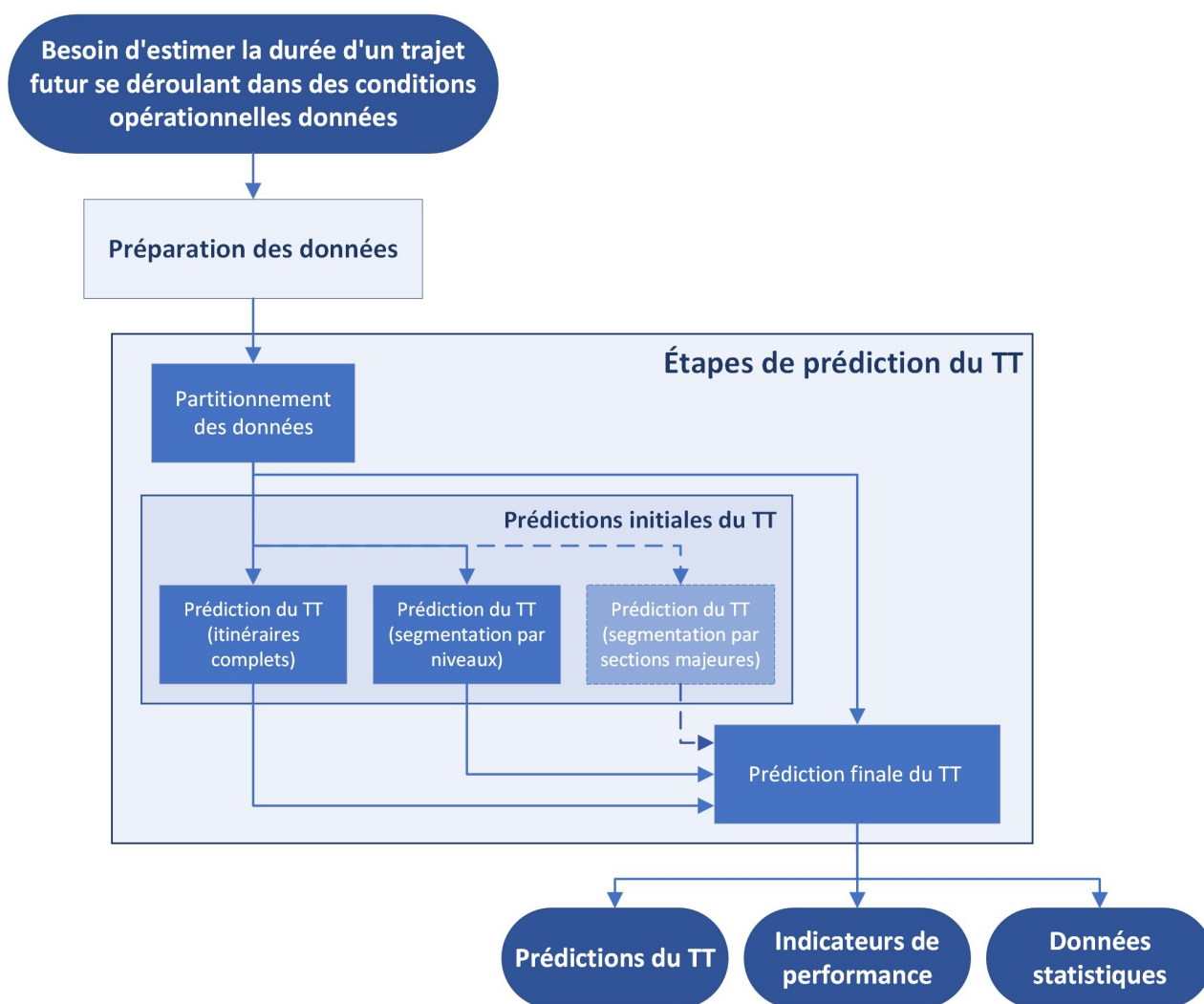


FIGURE 5.1 Organigramme décrivant l'architecture de notre modèle de prédiction de TT au sein de notre méthode globale

Les deux sous-sections suivantes permettront de justifier les éléments constitutifs de cette architecture, de présenter les sous-modèles associés et le fonctionnement attendu à chaque étape.

5.3.1 Sous-modèle de partitionnement

Placé en tête de notre modèle de prédiction de TT, le sous-modèle de partitionnement prend en entrée les données issues de notre modèle de préparation de données. Comme tout modèle de partitionnement, il vise à regrouper ces données d'entrée en n groupes (aussi appelés *clusters*) significatifs. Le choix du nombre de groupes dépend du contexte. Pour une distribution donnée, on peut trouver le nombre optimal de groupes en sélectionnant le meilleur score de partitionnement. Dans le contexte d'une tâche de prédiction, le modèle permettant de minimiser le « Critère d'Information d'Akaike » (AIC) est généralement considéré comme étant le meilleur [38]. Ainsi, ce critère fournira le score de partitionnement qui évalue l'adéquation du GMM avec les données d'entrée, en fonction du nombre de gaussiennes que nous demandons au GMM d'employer pour approcher la distribution observée des TT.

D'après les travaux menés par Fan et al. [18, 21, 22], détaillés dans le chapitre 3, on juge pertinente l'utilisation d'un modèle de partitionnement dans notre contexte opérationnel. En effet, l'équipe de chercheurs avait pu démontrer l'intérêt d'utiliser un tel modèle en amont de divers autres modèles de ML pour prédire la productivité horaire de HT dans des mines à ciel ouvert. Nous considérons que notre problématique actuelle présente des similarités considérables avec cet objectif et nous sélectionnons donc le même modèle que ces chercheurs : le GMM, dont ils avaient justifié l'utilisation. Comme eux, nous désirons que ce modèle sépare pertinemment les variables décrivant nos trajets en quelques groupes représentant chacun un type de trajet différent. Lesdits types de trajets pourraient inclure de potentiels ralentissements et arrêts selon des conditions opérationnelles plus ou moins favorables ; cette classification permettrait donc en aval de plus facilement prédire les TT correspondants.

Si le GMM ne réussit pas à trouver quelques groupes pertinents directement dans les données de prédiction, il est possible de l'associer à un modèle supplémentaire pour obtenir malgré tout un partitionnement intéressant. Dans ce cas, on commencera par simplement appliquer le GMM à la distribution des TT observés sur ce trajet pour identifier les gaussiennes latentes permettant de reconstruire le plus efficacement possible cette distribution de TT, toujours d'après l'AIC. On demandera ensuite au GMM de donner, pour chaque trajet observé, la probabilité d'appartenance du TT correspondant à chacune des gaussiennes identifiées, uniquement d'après la connaissance de ce TT. Ces probabilités serviront ensuite de fonction objectif pour un réseau de neurones de type perceptron multicouches (*Multi-*

Layer Perceptron (MLP)). Il tâchera de prédire la probabilité qu'un trajet donné, avec son lot de conditions opérationnelles, se retrouve dans chacune des gaussiennes prédites. Enfin, ces prédictions de probabilité seront jointes aux variables explicatives décrivant les conditions opérationnelles pour aider les sous-modèles de prédiction qui suivront. Pour évaluer la pertinence de cette combinaison d'un GMM avec un MLP de régression, on comparera les prédictions des sous-modèles à la fois sans, et avec, ces prédictions de probabilité.

Pour implémenter notre GMM en Python, nous utiliserons la classe `GaussianMixture` du module `mixture` de la bibliothèque Python `scikit-learn` (aussi appelée `sklearn`). Concernant l'implémentation du MLP de régression, nous utiliserons les classes du module `keras` de la bibliothèque `tensorflow`. L'optimisation de ses hyperparamètres fera appel à différentes classes de la bibliothèque `keras_tuner`, en particulier `RandomSearch`.

Nous pouvons affirmer que l'utilisation d'un GMM pour améliorer la prédiction de TT de HT en souterrain est inédite d'après notre revue de littérature systématique.

5.3.2 Sous-modèles de prédiction de TT

Présentons maintenant notre démarche globale de prédiction de TT et notre sélection de sous-modèles de prédiction.

5.3.2.1 Segmentations d'itinéraires adoptées et modèles associés

Nous avons pu constater dans la littérature qu'il était possible de prédire un TT sur un itinéraire pris dans son ensemble, mais aussi sur ce même itinéraire segmenté. Dans ce dernier cas, on a pu observer qu'il est d'usage de simplement sommer les prédictions de « TT par segment ». Ces deux méthodes fonctionnelles sont potentiellement complémentaires puisqu'elles ne présentent pas les mêmes avantages. En particulier, dans notre contexte précis :

- **Arguments favorables à étudier les itinéraires au complet** À nos yeux, le premier avantage à prédire des TT sur les itinéraires au complet réside dans le fait que la variance des TT prédits est alors considérablement réduite. Cela s'explique par le fait que ces TT complets sont eux-mêmes des sommes de TT par segment. Ces derniers trajets, plus courts, ont une durée qui varie amplement en fonction des ralentissements et accélérations du HT. Lorsqu'ils s'additionnent naturellement pour former un TT complet, ces variations tendent à se compenser pour faire converger la somme vers la moyenne du TT complet. Pendant la phase de prédiction du modèle, l'utilisation de cette méthode réduit ainsi considérablement la probabilité de produire des prédictions aberrantes, ce qui devrait contribuer à limiter l'erreur moyenne des prédictions. Un

autre phénomène, lié cette fois-ci à l'imprécision de la localisation estimée du HT via les balises, augmente l'intérêt d'utiliser les itinéraires complets pour nos prédictions. Effectivement, lorsque l'on relève un TT par segment d'un itinéraire, nous estimerions que l'incertitude sur la localisation du HT au moment de la détection peut équivaleoir à une variation du TT calculé allant de plusieurs secondes à plusieurs dizaines de secondes. Cette forte incertitude s'explique par le très grand rayon de détection des balises et parfois par la trajectoire potentiellement bien différente du HT par rapport à celle qu'il aurait lors du trajet complet que l'on cherche à étudier. Lorsque nous sommerons ces TT par segment, les incertitudes se cumuleront, ajoutant un bruit certainement handicapant au TT prédit. Pour un trajet complet, cette incertitude n'est observée qu'au tout début du trajet, et elle est ainsi mieux diluée dans le TT total, apportant une robustesse appréciable comparativement à la segmentation des itinéraires.

- **Arguments favorables à la segmentation des itinéraires** : pour autant, comme l'a montré la littérature, la segmentation des itinéraires en vue de produire des prédictions de TT complet peut permettre d'améliorer considérablement les résultats de prédiction finaux. Nous avons déjà segmenté les itinéraires lors de la préparation des données dans le but d'affiner notre compréhension des trajets et notre filtrage des TT les plus longs, car les durées aberrantes sont généralement mises en exergue sur les plus courts segments du trajet complet. Cette segmentation vise ici à aider le modèle à comprendre quelles variables d'entrée sont à même de faire considérablement varier les TT pouvant être observés sur chaque sous-segment. Une sous-segmentation pertinente devrait permettre au modèle de comprendre les phénomènes latents à l'origine de telles variations de TT, en particulier pour les segments où les congestions sont les plus probables (p. ex. aux abords du niveau 300 de la mine 1). Par ailleurs, notons que l'on observera généralement un bien plus grand nombre de trajets sur chacun des segments qui constituent un itinéraire quelconque que de trajets complets sur ce même itinéraire. Cette remarque est valable quelle que soit la finesse de la segmentation. Ainsi, sur de nouveaux itinéraires vers les niveaux les plus profonds, le nombre de trajets complets observables pourrait être bien trop faible pour prédire de manière satisfaisante les TT futurs, et la segmentation des itinéraires serait donc nécessaire à l'obtention de résultats acceptables. De manière générale, on a aussi bien moins d'observations de trajets sur les itinéraires les plus longs. Cela s'explique en partie par le fait que ces trajets sont naturellement plus rares, mais aussi car notre algorithme de préparation des données élimine le trajet complet étudié dès lors qu'un seul de tous les éléments de contrôle programmés n'est pas respecté. Un long trajet, même valide, a donc plus de chances d'être éliminé de l'étude, p. ex. si une balise interdite détecte malencontreusement le

HT à cause de son trop grand périmètre de détection ou, dans une moindre mesure, si une unique balise de changement de niveau manque de le détecter. Rappelons que ce devrait être quelquefois le cas pour les longs trajets et que notre algorithme, tel qu'il est configuré pour la mine 1, ne peut pas accepter cette situation. Dans le cas de la mine 1, nos stratégies de préparation de données et de prédiction s'appuient sur l'intégralité de ces détections particulières pour améliorer leur précision.

Nous supposons qu'en combinant pertinemment ces deux méthodes, leurs inconvénients respectifs pourraient être efficacement réduits : on pourrait théoriquement obtenir des prédictions rarement aberrantes, mais elles pourraient tout de même s'éloigner remarquablement du TT moyen observé sur un itinéraire (particulièrement en descente) lorsque les conditions de trajet paraissent exceptionnellement difficiles au modèle.

Concernant la méthode de segmentation des itinéraires, nous proposons deux logiques de segmentation bien différentes :

1. La première est valable pour n'importe quelle mine souterraine. Elle consiste à identifier tous les segments inter-niveaux inclus dans l'itinéraire. Ces segments inter-niveaux sont délimités par les balises de changement de niveau, ou bien par les balises qui font office de balises de changement de niveau dans les sections de rampe simple ; et
2. La seconde semble pertinente en raison de l'architecture observée dans la mine 1 (voir figure 5.2). Chaque itinéraire reliant la surface à un niveau de grande profondeur pourrait être considéré comme une succession de deux itinéraires de bonne longueur séparés par le niveau 300, qui est bien spécifique et pour lequel le TT du HT sont bien moins prévisibles. Précisons que le réglage inadéquat des balises ne nous permet assurément pas d'identifier tous les détours et pauses des HT au garage, et que la probabilité de rencontre avec d'autres HT est supposément très importante. Dans le cas des itinéraires qui passent par la rampe 2, on doit prendre en compte la traversée du niveau 300 comme un petit itinéraire de transition entre les deux itinéraires de bonne longueur, pour bien prendre en compte ces spécificités. Il est théoriquement tout à fait possible de parcourir ce petit segment sans être ralenti, et l'augmentation du TT liée à la congestion ne devrait donc pas être une constante facile à prédire pour les modèles. Cette deuxième segmentation peut être adaptée à des mines à l'architecture plus classique. Admettons p. ex. que l'on s'intéresse à un itinéraire qui relie l'un des niveaux les plus profonds d'un site minier depuis la surface. Admettons aussi que l'on sait par expérience que le modèle prédit très efficacement le TT mis par un HT pour relier un niveau intermédiaire donné depuis la surface. Dans ces conditions, il serait particulièrement pertinent d'utiliser ces prédictions ici pour favoriser l'obtention de prédictions de TT de haute précision sur

l'ensemble de l'itinéraire. Cette logique de segmentation, bien que facultative, pourrait donc s'avérer redoutablement efficace pour exploiter les connaissances que l'on a déjà concernant les performances du modèle sur des itinéraires intermédiaires.

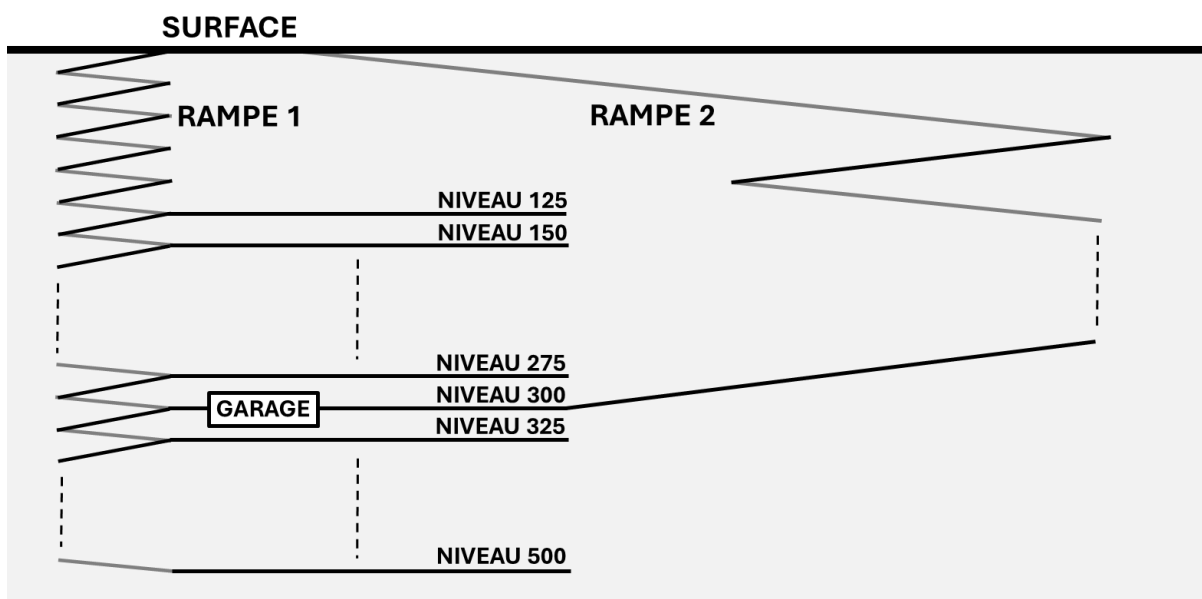


FIGURE 5.2 Architecture fondamentale de la mine 1

Avec ces deux à trois segmentations, on peut donc générer jusqu'à trois prédictions initiales de TT bien différentes. Les deux à trois sous-modèles de prédiction associés exploiteront chacun un degré différent de finesse dans l'analyse des trajets et des conditions opérationnelles associées.

Évoquons maintenant une nouvelle problématique pour garantir que notre méthodologie soit reproductible : nous devons fixer un seuil du nombre d'observations de trajets dont on devrait disposer (sur l'itinéraire au complet et sur les segments inter-niveaux) *a minima* pour garantir une certaine représentativité de l'échantillon. On évitera ainsi que nos sous-modèles ne divergent. Nous fixerons ce seuil à 100, mais c'est un strict minimum. Aussi, il serait parfois bien plus intéressant de s'intéresser à des itinéraires raccourcis avant de chercher à prédire les derniers segments inter-niveaux d'un long itinéraire donné, peu observé. Pour continuer en ce sens, dans le cas de l'étude d'un itinéraire complet inédit (resp. presque inédit), le nombre d'observations de trajets de HT sera nul (resp. insuffisant pour constituer un échantillon représentatif) sur l'itinéraire complet, mais aussi assurément sur son segment inter-niveaux le plus profond. En effet, c'est plus spécifiquement ce petit segment qui manquera d'observations, et non le reste de l'itinéraire. Dans une telle situation, on s'appuiera sur l'arbre de décision suivant :

- Si ce segment le plus profond est très similaire à celui qui le précède immédiatement :
 - Si ce segment précédent présente plus de 100 observations de trajets :
 - On prédit les TT sur ce segment puis on duplique ces prédictions en admettant qu'il s'agit aussi des prédictions de TT sur le segment problématique.
 - Sinon :
 - On retourne au début de l'arbre de décision en s'intéressant cette fois-ci au segment inter-niveaux antérieur.
- Sinon :
 - On divise la longueur du dernier segment par la vitesse moyenne du HT dans la rampe pour obtenir une estimation de la future moyenne des TT sur ce dernier segment.

On pourra procéder de même si une centaine d'observations ne suffit pas et que l'on désire disposer d'un millier d'observations de trajets sur l'itinéraire quasi-complet pour ensuite se ramener à l'itinéraire complet.

Dans tous les cas, on devra adapter un peu notre méthode de prédiction de TT. En effet, nous devrons appliquer nos modèles de prédiction à l'itinéraire quasi-complet, i.e. raccourci d'un segment inter-niveaux (ou de plusieurs si l'on a dû remonter plusieurs fois l'arbre de décision). Si notre itinéraire respecte les critères évoqués ultérieurement, on appliquera alors les méthodes décrites dans les prochains paragraphes aux prédictions obtenues. Enfin, on additionnera simplement les prédictions de TT alors obtenues avec les TT estimés via l'arbre de décision précédent pour le(s) dernier(s) segment(s). Grâce à cette stratégie de contournement, on devrait réussir à prédire convenablement les TT futurs sur des itinéraires presque inédits voire totalement inédits. La portion totalement inédite de ces itinéraires étant très courte comparativement à la longueur totale de ces derniers, on suppose même que la précision des prédictions devrait être notablement similaire à celle pouvant être obtenue sur le même trajet raccourci d'une ou plusieurs balises de changement de niveau.

Revenons maintenant au cas général. Quelle que soit la segmentation adoptée, on commencera par entraîner nos modèles de ML à prédire les TT sur chacun des segments puis on leur demandera de formuler leurs prédictions. Pour améliorer ces dernières, nous commençons par tenter de réduire l'erreur des prédictions sur chaque segment. En effet, la segmentation des itinéraires peut éventuellement induire des biais dans nos modèles au cours de l'entraînement du fait de l'utilisation potentiellement massive de trajets par segment non inclus dans les trajets complets, potentiellement légèrement différents de ceux qui y sont inclus. Si l'on dispose d'au moins une centaine d'observations de trajets **complets** sur l'itinéraire étudié,

on utilise alors une technique de régression des résidus. Sur chaque segment, elle visera à mieux réaligner les prédictions du sous-modèle sur les TT issus des conditions véritables des trajets complets. On testera ultimement l'utilité de cette régression des résidus sur les valeurs de test, pour éventuellement s'en défaire dans l'éventualité où elle ajouterait une complexité inutile au modèle.

Pour autant, la manipulation de ces prédictions de TT par segment ne sera pas terminée puisque l'agrégation des TT par segment n'aura pas encore été réalisée. Bien que la manière la plus simple d'y parvenir soit de sommer directement les TT corrigés par la régression des résidus, on préfère employer une méthode plus robuste pour minimiser le décalage éventuel entre les TT obtenus après agrégation et les TT réels. On a pour cela élaboré une fonction de coût personnalisée pour déterminer une régression linéaire multiple réellement pertinente entre les TT prédits pour chaque segment et la durée totale du trajet à prédire. Cette fonction de coût personnalisée nous permettra de pénaliser les régressions qui se distinguent trop d'une simple agrégation par sommation des TT par segment.

Tout comme la régression des résidus, cette régression linéaire multiple n'aura de sens que si l'on dispose d'un nombre suffisant d'observations de trajets complets sur l'itinéraire étudié. Pour les mêmes justifications, on désire disposer d'une centaine d'observations *a minima*. À défaut, on pourra appliquer ces deux régressions pour tous les segments situés sur le trajet quasi-complet, écourté de quelques balises problématiques situées aux extrémités pour retrouver quelques centaines d'observations. Alternativement, si l'implantation pratique de ces étapes facultatives venait à complexifier considérablement l'utilisation de notre méthodologie par les planificateurs, il serait possible d'agréger simplement par sommation les TT par segment, bien que la précision du modèle complet puisse potentiellement en pâtir.

Finalement, une fois que nous disposons des prédictions initiales de TT sur l'itinéraire complet et sur l'itinéraire segmenté d'une ou plusieurs manières, un sous-modèle supplémentaire sera dédié à effectuer la prédiction finale de TT. Il s'agira d'un modèle d'empilement se basant à la fois sur les prédictions initiales que nous venons d'évoquer et sur les probabilités respectives d'appartenance du trajet aux n groupes générés par le modèle de partitionnement. Nous tâcherons bien sûr de comparer les résultats obtenus par ce sous-modèle supplémentaire comparativement à chacun des trois autres modèles. Nous les comparerons aussi à une prédiction très basique de ce TT complet, i.e la moyenne des TT contenus dans les données d'entraînement. Cette prédiction sera donc identique pour tous les trajets. Elle correspond au modèle de référence répandu dans l'industrie minière selon la littérature. Ajoutons que le modèle d'empilement aura aussi accès à cette valeur en complément pour ses prédictions.

5.3.2.2 Modèles de ML testés

Les sous-modèles que nous destinons à produire des prédictions de TT dans le présent mémoire seront tous des modèles de ML puisqu'ils se sont montrés d'une part très supérieurs à la méthode de simulation TALPAC pour cette tâche, pourtant répandue dans le secteur minier, et d'autre part, puisqu'ils semblent de manière générale être les plus performants pour la prédiction de TT (et grandeurs associées) de HT selon la littérature.

Les modèles de ML sélectionnés pour obtenir des prédictions initiales de TT pour chaque groupe du partitionnement seront les six modèles qui semblaient supplanter tous les autres modèles de référence apparaissant dans notre revue de la littérature. Il ne semblait par ailleurs pas évident de départager ces six modèles puisque leurs performances dépendaient du contexte opérationnel étudié, ou bien parce qu'ils n'ont pas été directement comparés entre eux dans les articles étudiés. Nous avons listé ces modèles apparemment performants dans notre revue critique. On rappelle ci-dessous le nom de ces modèles, ainsi que la classe Python ou les ressources que nous utiliserons pour respectivement implémenter chacun de ces modèles :

- **XGBoost** de régression via la classe `XGBRegressor` de la bibliothèque `xgboost` ;
- **RF** de régression via la classe `RandomForestRegressor` du module `ensemble` de la bibliothèque `scikit-learn` ;
- **GBR** via la classe `GradientBoostingRegressor` du module `ensemble` de la bibliothèque `scikit-learn` ;
- **SVM** de régression via la classe `SVR` du module `svm` de la bibliothèque `sklearn` ;
- **BRNN** via le module `keras` de la bibliothèque `tensorflow` ;
- **LSTM** via le module `keras` de la bibliothèque `tensorflow` (il s'agit du RNN évoqué à plusieurs reprises au chapitre précédent) ;

Nous testerons chacun d'entre eux pour chacun des degrés de finesse de segmentation adoptés.

Concernant le BRNN, nous devons mentionner ici que le modèle que nous implémenterons n'en sera pas un *stricto sensu*. La complexité élevée d'un véritable BRNN lui permet de modéliser l'incertitude sur les poids des connexions neuronales et sur la variable de sortie (ici les TT à prédire). Cela lui permet de formuler non pas des prédictions déterministes, mais de véritables distributions de probabilité du TT attendu à chaque trajet. Or, notre méthodologie ne s'appuie pas sur de telles distributions de probabilité car nous désirons comparer les performances déterministes de six modèles pour les trois segmentations possibles. Aussi, pour des raisons de complexité, on cherchera simplement à incorporer une forme de régularisation, inspirée des réseaux bayésiens régularisés, à un MLP. Cette méthode devrait donner

selon nous des résultats déterministes significativement proches des valeurs centrales des distributions qui auraient été modélisées par un véritable BRNN. Nous pouvons approcher les principales caractéristiques d'un tel réseau de neurones via l'introduction combinée :

- De plusieurs couches cachées successives ;
- De paramètres de régularisation « L2 » au sein de ces couches, qui viseront à contrôler les poids de façon similaire à celle des distributions *a priori* gaussiennes centrées utilisées dans un véritable BRNN ;
- D'un *dropout* (ou « désactivation aléatoire de neurones ») pour ajouter un effet de régularisation en introduisant une incertitude dans les activations des neurones, imitant ainsi approximativement le fonctionnement d'un véritable BRNN.

Nous conservons par ailleurs le terme « BRNN » pour désigner le modèle que nous implémentons, malgré sa nature déterministe.

Pour implémenter le modèle de régression linéaire des résidus, nous utiliserons la classe `LinearRegression` du module `linear_model` de la bibliothèque `scikit-learn`.

L'application du modèle de régression linéaire multiple, qui vise à agréger les prédictions de TT par segment, reposera sur la classe `Minimize` du module `optimize` de la bibliothèque `scipy`, car elle permet l'intégration de notre fonction de coût personnalisée.

Enfin, pour trouver un modèle d'empilement adapté, nous testerons les performances de quatre modèles respectivement issus de quatre des classes les plus simples parmi celles présentées, à savoir `LinearRegression`, `RandomForestRegressor`, `GradientBoostingRegressor` et `SVR`.

Concernant l'optimisation des hyperparamètres de tous les modèles cités dans les paragraphes précédents, la régression linéaire multiple fera figure d'exception et on ne détaillera son optimisation que lors de son utilisation. Pour les autres modèles, voici les ressources ou les techniques utilisées pour mener à bien notre optimisation des hyperparamètres :

- **Pour le BRNN et le LSTM** : on utilisera les classes de la bibliothèque `keras_tuner`, en particulier `RandomSearch` ;
- **Pour le modèle d'empilement** : les classes de la bibliothèque `optuna` [39] nous permettront d'optimiser conjointement notre sélection du modèle de ML (parmi les quatre modèles proposés) à utiliser et la combinaison d'hyperparamètres utilisée par ce dernier ; et
- **Pour tous les autres modèles** : on emploiera la classe `RandomizedSearchCV` du module `model_selection` de la bibliothèque `scikit-learn`.

5.4 Sélection de l’itinéraire de test et identification du seuil de filtrage à privilégier

Au cours de la présente section, nous exposerons les derniers éléments préliminaires qu’il nous reste à fixer avant l’application concrète de notre modèle. On présente les réflexions qui nous permettent de sélectionner avec confiance un itinéraire pertinent pour réaliser nos prédictions de TT, puis celles permettant d’identifier un seuil de filtrage à privilégier.

5.4.1 Sélection de l’itinéraire de test

Dans le reste du présent chapitre, nous nous intéresserons à un seul itinéraire de grande longueur. Nous jugeons pertinent de justifier cette restriction importante dans le présent paragraphe. Tout d’abord, mentionnons le fait que notre méthodologie de préparation des données doit être adaptée à l’itinéraire que l’on décide d’étudier. Même appliquée à un unique site minier, elle ne permet pas, en effet, de disposer immédiatement de données prêtes à l’emploi pour tous les itinéraires du site. De surcroît, l’architecture particulière du réseau de la mine 1 et les réglages problématiques des balises nous obligent parfois à employer des efforts conséquents pour convenablement adapter les critères d’exclusion des trajets problématiques dans notre algorithme. Les itinéraires les plus longs et ceux qui incluent des segments que nous n’avons pas encore explorés complexifient là encore cette préparation des données : il suffit d’une unique balise régulièrement défaillante sélectionnée comme balise de changement de niveau, ou d’une seule balise interdite qui détecte souvent à très longue distance les HT circulant pourtant sur l’itinéraire prévu, pour faire exclure une vaste majorité des trajets observables. Ces contraintes liées à la préparation des données se cumulent bien sûr avec la complexité du modèle de prédiction proposé. Au vu de la difficulté clairement anticipable de la tâche de prédiction, on juge cette complexité pertinente et nécessaire à l’obtention d’un modèle de prédiction robuste. D’autres éléments achèvent de nous convaincre que l’étude d’un unique itinéraire de grande longueur semble déjà constituer une tâche d’une ampleur appréciable. Citons en particulier les multiples modèles de ML que nous allons tester et comparer entre eux après optimisation de leurs hyperparamètres respectifs pour les différentes segmentations, mais aussi le besoin de valider si chacune des nombreuses étapes additionnelles proposées dans notre méthodologie permet effectivement d’améliorer les résultats de prédiction de notre modèle.

Déterminons maintenant l’itinéraire de grande longueur sur lequel appliquer notre modèle intégré de prédiction, de manière à rendre ses performances générales plus facilement évaluable. Simultanément, on désire que toutes les étapes additionnelles de sa structure com-

plexe soient mises à l'épreuve. Pour satisfaire ce dernier point, qui maximise véritablement la complexité de notre méthodologie, on doit s'intéresser à un itinéraire de grande longueur qui traverse le niveau 300, donc qui passe nécessairement par la rampe 2. Ce trajet doit idéalement être connecté à la surface et aller jusqu'à un niveau considérablement profond de la mine 1 ; notons au passage qu'il s'agit d'un itinéraire similaire à ceux qui intéresseraient le plus les planificateurs. Aucune balise de changement de niveau de cet itinéraire ne doit se révéler régulièrement problématique et la balise A de ce trajet ne doit évidemment pas provoquer l'apparition d'une distribution à deux modes. Enfin, on doit disposer sur cet itinéraire d'au moins une centaine d'observations de trajets complets, i.e. des trajets jugés entièrement valides par notre algorithme de préparation des données, comme l'imposent certaines des méthodes facultatives présentées plus haut. Dans le contexte de ce mémoire, il faut bien noter que notre objectif est de valider les performances de notre modèle avec une confiance élevée, nous avons donc en réalité besoin de dépasser très largement les 100 observations de trajets complets. La disponibilité de quelques milliers de trajets complets observés est donc ici un critère déterminant (on aimerait disposer de plus de 1000 trajets de test). Après vérification, l'itinéraire **surface→niveau350 via la rampe 2** vérifie tous ces critères. Ajoutons qu'il s'agit de surcroît d'un itinéraire en descente favorisant les croisements de véhicules, qui sont particulièrement délicats à prédire et qui sont supposés capables d'augmenter très considérablement certains TT ordinaires. Bien que l'on dispose de plus de 5000 observations de trajets jugés valides sur cet itinéraire, qui permettent un meilleur ajustement de notre modèle à plusieurs reprises, la difficulté de prédiction des TT devrait rester notablement élevée. Lorsque l'on règle le seuil de filtrage à 30 minutes inter-niveaux, la variabilité des TT retenus par notre algorithme sur un tel itinéraire est visiblement conséquente comme en témoigne la figure 5.3. La traversée du niveau 300 amène un lot considérable de problématiques additionnelles de préparation de données et de prédiction de TT. Il faut aussi noter que l'on mélange des segments de longueur nécessairement variable (il n'y a que six balises exploitables dans la rampe 2) qui se situent alternativement dans une section de rampe simple présentant très peu de virages, puis dans les galeries horizontales d'un niveau, et enfin dans les sections inter-niveaux hélicoïdales d'une autre rampe. L'étude de cet itinéraire va aussi nous amener à étudier l'itinéraire en rampe simple **surface→niveau300 via la rampe 2** qui constitue le seul segment majeur de cet itinéraire et que nous avons eu l'occasion de largement analyser dans le chapitre 4. Lors de la segmentation par segments majeurs, on décide que la prédiction de TT pour le segment reliant le niveau 300 au niveau 350 sera réalisée ici en prédisant le TT sur chacun deux petits segments inter-niveaux correspondants. On juge en effet trop peu pertinente l'utilisation d'un segment majeur, fait pour être de grande longueur, en raison de la trop faible longueur de cette portion d'itinéraire justement.

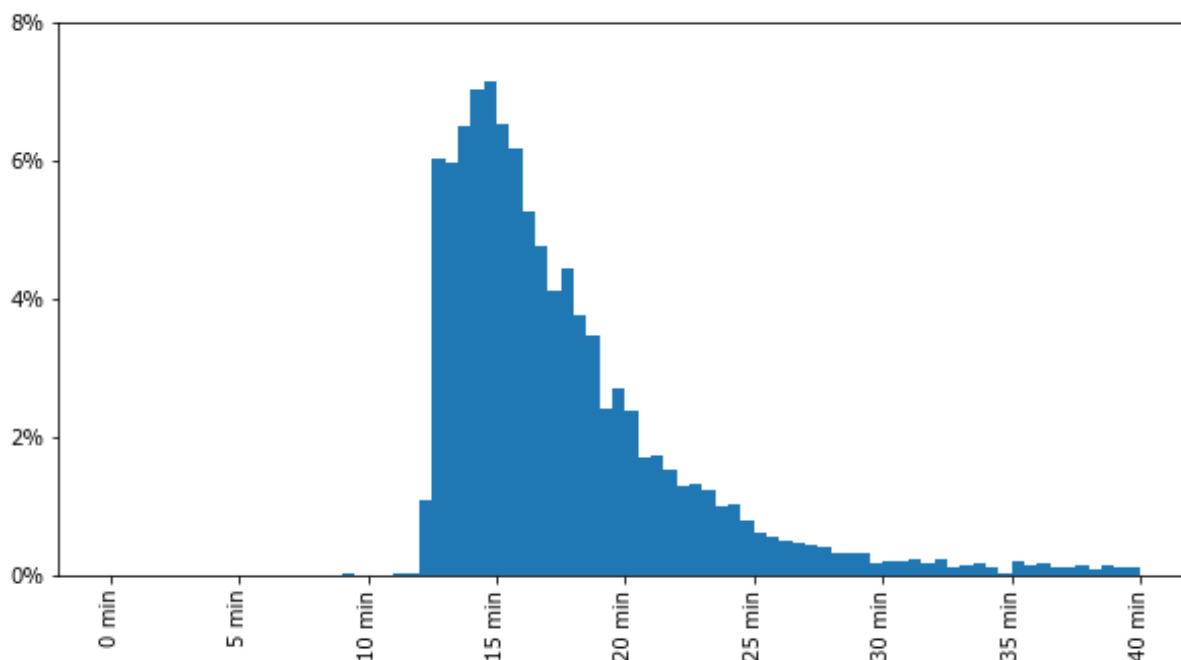


FIGURE 5.3 Histogramme des TT observés sur l'itinéraire surface→niveau350 via la rampe 2 dans la mine 1

5.4.2 Identification du seuil de filtrage de TT à privilégier

Concernant l'identification du seuil de filtrage des TT inter-niveaux à privilégier parmi ceux présentés au chapitre 4, plusieurs discussions doivent être menées pour aboutir à un choix éclairé.

Tout d'abord, retenons que les TT non ordinaires, que l'on cherche à exclure via des seuils de filtrage suffisamment sévères, sont essentiellement provoqués par des pauses non ordinaires dont la durée dépend globalement très peu des conditions opérationnelles du trajet en cours. Les meilleures conditions de trajet possibles peuvent ainsi aboutir à un TT non ordinaire atteignant une durée d'une heure, et non de moins d'un quart d'heure pour l'itinéraire surface→niveau350.

De ce fait, de bien meilleurs résultats de prédiction sont naturellement anticipables si l'on exclut de l'étude tous les trajets dont au moins un TT inter-niveaux était statistiquement considéré comme aberrant (cas d'un seuil fixé via la limite supérieure de la moustache de Tukey) ; ou si l'on exclut une partie non négligeable des TT inter-niveaux les plus longs (cas d'un seuil fixé via la méthode des cinq pour cent). Il faut bien noter que ce dernier seuil n'exclut globalement jamais cinq pour cent des trajets totaux, mais souvent près d'un quart ou d'un tiers des trajets totaux puisqu'il est appliqué à chaque segment inter-niveaux. Il

est d'ailleurs moins sévère que le premier seuil, qui exclut donc une proportion encore plus importante des trajets totaux. Pour ces raisons, ces seuils risquent d'exclure un nombre considérable de trajets qui auraient été jugés ordinaires par un observateur humain présent lors du trajet. Au cours du chapitre 4, rappelons que l'utilisation de chacun de ces deux seuils avait provoqué une perte non négligeable de significativité statistique de l'évolution du TT moyen en fonction du nombre de HT considérés comme actifs sur l'itinéraire surface→niveau300. Il s'agit selon nous d'une raison parfaitement valable de penser que les deux seuils mis ici en question excluent effectivement un nombre considérable de trajets ordinaires sur cet itinéraire. Ce dernier itinéraire est en outre inclus dans l'itinéraire surface→niveau350 sur lequel nous allons appliquer notre modèle de prédiction, cette observation ne doit donc pas être négligée.

Or, bien que le fait de se passer d'un tiers des trajets ordinaires puisse sembler gênant sans paraître inconcevable, il faut bien noter que tous ces trajets ordinaires exclus seraient uniquement les plus longs, avec vraisemblablement les conditions opérationnelles les plus difficiles. Leurs propriétés particulières impliquent que l'on perdrait immédiatement une vaste part de la significativité de nos indicateurs de performance. En effet, dans cette hypothèse, de bons indicateurs de performance signifieraient seulement que notre modèle sait prédire convenablement les TT des trajets les plus simples et ordinaires, dont les TT sont d'ailleurs plutôt stables. Lors de son utilisation en conditions réelles par les planificateurs, notre modèle prédirait alors potentiellement des TT insuffisamment fiables dès lors que les conditions opérationnelles pourraient être jugées plus difficiles pour les HT. La robustesse de ce modèle serait en effet potentiellement insuffisante pour générer des prédictions valides dans ces cas précis puisqu'il n'aurait jamais été entraîné sur les TT liés aux conditions opérationnelles les plus difficiles. Il serait de plus impossible de tester la performance du modèle sur ces trajets précis, puisqu'on ne peut pas les isoler des trajets non ordinaires que l'on avait initialement exclus via, justement, le seuil de filtrage.

En résumé, avec de tels seuils, nous ne pourrions pas dire si notre modèle est valide ou non en bout de ligne. Toutes ces réflexions nous mènent à conclure que nous devons nous résigner à utiliser plutôt le seuil de filtrage des 30 minutes inter-niveaux. Rappelons que si les balises de changement de niveau manquaient véritablement de fiabilité, on ajouterait un autre seuil de filtrage encore plus large (à adapter) pour laisser le temps au HT de parcourir plusieurs segments inter-niveaux successifs sans rencontrer de balises fonctionnelles de changement de niveau.

Bien sûr, le seuil des 30 minutes sera exagérément large comparativement au seuil nécessaire pour filtrer le maximum de TT non ordinaires sans filtrer un seul TT ordinaire ; mais il est

impossible d’optimiser ce seuil ici tant ces deux types de TT se confondent. Avec ce seuil, nous risquons en fait de notablement restreindre les performances de notre modèle à prédire un TT ordinaire, mais pas l’inverse. L’utilisation de la moyenne des TT comme modèle de référence amplifie quant à elle le défi : elle tient compte dans son calcul des longs TT non ordinaires sans être influencée par leurs conditions opérationnelles respectives et se débrouille donc aussi bien quelle que soit la nature, ordinaire ou non, des longs TT. Au contraire, uniquement à cause de ces trajets particuliers, notre modèle pourrait mal évaluer voire mal interpréter l’impact de certaines variables explicatives, attribuant à tort un effet positif à une variable qui, en réalité, a une corrélation négative avec le TT.

Au final, dans le cas où les performances de notre modèle se révéleraient tout de même acceptables, on pourra être certain que son architecture est indubitablement pertinente.

Suite à notre sélection du seuil de filtrage des 30 minutes, notablement handicapant, il nous semble enfin particulièrement intéressant de faire le constat suivant : même en restant entraîné sur nos données de test très insuffisamment filtrées, notre modèle pourrait donner des prédictions de TT significativement meilleures pour la mine 1 lors de son utilisation en conditions réelles. En effet, de notre côté, nous allons devoir évaluer les performances de notre modèle sur la base de données de test décrivant un ensemble de trajets ordinaires et de trajets non ordinaires, sans moyen de dissocier convenablement ces trajets entre eux. Le seuil de 30 minutes favorise l’apparition d’une proportion handicapante de trajets non ordinaires dans ces données. Comme explicité au chapitre précédent, ces derniers ne correspondent pourtant pas à un pur trajet $A \rightarrow B$, mais bien à une combinaison de plusieurs activités (ou inactivités) dont la première fait passer le HT en A et dont la dernière le fait passer en B. Les détections successives des balises nous laissent alors dubitatifs sur la nature du déplacement observé. Pour les planificateurs, les activités qui constituent un trajet non ordinaire pourront régulièrement être prévisibles en amont du quart de travail, en particulier tous types de ravitaillements et certaines pauses. Selon l’organisation adoptée dans un site minier donné, les planificateurs peuvent même être à l’origine d’une vaste part des trajets non ordinaires s’ils programment directement ces activités et communiquent aux conducteurs les horaires correspondants en amont du trajet concerné, afin de réduire les risques de congestion et d’attente aux abords des points de ravitaillement. Par conséquent, les planificateurs n’auraient pas de raison d’assimiler ces trajets non ordinaires prévisibles à de véritables trajets $A \rightarrow B$. Ils n’envisageront vraisemblablement jamais de faire prédire au modèle la durée de ces trajets non ordinaires prévisibles (i.e., de leur point de vue, une combinaison d’activités entière prévue) comme s’il s’agissait d’un pur trajet $A \rightarrow B$. Ils auront plutôt tendance, et nous les y encourageons, à additionner le TT prédit par le modèle pour ce trajet à la durée qu’ils prévoient pour les autres activités incluses. Par ailleurs, si un trajet devient non ordinaire alors qu’il avait

été présupposé ordinaire, il semblerait déraisonnable de considérer que le modèle a échoué à effectuer une prédiction fiable du TT, puisque nous présentons clairement notre solution comme un modèle de prédiction de TT, et non comme un modèle généraliste de prédiction de durées d'activités relatives aux HT. Ce trajet devrait donc être exclu de l'évaluation des performances réelles du modèle, ce que l'on admet notablement plus spontané et plus simple à exécuter lors du suivi des opérations en temps réel qu'*a posteriori*, comme nous avons essayé de le faire. Au final, on considère que, durant son application concrète, notre modèle serait essentiellement évalué sur ses performances à prédire les TT ordinaires comme initialement désiré. Or, contrairement aux conditions opérationnelles des trajets non ordinaires, celles des trajets ordinaires influent bien plus prévisiblement sur le TT à venir, limitant mécaniquement les erreurs de notre modèle. Bien sûr, cette affirmation n'est vraie qu'à la condition que notre modèle ait effectivement acquis un niveau de compréhension convenable des relations entre conditions opérationnelles et TT ordinaires, malgré l'important bruit qui devrait être généré par les trajets non ordinaires. Si ses performances sont prouvées acceptables dans le présent chapitre, on pourra considérer que cette condition est remplie.

En résumé, il nous semble raisonnable d'affirmer que, pour la mine 1, **les performances de notre modèle devraient être dans l'ensemble significativement meilleures en conditions réelles** sous les hypothèses suivantes :

- **H1** : dans la suite du présent chapitre, notre modèle générera des prédictions acceptables sur l'itinéraire surface→niveau350 ;
- **H2** : durant son utilisation réelle, notre modèle sera convenablement entraîné puis appliqué à des itinéraires propices à l'apparition de trajets non ordinaires, ce qui inclut la totalité des itinéraires qui relient la surface à des niveaux profonds de la mine y compris l'itinéraire surface→niveau350 ; et
- **H3** : le processus de filtrage des trajets non ordinaires, décrit ci-avant comme essentiellement spontané, sera effectivement en place dans la mine 1, permettant une évaluation adaptée des performances réelles du modèle.

5.5 Partitionnement des données

La présente section vise à mettre en application notre modèle de partitionnement des données, tandis que sa pertinence pourra être évaluée dans la section suivante.

Rappelons qu'à l'issue du prétraitement de données : « X » désigne la matrice contenant les données normalisées correspondant aux conditions opérationnelles relatives à chaque trajet (données d'entrée) ; tandis que « y » désigne le vecteur des TT correspondants (valeurs ob-

jectifs). Dans le présent mémoire, nous noterons « X_{appr} » les données d’entrée de l’ensemble d’apprentissage et « X_{test} » les données d’entrée de l’ensemble de test. Nous appliquerons le même principe pour y . Par ailleurs, pour alléger nos explications, on utilisera le terme « GMM » à la fois pour désigner le modèle général de ML et pour désigner une instance de la classe **GaussianMixture** avec ses attributs, ses méthodes et ses possibles prédictions. Selon le contexte, on pourra donc écrire à la fois « le GMM est un outil puissant » et « le GMM tente ici d’inférer trois groupes » en désignant deux concepts légèrement différents. On utilisera aussi ce type de raccourci pour les autres modèles de ML.

Intéressons-nous maintenant à l’application concrète du modèle de partitionnement aux données ayant été prétraitées lors du chapitre précédent. Notre modèle devra être adapté puis appliqué à chacun des segments (inter-niveaux ou majeurs) de l’itinéraire, ainsi à l’itinéraire complet. C’est l’itinéraire surface→niveau350 qui nous servira d’exemple pour appuyer la plupart de nos explications.

Pour commencer, on sépare ces données en un ensemble d’apprentissage et un ensemble de test, selon un ratio de 70 : 30. Comme expliqué à la sous-section 5.3.1, on tente dans un premier temps d’ajuster le GMM à X_{appr} , à la recherche de groupes significatifs. L’AIC nous apprend alors très vite que cette stratégie est inappropriée ici : on observe que la valeur de l’AIC diminue de manière quasi-constante avec le nombre de groupes inférés. Cette évolution semble monotone, ou du moins c’est ce que l’on a pu observer jusqu’à cinquante groupes. Par ailleurs, lorsque l’on tente d’inférer un très grand nombre de groupes, lesquels contiennent de moins en moins d’éléments, on soupçonne fortement le GMM de surapprendre toujours plus les caractéristiques de nos données. En bref, aucune solution réellement favorable ne se détache avec cette méthode. La figure 5.4 rend compte intelligiblement de la quasi-absence d’informations supplémentaires apportées par un partitionnement en cinq groupes de X_{appr} via le GMM (ici sur l’itinéraire surface→niveau300) ; la distribution originale des TT y a été colorée en noire, et on a adapté l’échelle verticale de chacune des distributions des cinq groupes (aux couleurs pastels) pour faciliter la comparaison. Les histogrammes observés sont très significativement similaires, ce qui est indésirable, et certains d’entre eux ne contiennent pas plus d’une centaine d’observations tandis qu’un autre contient près de deux tiers des observations totales. Il y a donc un déséquilibre extrêmement fort entre les classes, qui est aussi indésirable. Nos observations, qui tranchent avec celles de Fan et al. [18, 21, 22] lors de leur utilisation du GMM, mériteront une analyse plus approfondie à la section 5.8.

Au vu du caractère inapproprié de cette première stratégie dans notre contexte, on recourt à la seconde, plus complexe.

Pour chaque nombre de groupes à inférer jusqu’à 20, on ajuste le GMM sur y_{appr} en traçant

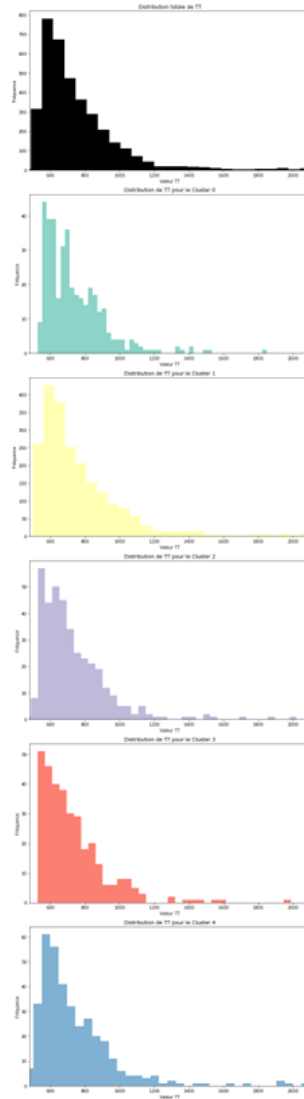


FIGURE 5.4 Illustration des résultats du GMM lors de l'inférence de cinq groupes dans X_{appr} pour l'itinéraire surface→niveau300 via la rampe 2 dans la mine 1

l'évolution de la valeur de l'AIC en fonction du nombre de groupes. D'après cette évolution, tout nous semble cohérent et prometteur, malgré le fait que nous utilisons les paramètres par défaut de **GaussianMixture**. On augmente donc le nombre d'itérations maximal du modèle à 1×10^4 en réduisant dans le même temps la tolérance admise à 1×10^{-9} . On trace à nouveau l'évolution de la valeur de l'AIC en fonction du nombre de groupes, et on fait apparaître le graphe obtenu à la figure 5.5. On retient que le nombre de groupes (et donc de gaussiennes) optimal est de cinq pour l'itinéraire surface→niveau350 via la rampe 2, puisque la pente des segments suivants devient positive, ou négligeable (on tient aussi à maximiser le nombre d'éléments par groupe, donc à minimiser le nombre de groupes).

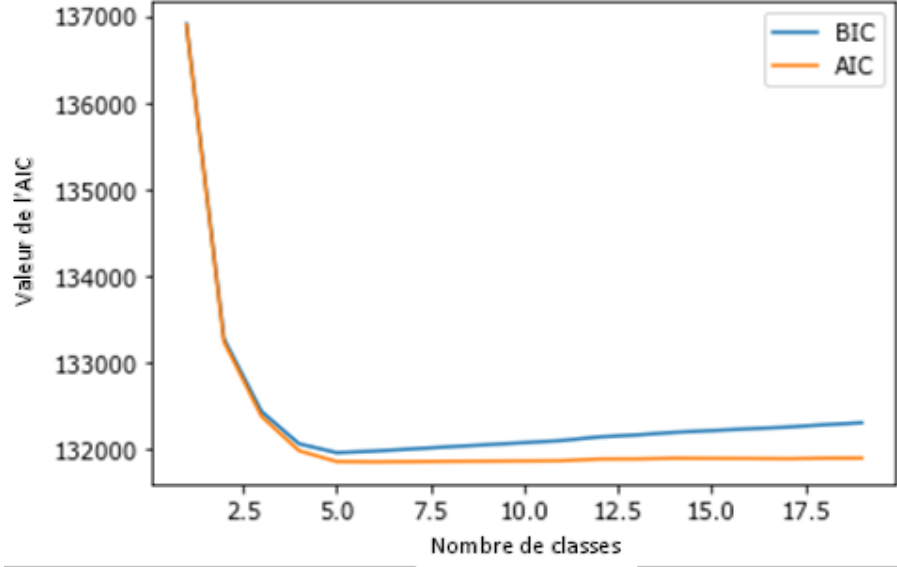


FIGURE 5.5 Évolution de la valeur de l'AIC en fonction du nombre de groupes inférés par le GMM dans y_{appr} pour l'itinéraire surface→niveau350 via la rampe 2 dans la mine 1

On extrait alors le poids, la moyenne et l'écart-type respectifs de chacune des cinq gaussiennes via les attributs `weights_`, `means_` et `covariances_` du GMM. On trie ensuite ces triplets de valeurs par moyenne croissante. Grâce à eux, il est possible de représenter ces cinq gaussiennes sur un seul et même graphe, en utilisant leur densité de probabilité respective (donnée par l'équation $f(x) = \frac{1}{\sigma\sqrt{2\pi}}e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$ avec μ la moyenne de la gaussienne et σ l'écart-type associé) et en pondérant cette dernière par le poids correspondant. On somme ensuite les densités de probabilité de ces cinq gaussiennes pour obtenir la distribution globale qu'elles décrivent, et on trace cette distribution, toujours sur le même graphe. Enfin, pour vérifier la cohérence visuelle de l'ensemble, on incorpore l'histogramme des TT d'apprentissage en arrière-plan du graphe. Comme le montre la figure 5.6, il est indubitable que le GMM a réussi à pertinemment décomposer notre distribution de y_{appr} en cinq groupes. Ceux-ci rassemblent des TT qui sont nettement plus longs lorsque l'on passe d'un groupe au suivant, par moyenne croissante ; on suppose donc que les conditions opérationnelles correspondantes, contenues dans X_{appr} , seront aussi nettement plus difficiles à chaque groupe pour les HT.

On mentionne par ailleurs que, comme expliqué dans les travaux de Fan et al. [18, 21, 22], on préfère habituellement utiliser le GMM sur des distributions pour lesquelles on visualise aisément l'existence d'une combinaison de lois normales symétriques sous-jacentes. Ce n'est pas le cas ici, en raison de l'asymétrie impressionnante de notre distribution de TT. Ainsi, on peut faire observer que l'extrémité gauche de notre distribution, très abrupte, a visiblement compliqué la tâche du GMM : il a été obligé d'utiliser une gaussienne à très faible écart-

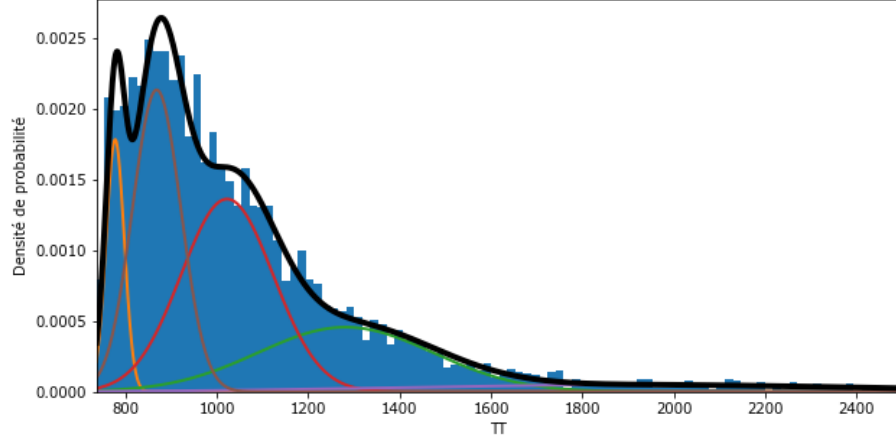


FIGURE 5.6 Inférence de cinq groupes par le GMM dans y_{appr} pour l’itinéraire surface→niveau350 via la rampe 2 dans la mine 1

type uniquement pour imiter cette caractéristique particulière de notre distribution. Si elle avait été un peu plus symétrique, on peut supposer que trois ou quatre groupes auraient éventuellement pu suffire à la décrire de manière optimale, au regard de l’AIC. Malgré tout, nous ne jugeons pas préoccupante l’utilisation d’un total de cinq groupes ici, et gageons que chacun de ces groupes contribuera effectivement à la significativité du partitionnement.

Passons maintenant à la deuxième étape qui, rappelons-le, consiste à entraîner un MLP de régression à prédire la probabilité d’appartenance d’un trajet donné à chacun des groupes qu’a déterminé auparavant le GMM.

Pour ce faire, on commence par appeler la méthode `predict_proba` de notre GMM. Elle permet, pour un TT de y_{appr} donné, de demander au GMM la probabilité que ce TT appartienne à chacun des groupes selon lui. Il s’appuie pour cela uniquement sur la densité de probabilité correspondant à chacune des gaussiennes qu’il a inférées. Nous appellerons $y_{GMM,appr}$ l’ensemble des probabilités ainsi générées par le GMM.

Pour continuer, nous passons à la construction de notre MLP de régression. Il ne sera pas précisé à chaque appel de classe ou de sous-module que, comme évoqué à la sous-section 5.3.1, toutes les classes et tous les sous-modules utilisés pour l’implémentation du MLP proviennent du module `keras` de la bibliothèque `tensorflow`. De même pour notre modèle d’optimisation du MLP, qui se base ici sur la bibliothèque `keras_tuner`.

On commence par créer un modèle séquentiel via la classe `Sequential`. C’est le cadre de notre MLP. On y intègre alors :

- Une couche d’entrée destinée à accueillir X_{appr} puis X_{test} , via la classe `Input` ;

- Deux couches cachées consécutives utilisant la fonction d'activation « *ReLU* », via la classe `Dense` ; et
- Une couche de sortie à n neurones destinée à renvoyer les probabilités d'appartenance de chacun des trajets à chacune des n classes et utilisant la fonction d'activation sigmoïde pour modéliser les probabilités de sortie entre 0 et 1, via la classe `Dense`.

Le modèle utilise l'écart quadratique moyen (non encore introduit jusqu'ici, on le notera « MSE » (*Mean Square Error*)) comme fonction de perte et il est compilé avec la classe `Adam` du sous-module `optimizers` en tant qu'optimiseur.

On entreprend alors d'optimiser la combinaison d'hyperparamètres que notre MLP intègre, via `RandomSearch`. On veut optimiser simultanément : le nombre de neurones de chacune des couches cachées, notés *nbUnits1MLP* et *nbUnits2MLP* ; le nombre d'échantillons par lot (que l'on impose au modèle de traiter avant que les poids du réseau ne soient mis à jour), noté *tailleLotsMLP* ; et le taux d'apprentissage placé en argument de l'optimiseur `Adam`, noté *tauxApprMLP*.

Avant toute chose, on crée une fonction personnalisée `ConstructeurOptiMLP`. Elle rassemble tous les paramètres qui cadrent l'optimisation et elle permet de simplifier le processus d'optimisation à venir, qui est complexifié par notre volonté d'optimiser la taille des lots. Son fonctionnement intègre deux étapes principales :

1. D'une part, elle construit notre MLP d'après tous les éléments évoqués jusqu'ici en configurant les valeurs qui pourront être prises par les hyperparamètres à optimiser qu'elle contient ; et
2. D'autre part, elle configure les valeurs qui pourront être prises *tailleLotsMLP* durant l'optimisation.

Pour ce faire, elle utilise deux classes de la bibliothèque `keras_tuner` :

- La classe `Choice` permettra de sélectionner aléatoirement une valeur du *tauxApprMLP* parmi un ensemble de valeurs discrètes, à chaque itération de `RandomSearch` ; et
- La classe `Int`, utilisée pour les trois autres variables, permettra quant à elle de sélectionner aléatoirement à chaque itération une valeur entière dans un intervalle donné, délimité par des bornes que nous noterons a et b , pour un espacement donné entre chaque valeur de l'intervalle, noté k ; i.e. l'ensemble de sélection sera $\{a, a + k, a + 2k, \dots, b\}$ avec a, b, k des entiers.

Ceci étant dit, nous décrivons les ensembles de sélection des hyperparamètres que nous avons retenu dans le tableau 5.2, ainsi que la classe que nous appliquerons à chacun d'eux dans notre fonction personnalisée.

TABLEAU 5.2 Ensembles de sélection des hyperparamètres du MLP et classe associée

Classe	Hyperparamètre	Ensemble de Sélection
Choice	<i>tauxApprMLP</i>	$\{1 \times 10^{-2}, 1 \times 10^{-3}, 1 \times 10^{-4}\}$
Int	<i>nbUnits1MLP</i>	$\{a = 32, b = 256, k = 32\}$
	<i>nbUnits2MLP</i>	$\{a = 8, b = 128, k = 8\}$
	<i>tailleLotsMLP</i>	$\{a = 16, b = 128, k = 16\}$

Il nous faut maintenant encapsuler `ConstructeurOptiMLP` dans un « modèle-enveloppe ». Cela permettra de faire simultanément varier, pendant l'optimisation, *tailleLotsMLP* ainsi que les hyperparamètres du MLP. Pour ce faire, on implémente une classe personnalisée mais rudimentaire que nous nommerons « **Capsule** ». Le code Python correspondant est présenté à la figure 5.7. **Capsule** hérite de la classe `HyperModel` de la bibliothèque `keras_tuner` (« `kt` » dans le code) et utilisera les hyperparamètres (« `hp` » dans le code) que l'on fournira pour construire le MLP (« `model` » dans le code) et retourner celui-ci ainsi que `tailleLots`. Elle utilisera pour cela notre fonction `ConstructeurOptiMLP` (« `model_func` » dans le code).

```
class Capsule(kt.HyperModel):
    def __init__(self, model_func):
        self.model_func = model_func

    def build(self, hp):
        model, tailleLots = self.model_func(hp)
        self.tailleLots = tailleLots
        return model

    def fit(self, hp, model, *args, **kwargs):
        kwargs['batch_size'] = self.tailleLots
        return model.fit(*args, **kwargs)
```

FIGURE 5.7 Code Python de la classe personnalisée **Capsule**

Finalement, nous disposons de tous les éléments nécessaires pour utiliser `RandomSearch`. On lui donne pour arguments :

- Un hypermodèle, qui est notre **Capsule** appliquée à `ConstructeurOptiMLP` ;
- Un critère de perte, qui sera la perte de validation (appelée via la chaîne de caractères prédéfinie « `val_loss` ») et qui utilisera le MSE puisque c'est ce que nous avons configuré dans `ConstructeurOptiMLP` ;
- Un nombre maximal d'essais, que l'on fixe à 100 car on cherche à considérablement explorer l'univers des combinaisons possibles ; et

- Un nombre d'exécutions par essai, que l'on fixe à 2 pour limiter un peu la variance des résultats tout en évitant d'allonger inutilement cette optimisation, que l'on ne considère pas comme la plus critique de notre méthodologie.

Pour démarrer la recherche des meilleurs hyperparamètres, on applique la méthode `search` de `RandomSearch` à X_{appr} et $y_{GMM,appr}$. On fixe le nombre d'époques maximal à 100. La validation croisée se fait quant à elle en conservant, à chaque itération, 20% des données pour la validation.

Une fois que le processus de recherche est terminé, on extrait les meilleurs hyperparamètres obtenus via la méthode `get_best_hyperparameters` de `RandomSearch`. Les valeurs optimales retournées sont décrites dans le tableau 5.3.

TABLEAU 5.3 Combinaison optimale des hyperparamètres du MLP

Hyperparamètre	Valeur optimale
$nbUnits1MLP$	192
$nbUnits2MLP$	16
$tailleLotsMLP$	48
$tauxApprMLP$	1×10^{-3}

On utilise ces valeurs optimales pour entraîner notre MLP sur X_{appr} et $y_{GMM,appr}$ via la méthode `fit` du MLP. On lui demande ensuite de formuler ses prédictions de classes d'après les X_{test} via la méthode `predict` du MLP, elles seront notées « $y_{MLP,pred}$ ». Nous lui faisons aussi générer, de même, les prédictions notées « $y_{MLP,appr}$ » d'après X_{train} , pour disposer de y_{MLP} au complet.

Tâchons maintenant d'évaluer les performances du MLP sur les données de test. Faute de critère totalement adapté à la comparaison de listes ordonnées de n probabilités complémentaires, nous décidons de simplifier ici les résultats obtenus. On veut repérer la probabilité maximale dans chaque liste de probabilités dans $y_{MLP,pred}$ puis extraire l'indice de la position qu'elle occupe dans la liste ordonnée, sous la forme d'un entier. La méthode `argmax` de la bibliothèque `numpy` permet de réaliser cette manipulation en une seule étape, tout en rassemblant chacun des indices dans une nouvelle liste. Pour chaque trajet de test, celle-ci donne ainsi la classe d'appartenance la plus probable d'après les conditions opérationnelles attendues, selon le MLP. Pour pouvoir la comparer à la liste de référence équivalente pour les données réelles, nous n'avons pas d'autre choix que d'appliquer directement un GMM sur les données de test, à titre de comparaison uniquement. Les résultats de partitionnement de ce GMM ne seront naturellement pas utilisés pour entraîner d'autres modèles. On utilise alors

la méthode `predict` de ce GMM pour directement obtenir la liste des indices correspondants à la classe d'appartenance la plus probable de chaque trajet, selon ce GMM.

On représente enfin le diagramme à barres combinant la liste finale issue de ce GMM (vraies valeurs) ainsi que celle issue du MLP (prédictions) à la figure 5.8.

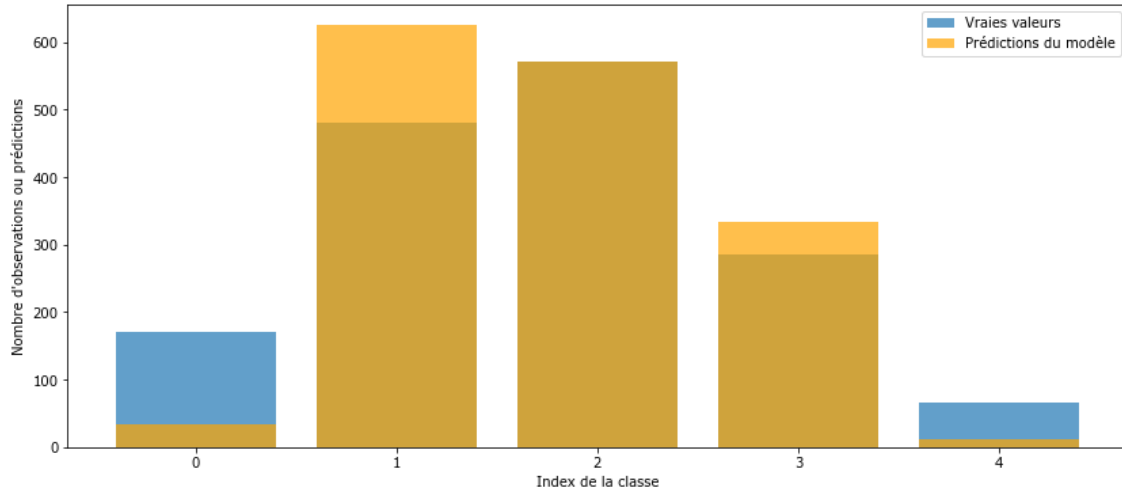


FIGURE 5.8 Diagramme à barres des vrais groupes annoncés par le GMM et des groupes prédits par le MLP pour l'itinéraire surface→niveau350 via la rampe 2 dans la mine 1

Le MLP semble donner des prédictions de probabilité maximale convenablement réparties entre les trois classes centrales, mais il est très rare qu'il tente de donner la probabilité maximale à l'une des classes extrêmes. Celles-ci ont par ailleurs nettement moins d'observations dans les valeurs vraies que les trois classes centrales, ce qui peut justifier en partie cette tendance. Étant donné que nous avons optimisé les hyperparamètres de notre modèle pour maximiser ses performances, on suppose que la distribution obtenue n'est pas trompeuse, dans le sens où notre MLP prédit sûrement bien des probabilités plus élevées pour les grandes classes lorsque les conditions du trajet sont défavorables ; on admet donc que la forte ressemblance entre les deux distributions, du moins sur les trois classes centrales, n'est pas liée au hasard. En bref, sans être infaillible, notre MLP ne prend pas de risques inconsidérés lors de ses prédictions et sait quand même faire varier de manière appréciable les probabilités qu'il prédit. On est donc plutôt confiant sur la qualité des résultats obtenus par le MLP et sur leur intérêt pour la suite de notre modèle.

Afin de permettre aux modèles de prédiction de TT d'utiliser les résultats intéressants de notre MLP pour appuyer leurs prédictions, une toute dernière manipulation doit clore cette section : la méthode `hstack` de la bibliothèque `numpy` nous permet de concaténer la matrice X avec la matrice y_{MLP} en une nouvelle matrice, que nous désignerons simplement par « X_{GMM} ».

On conserve aussi X pour comparer les performances des modèles de prédiction initiale de TT lorsqu'ils profitent ou non des résultats de notre modèle de partitionnement.

Résumons les grandes étapes de cette section :

- Nous avons configuré un GMM pour qu'il partitionne de manière optimale la distribution des TT d'apprentissage ;
- Sur la base de ce partitionnement de référence, nous avons implémenté un MLP pour qu'il prédise le partitionnement des trajets de test d'après leurs conditions opérationnelles, lesquelles seront aussi connues d'avance par les planificateurs ; et
- Nous avons optimisé la combinaison de quatre hyperparamètres cruciaux de ce MLP via une démarche personnalisée.

Au terme de ce travail, nous pensons disposer de données d'entrée bonifiées pour favoriser l'obtention de résultats probants par nos modèles de prédiction de TT.

5.6 Prédiction de TT

La présente section vise à appliquer les méthodes décrites à la sous-section 5.3.2 pour obtenir une prédiction robuste de chaque TT à venir.

La sous-section 5.6.1 présente en détail la démarche nécessaire à l'obtention de plusieurs prédictions initiales d'un même TT pour différentes segmentations, tandis que la sous-section 5.6.2 présente le sous-modèle d'empilement se basant sur ces dernières pour proposer une prédiction finale du TT à venir.

5.6.1 Prédictions initiales de TT sur un itinéraire donné

On présente successivement la démarche détaillée de prédiction de TT sur un itinéraire pris dans son ensemble, sur cet itinéraire sectionné en segments inter-niveaux successives et sur cet itinéraire sectionné en segments majeures.

Chacune de ces démarches pourra utiliser alternativement X et X_{GMM} mais on désignera toutes les données d'entrée par « X » dans la présente section, par souci d'uniformité et de concision.

5.6.1.1 Prédiction de TT sur l'itinéraire complet

La prédiction de TT sur l'itinéraire complet est bien moins complexe que la prédiction de TT sur un itinéraire segmenté, comme cela a pu transparaître à la sous-section 5.3.2. Pour

autant, comme nous traitons cette démarche en première, nous allons devoir expliciter ici les particularités de notre BRNN et de notre LSTM ainsi que l'optimisation des hyperparamètres de tous les modèles.

Pour commencer, on s'intéresse aux quatre modèles les plus simples de notre sélection, i.e. XGBoost, RF, GBR et SVM. Étant donné que X et y *log*-normalisé sont déjà prêts à être utilisés par ces modèles, on entreprend immédiatement d'optimiser les hyperparamètres de ces modèles. On décrit les ensembles de sélection des hyperparamètres de chacune des classes dans le tableau 5.4, avec $\mathcal{U}(a, b)$ l'ensemble des réels entre les valeurs a et b et $\mathcal{U}\{c, d\}$ l'ensemble des entiers entre les valeurs c et d .

TABLEAU 5.4 Ensembles de sélection des hyperparamètres pour XGBoost, RF, GBR et SVM

Modèle	Hyperparamètre	Ensemble de sélection
XGBRegressor	learning_rate	$\mathcal{U}(1 \times 10^{-3}, 1 \times 10^{-2})$
	max_depth	$\mathcal{U}\{2, 10\}$
	n_estimators	$\mathcal{U}\{100, 500\}$
	subsample	$\mathcal{U}(0,5, 0,7)$
RandomForestRegressor	max_depth	$\{\text{None}, 10, 20, 30, 40, 50\}$
	max_features	$\{\text{'auto'}, \text{'sqrt'}\}$
	min_samples_leaf	$\mathcal{U}\{1, 10\}$
	min_samples_split	$\mathcal{U}\{2, 15\}$
	n_estimators	$\mathcal{U}\{100, 200\}$
GradientBoostingRegressor	learning_rate	$\mathcal{U}(1 \times 10^{-3}, 5 \times 10^{-2})$
	max_depth	$\mathcal{U}\{3, 10\}$
	n_estimators	$\mathcal{U}\{100, 300\}$
	max_features	$\{\text{'auto'}, \text{'sqrt'}\}$
	min_samples_leaf	$\mathcal{U}\{1, 10\}$
	min_samples_split	$\mathcal{U}\{2, 15\}$
SVR	C	$\mathcal{U}(1, 100)$
	epsilon	$\mathcal{U}(1 \times 10^{-2}, 1 \times 10^{-1})$
	kernel	$\{\text{'linear'}, \text{'poly'}, \text{'rbf'}, \text{'sigmoid'}\}$

La documentation Python de chacune de ces classes étant très complète et en raison de la très grande variété d'hyperparamètres à optimiser, nous n'expliquerons pas chacun d'entre

eux.

Pour trouver la combinaison d'hyperparamètres optimale de chacun de ces modèles, on donne pour arguments à `RandomizedSearchCV` :

- La classe du modèle ;
- Les ensembles de sélection de chacun de ses hyperparamètres ;
- Un nombre d'itérations, fixé à 100 pour bien explorer les combinaisons mais limiter le temps total de calcul de notre modèle de prédiction ;
- Un nombre de plis pour la validation croisée, fixé à 2 pour réduire la variance et limiter le temps total de calcul de notre modèle de prédiction ; et
- Le nombre de cœurs de notre ordinateur que l'on autorise `RandomizedSearchCV` à exploiter, que l'on règle au maximum ($n_jobs = -1$).

L'optimisation respective de chacun des modèles sur X_{train} et y_{train} aboutit aux hyperparamètres donnés dans le tableau 5.5.

On ajuste alors chacun des modèles sur X_{train} et y_{train} avec leurs hyperparamètres optimaux respectifs via la méthode `fit`.

Pour terminer, on utilise la méthode `predict` de chacun de nos quatre modèles sur X_{test} pour obtenir leurs prédictions respectives de TT, que nous nommerons respectivement de $y_{XGB,pred}$, $y_{RF,pred}$, $y_{GBR,pred}$, $y_{SVM,pred}$. Afin de préparer toutes les données nécessaires pour notre modèle d'empilement, on leur fait aussi formuler leurs prédictions de y_{appr} d'après X_{appr} .

Intéressons-nous maintenant au BRNN. Rappelons que toutes les classes et tous les sous-modules utilisés pour l'implémentation de ce modèle proviennent du module `keras` de la bibliothèque `tensorflow`. Pour le construire, on commence par créer un modèle séquentiel via la classe `Sequential` puis on y intègre :

- Une couche d'entrée destinée à accueillir X_{appr} puis X_{test} , via la classe `Input` ;
- Une première couche cachée utilisant la fonction d'activation `ReLU` et intégrant une régularisation L2 des poids, via la classe `Dense` ;
- Une couche de réduction aléatoire du surajustement à un taux de 50%, via la classe `Dropout` du sous-module `layers` ;
- Une seconde couche cachée, utilisant également la fonction d'activation `ReLU` et intégrant une régularisation L2 des poids, via la classe `Dense` ; et
- Une couche de sortie à un seul neurone utilisant une fonction d'activation linéaire pour prédire la valeur continue du TT.

TABLEAU 5.5 Combinaisons optimales d'hyperparamètres de XGBoost, RF, GBR et SVM

Modèle	Hyperparamètre	Valeur optimale
XGBRegressor	learning_rate	5×10^{-3}
	max_depth	9
	n_estimators	400
	subsample	0,6
RandomForestRegressor	max_depth	None
	max_features	'sqrt'
	min_samples_leaf	6
	min_samples_split	10
	n_estimators	150
GradientBoostingRegressor	learning_rate	1×10^{-2}
	max_depth	6
	n_estimators	256
	max_features	'sqrt'
	min_samples_leaf	6
	min_samples_split	10
SVR	C	20
	epsilon	5×10^{-2}
	kernel	'rbf'

Notre BRNN utilise le MSE comme fonction de perte et la classe **Adam** comme optimiseur.

Comme nous avons pu l'observer, notre BRNN personnalisé a une structure proche du MLP que nous avons décrit précédent. Par conséquent, l'optimisation de ses hyperparamètres sera elle aussi très similaire. Aussi, nous allons pouvoir nous concentrer ici sur les points essentiels.

Voici les hyperparamètres sélectionnés pour l'optimisation :

- Le nombre de neurones de chacune des deux couches cachées, notés *nbUnits1BRNN* et *nbUnits2BRNN* ;
- Les coefficients de régularisation L2 de chacune de ces couches, notés *L2_1BRNN* et *L2_2BRNN* ;
- Le nombre d'échantillons par lot, noté *tailleLotsBRNN* ; et
- Le taux d'apprentissage placé en argument de l'optimiseur **Adam**, noté *tauxApprBRNN*.

Nous créons maintenant une fonction personnalisée `ConstructeurOptiBRNN`, dont le fonctionnement est quasiment identique à `ConstructeurOptiMLP` : on remplace simplement la structure du MLP par celle de notre BRNN et on rajoute deux instances supplémentaires de la classe `Choice` pour sélectionner aléatoirement, à chaque itération de `RandomSearch`, une valeur de `L2_1BRNN` et une valeur de `L2_2BRNN` parmi un ensemble de valeurs discrètes. Ces deux hyperparamètres ne prendront donc pas forcément des valeurs identiques.

Les ensembles de sélection des hyperparamètres que nous avons retenus seront décrits dans le tableau 5.6, ainsi que la classe (`Choice` ou `Int`) que nous appliquerons à chacun d’eux.

TABLEAU 5.6 Ensembles de sélection des hyperparamètres du BRNN et classe associée

Classe	Hyperparamètre	Ensemble de Sélection
Choice	<code>L2_1BRNN</code>	$\{1 \times 10^{-1}, 3 \times 10^{-2}, 1 \times 10^{-2}, 3 \times 10^{-3}, 1 \times 10^{-3}\}$
	<code>L2_2BRNN</code>	$\{1 \times 10^{-1}, 3 \times 10^{-2}, 1 \times 10^{-2}, 3 \times 10^{-3}, 1 \times 10^{-3}\}$
	<code>tauxApprBRNN</code>	$\{1 \times 10^{-2}, 1 \times 10^{-3}, 1 \times 10^{-4}\}$
Int	<code>tailleLotsBRNN</code>	$\{a = 16, b = 256, k = 16\}$
	<code>nbUnits1BRNN</code>	$\{a = 32, b = 512, k = 32\}$
	<code>nbUnits2BRNN</code>	$\{a = 4, b = 128, k = 4\}$

On encapsule enfin `ConstructeurOptiBRNN` dans `Capsule`, notre modèle-enveloppe déjà créé.

Nous disposons alors déjà de tous les éléments nécessaires pour utiliser `RandomSearch`, on l’applique donc aux arguments suivants :

- Un hypermodèle, qui est notre `Capsule` appliquée à `ConstructeurOptiBRNN` ;
- Un critère de perte, qui sera `val_loss` et qui utilisera donc le MSE ;
- Un nombre maximal d’essais, que l’on fixe à 200 car on cherche à explorer suffisamment l’univers des combinaisons possibles, particulièrement vaste ici ; et
- Un nombre d’exécutions par essai, que l’on limite à 2 pour éviter que cette optimisation ne dépasse une durée d’une heure.

On applique la méthode `search` de `RandomSearch` à `Xappr` et `yappr` en fixant le nombre d’époques maximal à 100 et en utilisant 20% de ces données pour chaque étape de la validation croisée.

Les valeurs optimales, retournées par la méthode `get_best_hyperparameters`, sont données dans le tableau 5.7.

TABLEAU 5.7 Combinaison optimale des hyperparamètres de notre BRNN

Hyperparamètre	Valeur optimale
$nbUnits1BRNN$	256
$L2_1BRNN$	3×10^{-3}
$nbUnits2BRNN$	20
$L2_2BRNN$	3×10^{-3}
$tailleLotsBRNN$	128
$tauxApprBRNN$	1×10^{-3}

On utilise cette combinaison optimale pour entraîner notre BRNN sur X_{appr} et y_{appr} via la méthode `fit`, en fixant là encore le nombre d'époques maximal à 100 et en utilisant 20% de ces données pour sa validation croisée. On lui demande de formuler ses prédictions de TT, $y_{BRNN,pred}$, d'après X_{test} via la méthode `predict`. On lui fera aussi prédire y_{appr} d'après X_{appr} , ce qui donnera $y_{BRNN,pred,appr}$.

Intéressons-nous enfin au LSTM. Rappelons que toutes les classes et tous les sous-modules utilisés pour l'implémentation de ce modèle proviennent du module `keras` de la bibliothèque `tensorflow`.

L'architecture particulière de ce modèle paraît prometteuse dans notre contexte, puisqu'il pourrait exploiter les évidentes relations qui lient un trajet futur aux trajets récents observés sur le même itinéraire. Pour rester en cohérence avec un cas d'utilisation réel de notre modèle par les planificateurs, on considérera ici pour un trajet donné que les « trajets récents observés » évoqués sont les trajets qui précèdent immédiatement le quart de travail du trajet auquel on s'intéresse. Il s'agira donc *a minima* de trajets qui se sont déroulés au quart de travail précédent, et non dans les minutes précédentes. Ajoutons que nous disposons évidemment aussi des trajets antérieurs qui se sont déroulés durant le même quart et qu'il serait tout à fait intéressant de comparer la précision des résultats du LSTM lorsqu'ils sont utilisés ou non. On ne peut en effet pas exclure la possibilité que notre modèle soit utilisé pour prédire aussi à très court terme les TT à venir d'un itinéraire donné, p. ex. pour ajuster en temps quasi-réel l'ordonnancement des opérations d'un quart de travail.

Commençons par établir un nombre de « derniers trajets observés » à exploiter par notre LSTM dès lors qu'il veut prédire le trajet suivant. Dans le cas d'itinéraires peu fréquentés, le choix d'un nombre trop élevé de ces trajets risquerait d'amener le modèle à s'appuyer

sur une part considérable de trajets trop anciens pour refléter correctement les conditions opérationnelles du trajet en cours, ce qui pourrait dégrader ses prédictions. Au contraire, le choix d'un nombre trop faible risquerait de pousser notre modèle à essayer de donner trop d'importance à chacun des trajets récents fournis, ce qui dégraderait inévitablement sa robustesse aux TT aberrants. Après quelques tests, nous fixons finalement ce nombre à 10 car il paraît particulièrement laborieux de l'optimiser en plus des autres hyperparamètres, bien que le LSTM y gagnerait sûrement en précision (précisons tout de même que fixer ce nombre à 30 dégrade visiblement les prédictions pour l'itinéraire complet).

L'architecture de tout LSTM va nous obliger ici à créer au nôtre un ensemble de données d'entrée sur mesure, X_{LSTM} et, plus anecdotiquement, à écourter y pour obtenir y_{LSTM} .

Commençons par nous occuper de y_{LSTM} . Étant donné que chacun des trajets qui s'y trouvera devra être précédé de 10 autres trajets ayant eu lieu dans les quarts de travail précédents, on adapte la liste en conséquence : à partir du quart de travail le plus ancien d'une copie de y , on retire tous les trajets de chaque quart de travail jusqu'à cumuler au moins 10 trajets. Les trajets restants forment y_{LSTM} . Sa longueur sera notée « $nbTrajLSTM$ », il s'agit du nombre de trajets à prédire avec le LSTM.

Passons maintenant à la création de X_{LSTM} à partir de X et de y , dont nous utiliserons des copies ici. On pose « $nbCaract$ » le nombre de colonnes de X , i.e. le nombre de caractéristiques sur lesquelles tous les autres modèles prédisent les valeurs de y .

Dans un premier temps, on crée une matrice 3D de dimensions $nbTrajLSTM \times (10 + 1) \times (nbCaract + 2)$ pour l'instant vide, il s'agit de X_{LSTM} .

On concatène y à la fin de X pour former $X_{augmentée}$.

On introduit une nouvelle caractéristique, qui sera spécifique au LSTM, car elle n'est pas exploitable par nos autres modèles : il s'agit du nombre de quarts de travail d'écart entre celui du trajet en cours et celui d'un trajet antérieur. Sa valeur minimale sera donc de 1 (resp. de 0 dans le cas particulier où l'on accepte aussi l'utilisation des trajets du quart de travail en cours). Pour la normaliser, on calcule à proprement parler son inverse i.e. $f(x) = 1/x$ (resp. on calcule $f(x) = 1/(x + 1)$). Ainsi, cette valeur chutant très rapidement aux abords de 1, elle devrait permettre d'inciter le LSTM à pondérer plus généreusement les relations entre les trajets les plus récents et le trajet à prédire. Pour chacun des trajets de y_{LSTM} , on utilise y pour retrouver les 10 trajets qui le précèdent immédiatement sans avoir eu lieu durant le même quart de travail (resp. ayant eu lieu ou non durant le même quart de travail). On calcule le nombre de quarts de travail d'écart normalisé de chacun de ces trajets au trajet de y_{LSTM} et on stocke chacune de ces valeurs dans un vecteur de longueur 10. En itérant

sur chacun des trajets de y_{LSTM} , on obtient un total de $nbTrajLSTM$ vecteurs que l'on place alors dans les 10 premières colonnes de X_{LSTM} , leur faisant occuper la dernière position de chacun des canaux, le long de toutes les lignes.

De même, pour chaque trajet de y_{LSTM} , on crée nos séquences agrégées d'après les 10 trajets précédents dans $X_{augmentée}$: chaque caractéristique de chacun des 10 trajets précédents sera placée dans un vecteur de longueur 10, et chacun de ces vecteurs sera placé dans la ligne correspondante des 10 premières colonnes de X_{LSTM} , occupant ainsi toutes les lignes et les $nbCaract + 1$ premières positions des canaux.

Concernant la dernière colonne de X_{LSTM} , elle permettra au LSTM d'aussi prendre en compte les conditions opérationnelles du trajet dont on veut prédire le TT. Les $nbCaract$ premières positions des canaux seront remplis le long de toutes les lignes en y insérant toutes les valeurs des $nbTrajLSTM$ dernières lignes de la matrice 2D X . Cela signifie que les trajets qui avaient été retirés à y pour obtenir y_{LSTM} ne seront pas transférés dans X_{LSTM} . On appliquera un *padding* constant aux deux dernières positions de chacun des canaux de cette colonne via la méthode `pad` de la bibliothèque `numpy`. La valeur attribuée sera 1 pour l'avant-dernière position, laquelle correspond au nombre normalisé de quarts de travail d'écart. 1 correspond alors aux trajets les plus récents et les plus influents, on incite donc le LSTM à donner un poids élevé aux valeurs de cette colonne. Pour la dernière position, qui correspond au TT *log*-normalisé, la valeur attribuée sera de 0, qui correspondra alors théoriquement à un trajet moyen pour le LSTM, ce qui évite de le déstabiliser. Ceci nous évite aussi de créer un masque pour l'entraînement, qui n'aurait plus fait effet lors de la phase de prédiction. Enfin, on surpondère cette colonne, évidemment plus significative que les autres, en la multipliant simplement par un coefficient k , supérieur à 1. Nous décidons d'optimiser manuellement cette valeur, une fois que nos hyperparamètres auront été déterminés. On admet que cette valeur donnera toujours un bon aperçu du ratio entre l'importance des conditions opérationnelles en cours par rapport aux caractéristiques des dix trajets passés, quel que soit l'itinéraire et le site minier étudié. On la fixera donc définitivement par la suite pour toutes les prédictions utilisant un LSTM.

On sépare enfin X_{LSTM} et y_{LSTM} en un ensemble d'apprentissage et un ensemble de test, selon un ratio de 70 : 30.

Passons maintenant à l'implémentation du LSTM. Nous le construisons à l'aide de l'API fonctionnelle de `keras`, car nous désirons implémenter un mécanisme d'attention. Nous avons donc besoin d'une plus grande flexibilité dans la définition des architectures de réseau que celle permise par la classe `Sequential`. Voici les éléments successifs composant la structure de notre LSTM :

- Une couche d'entrée destinée à prendre en entrée X_{LSTM} , via la classe **Input** ;
- Une couche de neurones LSTM configurée pour retourner les séquences entières et qui intégrera une régularisation L2, via la classe **LSTM** ;
- Un mécanisme d'attention, qui retourne des poids d'attention et que nous décrirons plus tard ;
- Un calcul de produit scalaire des sorties LSTM et des poids d'attention, via la classe `texttt{dot}` ;
- Une fonction d'aplatissement du résultat obtenu, via la classe **Flatten** ; et
- Une couche de sortie utilisant une régularisation L2, via la classe **Dense**.

Le mécanisme d'attention évoqué plus tôt a pour unique objectif d'augmenter la pondération des pas de temps les plus pertinents de la séquence temporelle. Ici, on suppose que les TT les plus récents ont une plus grande influence sur le TT à venir, on estime donc que les tous derniers trajets parmi les 10 trajets de notre fenêtre temporelle devraient voir leur pondération augmenter. Le mécanisme d'attention repose sur la succession d'éléments et de fonctions suivante :

- Une couche de neurones avec activation *tanh* qui calcule un score d'attention à partir de la sortie LSTM, via la classe **Dense** ;
- Une fonction d'aplatissement des scores d'attention, via la classe **Flatten** ;
- Une fonction d'activation *softmax* pour normaliser les poids d'attention, via la classe **Activation** ;
- Une fonction de répétition des poids d'attention normalisés pour chaque pas de temps (ou « unité LSTM »), via la classe **RepeatVector** ; et
- Une fonction de permutation pour adapter les dimensions à la suite du processus, via la classe **Permute**.

Tous ces éléments sont unifiés pour former notre LSTM via la classe **Model**. Il sera compilé avec l'optimiseur **Adam** et utilise le MSE comme fonction de perte.

Similairement au BRNN et au MLP, voici les hyperparamètres sélectionnés pour l'optimisation :

- Le nombre de neurones de la couche de neurones LSTM, noté *nbUnitsLSTM* ;
- Les coefficients de régularisation L2 de chacune de ces couches, notés *L2_1LSTM* et *L2_2LSTM* ;
- Le nombre d'échantillons par lot, noté *tailleLotsLSTM* ; et
- Le taux d'apprentissage placé en argument de l'optimiseur **Adam**, noté *tauxApprLSTM*.

Nous créons maintenant une fonction personnalisée `ConstructeurOptiLSTM`, identique à `ConstructeurOptiBRNN` une fois que l'on a remplacé la structure du BRNN par celle de notre LSTM et que l'on retire la sélection du nombre de neurones de la seconde couche cachée du BRNN.

Les ensembles de sélection des hyperparamètres que nous avons retenus seront décrits dans le tableau 5.8, ainsi que la classe (`Choice` ou `Int`) que nous appliquerons à chacun d'eux.

TABLEAU 5.8 Ensembles de sélection des hyperparamètres du LSTM et classe associée

Classe	Hyperparamètre	Ensemble de Sélection
Choice	$L2_1LSTM$	$\{1 \times 10^{-2}, 3 \times 10^{-3}, 1 \times 10^{-3}, 3 \times 10^{-4}, 1 \times 10^{-4}\}$
	$L2_2LSTM$	$\{1 \times 10^{-2}, 3 \times 10^{-3}, 1 \times 10^{-3}, 3 \times 10^{-4}, 1 \times 10^{-4}\}$
	$tauxApprLSTM$	$\{1 \times 10^{-2}, 3 \times 10^{-3}, 1 \times 10^{-3}, 3 \times 10^{-4}, 1 \times 10^{-4}\}$
Int	$tailleLotsLSTM$	$\{a = 32, b = 512, k = 32\}$
	$nbUnits1LSTM$	$\{a = 16, b = 128, k = 16\}$

On encapsule ensuite `ConstructeurOptiLSTM` dans `Capsule`.

On applique `RandomSearch` aux arguments suivants :

- Un hypermodèle, qui est notre `Capsule` appliquée à `ConstructeurOptiLSTM`;
- Un critère de perte, qui sera `val_loss` et qui utilisera donc le MSE;
- Un nombre maximal d'essais, que l'on limite à 100 pour des raisons computationnelles ;
et
- Un nombre d'exécutions par essai, que l'on limite à 2 pour éviter que cette optimisation ne dépasse une durée d'une heure.

On applique la méthode `search` de `RandomSearch` à $X_{LSTM,appr}$ et $y_{LSTM,appr}$ en fixant le nombre d'époques maximal à 100 et en utilisant 20% de ces données pour chaque étape de la validation croisée. Les valeurs optimales, retournées par la méthode `get_best_hyperparameters`, sont données dans le tableau 5.9.

Enfin, comme évoqué, on teste différentes valeurs du coefficient k pour surpondérer pertinemment les conditions opérationnelles du trajet en cours. On entraîne alors successivement le LSTM pour chacune des valeurs suivantes : $\{1; 3; 7; 10; 25; 100\}$ et on observe l'évolution des critères de performance. Il semble alors clair que $k = 7$ est l'option la plus pertinente. On teste donc plusieurs fois notre LSTM avec les valeurs $\{3; 7; 10\}$ pour différentes séparations des échantillons d'apprentissage/test et on confirme qu'en général, les prédictions sont

TABLEAU 5.9 Combinaison optimale des hyperparamètres de notre LSTM

Hyperparamètre	Valeur optimale
$nbUnitsLSTM$	48
$L2_1LSTM$	3×10^{-3}
$L2_2LSTM$	1×10^{-4}
$tailleLotsLSTM$	416
$tauxApprLSTM$	3×10^{-3}

meilleures avec $k = 7$. On fixe définitivement cette valeur. On remarque par ailleurs que, en plus de l'amélioration des critères de performance, l'histogramme des TT prédits voit sa variance augmenter. Cela pourra faciliter le travail du modèle d'empilement. En effet, les conditions opérationnelles les plus difficiles sont alors mieux mises en évidence à l'avance par le LSTM.

Explicitons plus en détail l'obtention des résultats finaux. Pour cela, on utilise dans un premier temps la combinaison optimale obtenue pour entraîner notre LSTM sur $X_{LSTM,appr}$ et $y_{LSTM,appr}$ via la méthode `fit`, en fixant là encore le nombre d'époques maximal à 100 et en utilisant 20% de ces données pour sa validation croisée. On lui demande ensuite de formuler ses prédictions de TT, notées $y_{LSTM,pred}$, d'après $X_{LSTM,test}$ via la méthode `predict`. On lui fera aussi formuler ses prédictions de $y_{LSTM,appr}$ d'après $X_{LSTM,appr}$.

Tous nos modèles de ML ont désormais formulé leurs prédictions de TT, après optimisation de leurs hyperparamètres. Nous comparerons leur précision respective à la section 5.7. Forts de notre expérience de prédiction de TT sur les itinéraires au complet, nous allons maintenant nous pencher sur la prédiction de TT sur des itinéraires segmentés.

5.6.1.2 Prédiction de TT sur l'itinéraire sectionné en segments inter-niveaux

La segmentation des itinéraires provoque des problématiques additionnelles, décrites dans la sous-section 5.3.2. Nous mettrons ici en application la combinaison de solutions qui y est évoquée pour améliorer la prédiction agrégée du TT complet.

En amont de toute autre chose, rappelons que nous allons maximiser le volume de données d'apprentissage sur chacun des segments de l'itinéraire complet étudié. Nous allons donc tâcher de détecter la totalité des trajets ayant traversé chacun de ces segments dans le sens

désiré, qu'ils aient fait partie d'un trajet sur l'itinéraire complet ou non. Pour autant, nous désirons uniquement prédire des TT par segment qui se sont déroulés durant des trajets complets. Nous devons donc disposer aussi des ensembles de données relatifs à ces trajets complets, qui devront en particulier inclure chacun des TT par segment observés.

On applique ainsi notre modèle de préparation des données à l'itinéraire complet surface→niveau350 en relevant, sur chacun des segments inter-niveaux, les TT par segment réels correspondants et les données opérationnelles associées. Dans le cas de sites miniers pour lesquels les balises de changement de niveau sont peu fiables, ces TT par segment n'existent pas nécessairement et on laissera donc le cas échéant une valeur vide (« *None* » en Python) dans la liste de TT et dans la liste des conditions opérationnelles associées. Pour le i -ème segment, les TT par segment réels seront appelés « y_i » et leurs conditions opérationnelles « X_i ». On relève au passage aussi X et y . On sépare alors simultanément X , y , les X_i et les y_i en ensembles d'entraînement et de test, de manière à faire correspondre les éléments de $X_{i,appr}$ avec ceux de X_{appr} (de même pour les autres ensembles). On repère chaque trajet des ensembles de test par un identifiant unique, qui correspond à la concaténation de l'horodatage du trajet et de l'identifiant du HT à l'origine de ce trajet.

Pour obtenir les autres ensembles de données, on applique cette fois-ci notre modèle de préparation de données à chacun des segments inter-niveaux de la même manière que nous l'utiliserions sur un itinéraire complet. Seulement, on tâche d'appliquer des contraintes raisonnées aux trajets qui traversent le segment. On veut en effet s'assurer que les HT aient adopté une trajectoire significativement similaire à celle qu'ils auraient adoptée s'ils réalisaient un trajet sur l'itinéraire au complet. On impose donc des détections obligatoires de balises précédentes cohérentes pour chaque segment, comme on le ferait avant la détection de la balise A d'un itinéraire complet. On rassemblera ainsi généralement les données de plus d'une dizaine de milliers d'observations de trajets sur chacun des segments de l'itinéraire surface→niveau350, avec leurs conditions opérationnelles respectives et leur TT respectif. Étape incontournable, on exclut de ces observations tous les trajets qui correspondent aux identifiants uniques créés auparavant d'après les données de test. On dispose alors de matrices de TT que l'on dénomme « $y_{i,tot}$ » et de matrices de conditions opérationnelles « $X_{i,tot}$ ».

Alternativement, nous devons bonifier ces dernières matrices via l'utilisation du GMM et du MLP pour obtenir les « $X_{i,tot,GMM}$ ». L'utilisation de l'AIC pour le GMM nous montre qu'un partitionnement à quatre classes est généralement optimal sur $y_{i,tot}$. On fixe donc le nombre de classes à quatre quel que soit le segment inter-niveaux de la mine 1 que l'on cherche à partitionner, pour simplifier notre modèle. On optimise alors les hyperparamètres de notre MLP pour générer de bonnes prédictions de classes d'appartenance. Les valeurs optimales

retournées, que l'on jugera aussi valables pour tout segment inter-niveaux de la mine 1, sont décrites dans le tableau 5.10.

TABLEAU 5.10 Combinaison optimale des hyperparamètres du MLP

Hyperparamètre	Valeur optimale
$nbUnits1MLP$	192
$nbUnits2MLP$	8
$tailleLotsMLP$	32
$tauxApprMLP$	1×10^{-2}

En concaténant les prédictions du MLP optimisé avec nos données d'entrée déjà générées, on obtient finalement les $X_{i,tot}$. On utilise aussi les ensembles $X_{i,tot}$ et $y_{i,tot}$ pour générer ceux dédiés au LSTM. On ne désignera pas explicitement ces ensembles par la suite, puisqu'ils se manipulent de la même manière que $X_{i,tot}$ et $y_{i,tot}$.

Nous postulons que l'optimisation des hyperparamètres de nos modèles de ML sur les données d'un segment bien choisi les rendra suffisamment robustes pour généraliser aux autres segments. Ainsi, chaque modèle sera optimisé une seule fois sur les ensembles $X_{i,tot}$ et $y_{i,tot}$ du segment i , jugé le plus représentatif des segments inter-niveaux. Naturellement, cette décision n'est pas anodine. Comparativement à une optimisation des hyperparamètres sur chaque segment, nous risquons de faire baisser la précision des prédictions de nos modèles de ML sur les segments les plus spécifiques. Pour autant, cette décision sera la bienvenue pour limiter la durée de compilation totale de notre modèle de prédiction de TT en dessous de quatre heures sur un nouvel itinéraire donné. Elle simplifiera aussi appréciablement la structure de notre modèle. Ainsi, notre optimisation des mêmes combinaisons d'hyperparamètres que sur les itinéraires complets aboutit aux valeurs indiquées dans le tableau 5.11.

Pour chaque segment, on ajuste alors nos modèles de ML avec leurs hyperparamètres respectifs sur les $X_{i,tot}$ et les $y_{i,tot}$ puis on leur fait formuler leurs prédictions $y_{i,pred}$ d'après les $X_{i,test}$.

Bien que cette dernière étape soit triviale dans le cas général, elle devient beaucoup plus délicate dès lors que l'on utilise un LSTM dans le cas particulier où l'on tolère que certaines balises de changement de niveau ne seraient pas fiables. En effet, alors que les autres modèles de ML peuvent simplement éviter les données manquantes qui apparaîtraient dans les données d'entraînement, le LSTM a besoin d'avoir accès aux TT des trajets précédents. Si l'un d'entre eux est vide, le LSTM prédit une valeur vide. Dans une telle situation, on devra

TABLEAU 5.11 Combinaisons optimales d'hyperparamètres pour nos six modèles appliqués à un segment inter-niveaux

Modèle	Hyperparamètre	Valeur optimale
XGBRegressor	learning_rate	7×10^{-3}
	max_depth	6
	n_estimators	364
	subsample	0,5
RandomForestRegressor	max_depth	10
	max_features	'log2'
	min_samples_leaf	7
	min_samples_split	12
	n_estimators	175
GradientBoostingRegressor	learning_rate	1×10^{-2}
	max_depth	4
	n_estimators	196
	max_features	'log2'
	min_samples_leaf	3
	min_samples_split	5
SVR	C	95
	epsilon	1×10^{-1}
	kernel	'linear'
BRNN	<i>nbUnits1BRNN</i>	32
	<i>L2_1BRNN</i>	1×10^{-3}
	<i>nbUnits2BRNN</i>	32
	<i>L2_2BRNN</i>	1×10^{-3}
	<i>tailleLotsBRNN</i>	160
	<i>tauxApprBRNN</i>	1×10^{-1}
LSTM	<i>nbUnitsLSTM</i>	32
	<i>L2_1LSTM</i>	1×10^{-3}
	<i>L2_2LSTM</i>	1×10^{-3}
	<i>tailleLotsLSTM</i>	408
	<i>tauxApprLSTM</i>	5×10^{-4}

donc nous-même compléter ces TT dans $X_{i,test}$. Pour cela, il serait théoriquement possible de remplacer tous les TT manquants par la valeur moyenne des TT d'apprentissage observés sur le i-ème segment. Pour autant, nous jugeons que notre perte en précision pourrait être inacceptable dans le cas de dysfonctionnements majeurs de certaines balises. Pour pallier cette problématique, et étant donné que notre LSTM est pré-entraîné sur un ensemble plus vaste, nous décidons de compléter toutes les valeurs vides par des prédictions de TT (excepté

les tous premiers trajets de la période étudiée, dont on ne peut pas prédire le TT car ils ne sont pas précédés d'un nombre suffisant de trajets, on prédit donc la moyenne des TT pour ces quelques trajets qui figurent dans $X_{i,test}$). Pour un i -ème segment donné, il est nécessaire de partir du tout début de X_i , de prédire uniquement la valeur du TT suivant, de compléter X_i avec cette valeur (lorsque cela a un sens) et ainsi de suite. Progressivement, en suivant l'ordre des milliers de trajets complets observés, on peut remplir tous les TT inter-niveaux manquants du i -ème segment. De même pour chaque segment.

Pour en revenir au cas de la mine 1, étant donné que nous disposons de plus d'une centaine d'observations de trajets complets sur l'itinéraire étudié, nous considérons que l'application d'une régression des résidus est pertinente. Comme mentionné à la sous-section 5.3.2, l'objectif est de réduire les biais éventuellement induits dans nos modèles au cours de l'entraînement du fait de l'utilisation sûrement massive de trajets par segment non inclus dans les itinéraires complets. On s'appuie donc sur les données des trajets par segment qui sont inclus dans les itinéraires complets pour réajuster nos modèles. Tâchons de détailler cette méthode sur p. ex. le i -ème segment :

1. On fait prédire à nos modèles de ML les TT par segments contenus dans les données d'**entraînement** correspondant aux observations de trajets complets, soit les $y_{i,appr}$ d'après les $X_{i,appr}$;
2. On détermine les « résidus », que l'on calcule comme les différences entre les prédictions que l'on a généré au point précédent et les valeurs des $y_{i,appr}$;
3. On entraîne une instance par défaut de la classe `LinearRegression` à prédire ces résidus d'après $X_{i,appr}$;
4. On fait prédire à cette instance les résidus qui devraient se trouver dans les $y_{i,pred}$ d'après les $X_{i,test}$; et
5. On corrige enfin entièrement les $y_{i,pred}$ de tous nos modèles de ML d'après les prédictions de résidus finalement obtenues.

On conserve la dénomination « $y_{i,pred}$ » pour décrire l'ensemble des prédictions, maintenant modifiées. Dans le cas d'un site minier dans lequel les balises de changement de niveau ne seraient pas toujours fiables, on prédirait aussi à ce stade les valeurs vides des $y_{i,appr}$ en combinant le LSTM et le modèle de régression des résidus. On remplirait ainsi tous les espaces vides de ces ensembles pour permettre le bon fonctionnement du modèle suivant.

Passons enfin à l'agrégation des $y_{i,pred}$ pour former y_{pred} via un modèle de régression linéaire multiple, évoqué au cours de la sous-section 5.3.2. L'utilisation de ce modèle est là aussi rendue pertinente par l'existence de plus d'une centaine d'observations de trajets complets sur l'itinéraire étudié.

Pour un itinéraire donné, découpé en n segments, l'équation de régression d'une régression linéaire multiple peut s'écrire :

$$TT_{complet}(TT_1, TT_2, \dots, TT_n) = \sum_{i=1}^n k_i TT_i + C$$

avec k_1, k_2, \dots, k_n des coefficients réels, ici positifs, et C une constante réelle.

On calcule le MSE entre les valeurs prédites et les valeurs vraies, i.e. ici pour un nombre total de trajets complets noté « N » : $MSE = \frac{1}{N} \sum_{j=1}^N (\sum_{i=1}^n k_i TT_{i,j} + C - TT_{complet,j})^2$

Rappelons que la subtilité de notre modèle de régression linéaire multiple réside dans le fait que nous allons chercher à pénaliser les régressions qui se distinguent trop d'une simple somme des TT observés sur chaque segment. On notera ainsi « Pénalité » la somme de la distance quadratique de chaque coefficient k_i au coefficient théoriquement attendu devant chaque segment, i.e. donc 1. On peut alors écrire : $Pénalité = \sum_{i=1}^n (k_i - 1)^2$

Ainsi, notre fonction de coût personnalisée pourra être donnée par l'expression suivante :

$$\text{Coût}(TT_{complet,[1,N]}, TT_{[1,n],[1,N]}, k_{[1,n]}, \lambda_{reg}) = \text{MSE} + \lambda_{reg} \times \text{Pénalité} \quad (5.1)$$

avec λ_{reg} le paramètre de régularisation. Plus la valeur de ce paramètre est élevée et plus notre modèle devra chercher des valeurs des k_i proches de 1. Si ce coefficient tendait vers l'infini, il transformerait donc notre modèle en simple somme avec une constante additionnelle.

On cherche maintenant une valeur de λ_{reg} qui permette d'aboutir à des coefficients k_i notablement proches de 1 pour la plupart, tout en laissant une marge de manoeuvre au modèle pour lui permettre de générer des prédictions plus précises. Cette recherche d'une valeur pertinente pour ce paramètre est en grande partie subjective puisqu'on ne cherche pas nécessairement à atteindre la meilleure performance de régression. On préfère donc adapter λ_{reg} manuellement, et on trouve ici qu'une valeur de 1×10^4 donne globalement les résultats escomptés, avec de très bons résultats de régression et des coefficients k_i tous situés entre 0,90 et 1,03, dont près de la moitié sont pour ainsi dire égaux à 1. Étant donné que la valeur de ce paramètre aura été adaptée subjectivement et qu'elle n'aura donc pas été « optimisée », on considère qu'elle sera adaptée à une utilisation générale pour tous les itinéraires, lorsqu'ils sont segmentés par sous-niveaux. Plus concrètement, pour tous les itinéraires, on admet que la précision des prédictions de notre modèle ne risque aucunement d'être dégradée par cette unique adaptation de la valeur de λ_{reg} . Ce choix définitif vise bien évidemment à automatiser l'utilisation de notre solution. On comparera les prédictions obtenues avec cette technique d'agrégation personnalisée comparativement à celles d'une simple somme des TT par segment.

Pour chacun de nos modèles de ML, on entraîne le modèle de régression linéaire multiple à prédire y_{appr} d'après les $y_{i,appr}$. On lui demande enfin de formuler sa prédiction finale de TT pour chacun de nos modèles de ML d'après les $y_{i,pred}$ issus du modèle de régression des résidus. On lui fera aussi prédire y_{appr} d'après $y_{i,appr}$.

Dans le cas où certaines balises de changement de niveau ne seraient pas fiables et où l'on ne dispose donc pas de tous les TT inter-niveaux, on devra nous-même compléter $y_{i,appr}$ via la combinaison du LSTM puis lui appliquer le modèle de régression des résidus pour enfin alimenter le modèle de régression linéaire multiple. Ce dernier fonctionne de manière identique.

5.6.1.3 Prédiction de TT sur l'itinéraire sectionné en segments majeurs

Cette étape facultative est globalement très similaire à la précédente. Voici les différences que nous avons pu constater :

- Contrairement aux segments inter-niveaux (et équivalents dans les rampes simples), généralement prédéfinis, les segments majeurs que l'on juge pertinent d'utiliser doivent être désignés explicitement pour cette étape ;
- Le modèle de préparation de données s'appliquera à chacun des segments majeurs, ce qui demande de jalonner les trajets avec des balises inter-niveaux comme sur un trajet complet (les mêmes balises sont aussi utilisées pour la collecte des données du trajet complet) ;
- On doit appliquer le modèle de préparation de données à l'itinéraire complet en récupérant au passage des données sur les segments majeurs, ce qui n'est pas plus laborieux que de récupérer des données sur les segments inter-niveaux mais ce qui nous force ici à implémenter manuellement cette collecte ;
- Le paramètre de régularisation λ_{reg} de la régression linéaire multiple doit être adapté aux segments majeurs pour limiter encore plus la variation des k_i puisque l'utilisation de plus longs segments, moins nombreux, réduit indéniablement les problématiques d'incertitudes locales ; et
- Enfin, point non négligeable, on peut récupérer directement toutes les données de prédiction obtenues sur les segments inter-niveaux que l'on désire combiner à celles obtenues sur le(s) segment(s) majeur(s) pour reformer le TT complet.

Nous utiliserons ici ce dernier point pour récupérer directement les données de prédiction des trois derniers segments inter-niveaux de l'itinéraire surface→niveau350. On remplacera ainsi le segment de l'itinéraire allant de l'entrée du niveau 300, côté rampe 2, à l'entrée du

niveau 350, soit la traversée du niveau 300 et deux véritables changements de niveau. On prédit ainsi l'agglomération de segments de différentes longueurs.

L'utilisation de l'AIC pour le GMM nous montre qu'un partitionnement à six classes est optimal pour l'itinéraire surface→niveau300. Le partitionnement résultant est représenté à la figure 5.9. On fixera le nombre de classes à six quel que soit le segment majeur de la mine 1 que l'on cherche à partitionner, pour simplifier notre modèle.

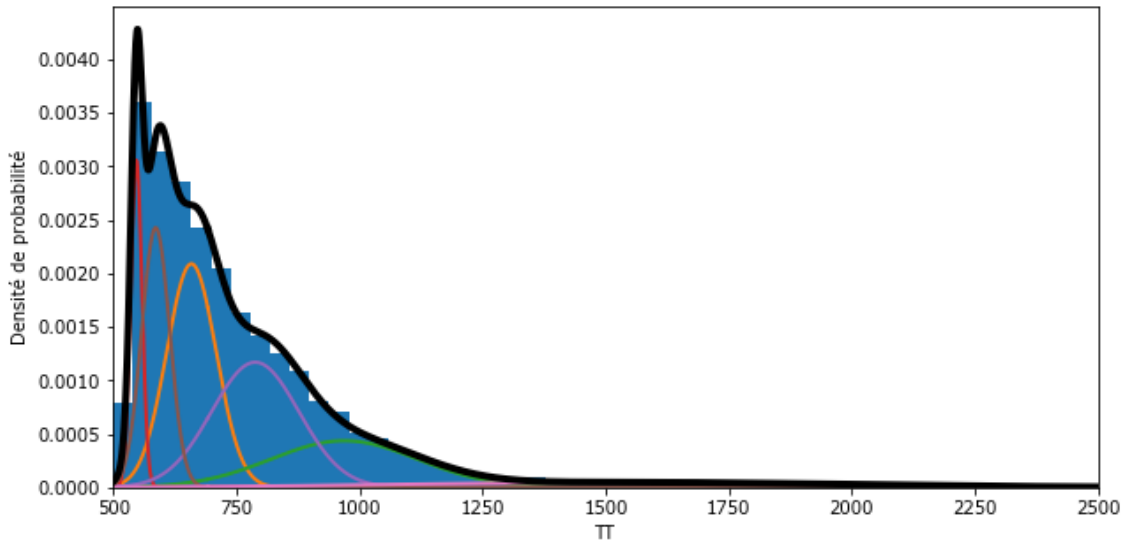


FIGURE 5.9 Inférence de six groupes par le GMM dans y_{appr} pour le segment majeur surface→niveau300 via la rampe 2 dans la mine 1

On optimise là encore les hyperparamètres du MLP. Les valeurs optimales retournées, que l'on jugera aussi valables pour tout segment majeur de la mine 1, sont décrites dans le tableau 5.12.

TABLEAU 5.12 Combinaison optimale des hyperparamètres du MLP

Hyperparamètre	Valeur optimale
$nbUnits1MLP$	160
$nbUnits2MLP$	8
$tailleLotsMLP$	48
$tauxApprMLP$	1×10^{-2}

On représente aussi la comparaison des diagrammes à barres respectifs des proportions de TT dans chacune des classes selon le GMM (vraies valeurs) et selon le MLP (prédictions) à

la figure 5.10.

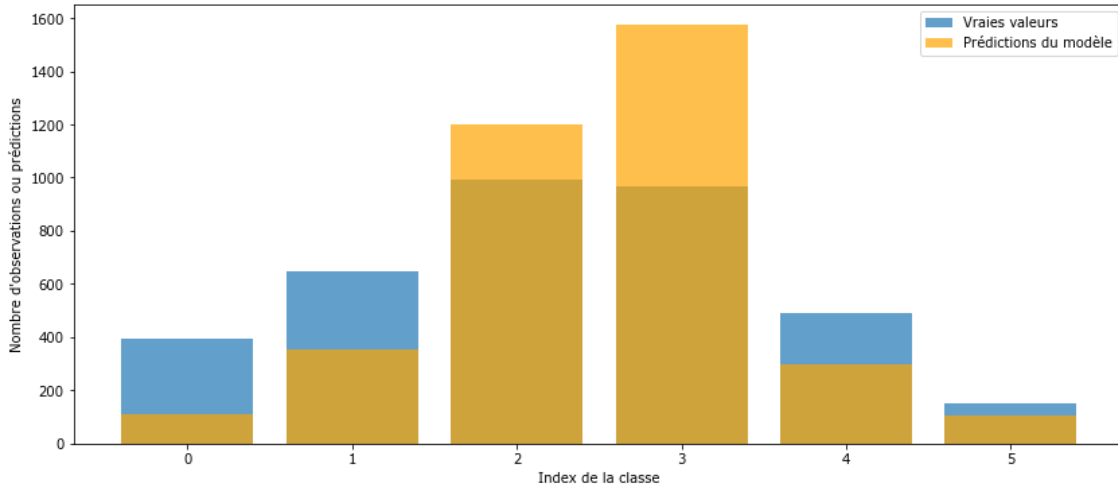


FIGURE 5.10 Diagramme à barres des vrais groupes annoncés par le GMM et des groupes prédits par le MLP pour l'itinéraire surface→niveau350 via la rampe 2 dans la mine 1

L'optimisation des hyperparamètres de nos modèles de ML sur le segment majeur qui nous intéresse, i.e. l'itinéraire surface→niveau300, aboutit aux valeurs inventoriées dans le tableau 5.13.

Par ailleurs, on trouve empiriquement que λ_{reg} provoque une régression que l'on juge totalement satisfaisante lorsqu'il adopte une valeur de 1×10^5 : le coefficient k_1 devant notre segment majeur est quasiment égal à 1 et les autres k_i sont situés entre 0,95 et 1. On fixe aussi cette valeur pour tous les itinéraires sectionnés en segments majeurs de la mine 1.

En prenant toutes les remarques précédentes en compte et en suivant le cheminement que nous avons suivi pour les itinéraires sectionnés en segments inter-niveaux, on obtient les ensembles de données de prédiction de chacun de nos modèles de ML pour notre itinéraire sectionné avec un segment majeur.

À l'issue de la présente sous-section, on sélectionne les modèles de ML qui ont respectivement montré les meilleures performances pour chaque type de segmentation. Nous récupérons ensuite les ensembles de prédictions pour les trajets de test ainsi que pour ceux d'apprentissage.

Nous allons pouvoir à présent nous intéresser à notre modèle d'empilement, qui sera le sujet de la prochaine sous-section.

L'analyse des performances sera quant à elle présentée à la section 5.7.

TABLEAU 5.13 Combinaisons optimales d'hyperparamètres pour nos six modèles appliqués à un segment majeur

Modèle	Hyperparamètre	Valeur optimale
XGBRegressor	learning_rate	1×10^{-2}
	max_depth	5
	n_estimators	428
	subsample	0,8
RandomForestRegressor	max_depth	10
	max_features	'sqrt'
	min_samples_leaf	1
	min_samples_split	12
	n_estimators	192
GradientBoostingRegressor	learning_rate	2×10^{-2}
	max_depth	4
	n_estimators	204
	max_features	'sqrt'
	min_samples_leaf	2
	min_samples_split	4
SVR	C	95
	epsilon	1×10^{-1}
	kernel	'linear'
BRNN	<i>nbUnits1BRNN</i>	128
	<i>L2_1BRNN</i>	3×10^{-2}
	<i>nbUnits2BRNN</i>	28
	<i>L2_2BRNN</i>	1×10^{-2}
	<i>tailleLotsBRNN</i>	480
	<i>tauxApprBRNN</i>	1×10^{-1}
LSTM	<i>nbUnitsLSTM</i>	272
	<i>L2_1SLTM</i>	4×10^{-3}
	<i>L2_2LSTM</i>	2×10^{-1}
	<i>tailleLotsLSTM</i>	408
	<i>tauxApprLSTM</i>	1×10^{-3}

5.6.2 Prédiction finale du TT par un modèle d'empilement

Intéressons-nous maintenant à l'implémentation puis à l'optimisation de notre modèle d'empilement, lequel constitue indéniablement une pièce maîtresse de notre méthodologie.

Dans un premier temps, nous devons structurer convenablement l'ensemble des données d'entrée qu'il prendra en argument. Pour cela, on concatène les vecteurs et matrices suivants :

- L'ensemble des prédictions générées par le meilleur modèle de prédiction de TT d'un itinéraire pris dans son ensemble ;
- L'ensemble des prédictions générées par le meilleur modèle de prédiction de TT d'un itinéraire sectionné en segments inter-niveaux ;
- L'ensemble des prédictions générées par le meilleur modèle de prédiction de TT d'un itinéraire sectionné en segments majeurs ;
- L'ensemble des prédictions de classes d'appartenance de chaque trajet générées par le MLP ;
- Un vecteur de longueur identique aux cinq autres, mais dont chaque élément correspondra à la moyenne des trajets d'entraînement. Point d'attention : on aura *log*-normalisé aussi cette valeur et cela ne donnera pas une valeur nulle puisque c'est la distribution transformée par le logarithme qui aura été normalisée. La moyenne de cette dernière est différente. L'idée est de créer un point de repère fixe dont la valeur est suffisamment fiable pour que le modèle d'empilement s'y rattache lorsque les prédictions initiales des autres modèles entrent fortement en contradiction. Si le modèle d'empilement a tendance à se fier excessivement à ce point de repère, au point de prédire quasiment toujours la moyenne des TT, on l'enlèvera des données de prédiction.

Ensuite, on tâche d'optimiser les hyperparamètres des quatre modèles de régression envisagés, qui correspondent aux classes `LinearRegression`, `GradientBoostingRegressor`, `RandomForestRegressor`, et `SVR`. Comme énoncé au cours de la sous-section 5.3.2, nous utiliserons pour ce faire la bibliothèque `optuna`.

On crée une fonction principale, `Objective`, qui configure `optuna` pour tester différentes combinaisons d'hyperparamètres pour chacun des quatre modèles de ML testés. Ces combinaisons sont rassemblées dans le tableau 5.14.

On intègre à `Objective` une validation croisée à trois plis sur l'ensemble d'entraînement pour rechercher un modèle avec une très bonne capacité de généralisation, via la méthode `cross_val_score` du module `model_selection` de la bibliothèque `scikit-learn`. On demande à `Objective` de retourner finalement la RMSE du score de validation obtenu à l'issue de cette validation croisée.

On crée une « étude » Optuna visant à minimiser la métrique en sortie d'`Objective`, ici la RMSE, via la méthode `create_study` et l'argument `direction="minimize"`. On utilise alors la méthode `optimize` de l'étude Optuna en lui donnant pour arguments notre fonction principale `Objective` et un nombre d'essais, fixé à 200.

Après la compilation, on extrait le meilleur essai de notre étude Optuna via la méthode

TABLEAU 5.14 Ensembles de sélection des hyperparamètres pour les modèles utilisés

Modèle	Hyperparamètre	Ensemble de sélection
Linear Regression	fit_intercept	{True, False}
	learning_rate	$\mathcal{U}(1 \times 10^{-3}, 1 \times 10^{-1})$
Gradient Boosting	max_depth	$\mathcal{U}\{2, 5\}$
	n_estimators	$\mathcal{U}\{50, 500\}$
	subsample	$\mathcal{U}(0,5, 1)$
Random Forest	max_depth	$\mathcal{U}\{2, 5\}$
	max_features	{None, 'sqrt', 'log2'}
	n_estimators	$\mathcal{U}\{50, 500\}$
SVR	C	$\mathcal{U}(0,1, 100)$
	epsilon	$\mathcal{U}(1 \times 10^{-2}, 1)$
	kernel	{'linear', 'poly', 'rbf', 'sigmoid'}

`best_trial`. L’affichage des paramètres correspondants à la meilleure étude nous indiquent alors que le meilleur modèle est le **GradientBoostingRegressor**. Ses hyperparamètres optimaux sont quant à eux indiqués dans le tableau 5.15.

TABLEAU 5.15 Combinaison optimale des hyperparamètres du GBR, modèle d’empilement optimal

Hyperparamètre	Valeur optimale
<i>n_estimators</i>	445
<i>max_depth</i>	2
<i>learning_rate</i>	$1,35 \times 10^{-2}$
<i>subsample</i>	0,8

On entraîne alors ce modèle optimal avec ses hyperparamètres optimaux sur notre ensemble d’entraînement. Pour finir, on lui demande de formuler nos prédictions finales de TT sur l’ensemble de données de test, ce qui complète la mise en œuvre de notre modèle global de prédiction de TT.

Comme les autres points de questionnement, la pertinence de notre modèle d’empilement sera évaluée dans la prochaine section.

5.7 Évaluation de la qualité des prédictions de TT

Passons maintenant en revue la grande variété de modèles de ML et de sous-versions différentes de notre modèle de prédiction de TT. Comme évoqué lors de notre présentation de notre méthodologie, au chapitre 3, nous nous appuyerons sur la RMSE pour évaluer la qualité des prédictions. L'utilisation de la MAE pourrait être pertinente pour les planificateurs et nous l'utiliserons donc conjointement pour éviter de tirer des conclusions trop hâtives. Par ailleurs, rappelons que le modèle qualifié de « modèle de référence » consiste à prédire invariablement la moyenne du TT de la distribution.

Commençons directement par l'évaluation de la pertinence d'utiliser le tandem GMM-MLP en amont de nos trois sous-modèles de prédictions initiales de TT. Pour cela, on compare les performances globales de l'ensemble de nos six modèles de ML lorsqu'ils utilisent les données correspondant aux classes d'appartenance prédites, et lorsque ce n'est pas le cas. Tous les modèles voient une amélioration *a minima* légère de leurs performances grâce aux prédictions de probabilités des classes d'appartenance par le MLP, qu'ils aient été optimisés ou non pour utiliser ces données. Le LSTM, en particulier, profite d'une plus grande amélioration, qui dépasse 5% de réduction relative de la MAE et plus de 3% de réduction relative de la RMSE. Rappelons que le tandem GMM-MLP pourra aussi permettre d'améliorer les prédictions finales du modèle d'empilement. Sa pertinence pour le partitionnement des données n'étant plus à prouver, nous n'avons pas évalué l'amplitude de son impact sur ces dernières.

Tâchons maintenant de comparer les résultats des modèles de ML tâchés de formuler les prédictions initiales de TT, aidés des prédictions du MLP.

Pour les prédictions sur l'itinéraire complet, le SVM se distingue particulièrement des autres modèles par ses résultats médiocres. Ces derniers ont des performances plus proches, qui brillent très peu comparativement au modèle de référence comme le montre le tableau 5.16. On y fait aussi figurer les résultats du SVM, à titre de comparaison. Mis à part le BRNN, nos autres modèles n'ont pas la fâcheuse tendance de prédire uniquement des TT que l'on pourrait juger excessivement rassemblés autour de la moyenne. En particulier, le SVM reproduit remarquablement bien la distribution réelle des TT de test au cours de ses prédictions (voir fig. 5.11), alors qu'il donne de loin les pires résultats de prédiction.

Pour l'itinéraire segmenté avec un segment majeur, le constat est bien différent. En effet, sur l'itinéraire surface→niveau300, le LSTM surperforme largement tous les autres modèles, avec près de 10% de réduction relative de la RMSE et de la MAE. Il n'y a aucun doute qu'il s'agit ici du meilleur modèle d'après nos critères de performance.

Pour l'itinéraire segmenté en sous-segments, le constat change encore. Le BRNN est sous-

TABLEAU 5.16 Performances de nos six modèles de ML et du modèle de référence (« Référence ») pour l’itinéraire complet

Modèle	MAE	RMSE
Référence	213	307
BRNN	204	307
XGBoost	214	312
RF	210	312
GBR	208	309
SVM	226	335
LSTM	199	312

performant, tout segment confondu, et les performances relatives des cinq autres modèles varient assez amplement les unes par rapport aux autres, donnant parfois l’avantage à XGBoost et parfois au LSTM. Cela s’explique certainement par la variété de segments rencontrés, qui seraient plus difficilement généralisables que ce que nous avons espéré.

Au final, nous ne retiendrons qu’un unique modèle de ML pour réaliser toutes les prédictions initiales : il s’agira du LSTM, mais ce choix ne repose pas uniquement sur l’ensemble des performances présentées ci-avant. En effet, lors de nos expérimentations, lorsque nous appliquions nos modèles de ML (déjà optimisés sur un itinéraire quelconque) sur un itinéraire de longueur intermédiaire sur lequel leurs hyperparamètres n’avaient pas été optimisés, le LSTM avait un net avantage. En particulier, sur l’itinéraire surface→niveau300, ses résultats de prédiction étaient meilleurs que sur l’itinéraire surface→niveau350 par rapport au modèle de référence, alors même qu’il avait été optimisé sur ce dernier. Cette observation est conforme à l’hypothèse selon laquelle les TT des itinéraires les plus longs sont les plus difficiles à prédire. Pour autant, tous les autres modèles de ML voyaient leurs performances s’effondrer dans les mêmes conditions. Le LSTM serait donc capable de mieux généraliser qu’eux à d’autres itinéraires, sans qu’une optimisation des hyperparamètres soit absolument nécessaire. Cette robustesse caractéristique joue en sa faveur, et, combinée à ses résultats souvent satisfaisants après optimisation, elle nous convainc de ne retenir que lui pour ces trois étapes de notre modèle global, et de se défaire donc des cinq autres modèles de ML.

Revenons-en maintenant aux itinéraires segmentés. Il s’agirait d’observer les performances respectives du modèle de régression des résidus, puis du modèle de régression linéaire multiple

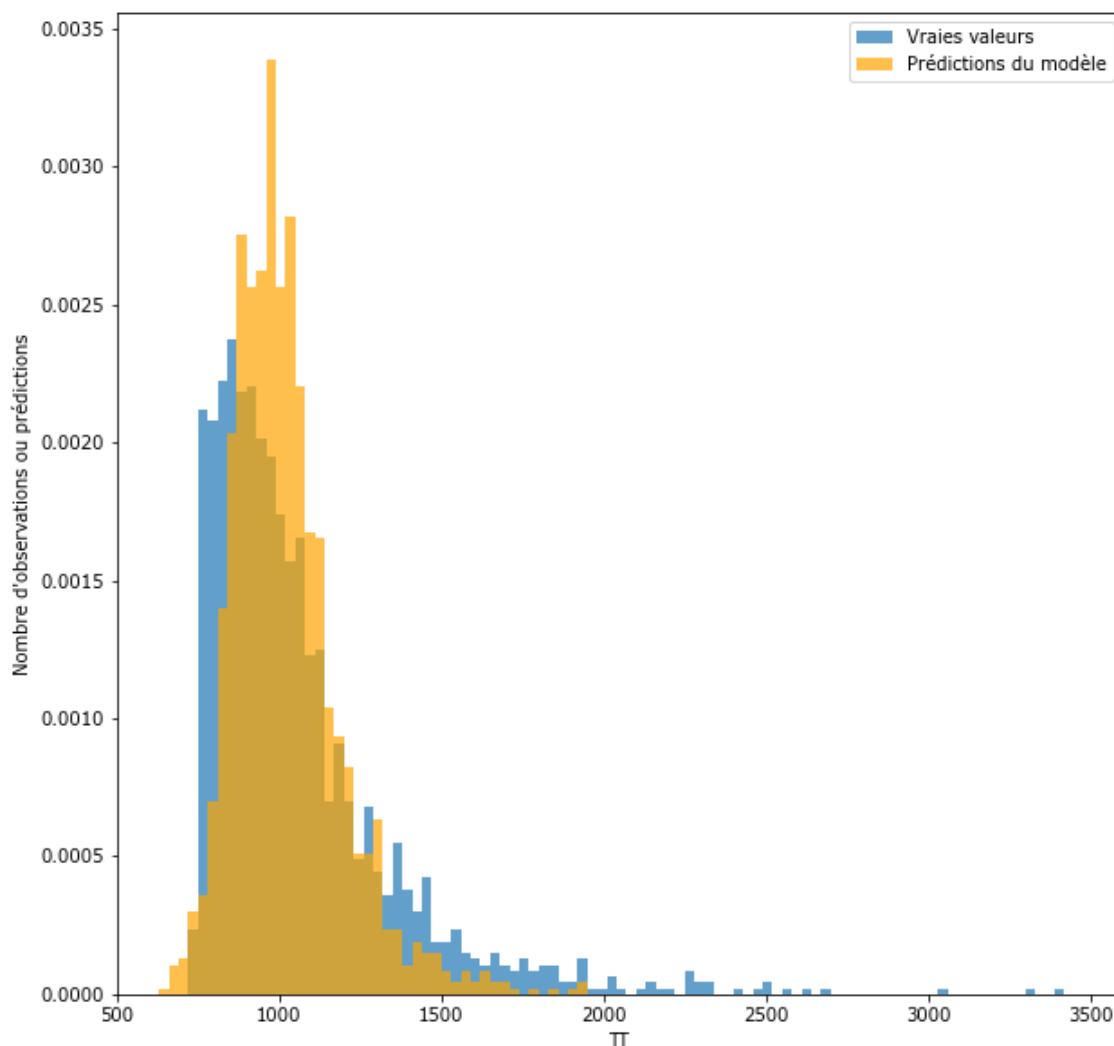


FIGURE 5.11 Histogrammes des TT réels et des TT prédits par le SVM pour l’itinéraire complet

chargé de l’agrégation.

Concernant le premier, on obtient des résultats tout à fait satisfaisants malgré le fait que, pour la seconde stratégie de segmentation, ces résultats soient très considérablement dépendants de la séparation aléatoire des données en ensembles de test et d’apprentissage. Pour évaluer la qualité de ces résultats, on a en fait comparé les résultats de prédiction donnés par deux sommations différentes : l’une utilisant les prédictions brutes de TT et l’autre utilisant ces valeurs de prédiction auxquelles on a fait subir une correction des résidus. Pour la première segmentation, on observe invariablement une réduction relative de la RMSE située aux alentours de 2 à 3%. La réduction relative de la MAE peut quant à elle varier entre $-0,5\%$ et 2% . Dans le même temps, pour la seconde segmentation, il a été déconcertant de

constater que la réduction relative de la RMSE se situait le plus souvent aux alentours de 13% mais pouvait descendre jusqu'à 2%, dans les mêmes conditions de compilation. La réduction relative de la MAE est encore plus variable puisqu'elle atteint régulièrement plus de 15% mais peut parfois descendre à 0%, là aussi dans les mêmes conditions. Nous tâcherons de discuter ces résultats dans la section suivante, et nous mettrons fortement en doute la représentativité des résultats les plus faibles. Quoi qu'il en soit, une amélioration globale est assurément observée avec la régression des résidus.

Concernant la régression linéaire multiple, les résultats sont beaucoup plus stables, et les performances de ce modèle sont bonnes. On compare bien sûr les résultats de prédiction de cette méthode avec ceux obtenus par la sommation des prédictions issues de la régression des résidus, puisque cette dernière est pertinente malgré l'instabilité latente décrite plus haut. On remarque une amélioration conséquente de la RMSE pour la prédiction par segments inter-niveaux, avec une réduction relative de plus de 6%. Pour le sectionnement avec un segment majeur, la réduction relative de cet écart est de l'ordre de 2 à 3%, quelle que soit la qualité des résultats de la régression des résidus. La MAE, en revanche, a tendance à empirer plus ou moins fortement lorsque l'on utilise ce modèle. Elle subit une augmentation relative inférieure à 1% pour le premier sectionnement, et une augmentation relative atteignant généralement 3% lorsque la régression des résidus améliore déjà très fortement les résultats de prédiction. Dans le cas contraire, plus rare, elle subit une augmentation relative inférieure à 1%. Globalement, on juge que la régression linéaire multiple est plutôt pertinente, mais il est clair qu'elle sacrifie une part de la précision obtenue par la régression des résidus pour la MAE afin de favoriser l'amélioration de la RMSE.

Comparons maintenant les prédictions initiales de TT issues du LSTM pour l'itinéraire complet et pour chacune des stratégies de segmentation de l'itinéraire surface→niveau350 (naturellement après l'application des régressions dans ces cas-ci). Les performances de prédiction de TT sur les itinéraires sectionnés en segments inter-niveaux approchent très difficilement les performances de prédiction du LSTM sur l'itinéraire complet. On observe ainsi une réduction relative de la RMSE allant de 0 à plus de 6% et une réduction relative de la MAE allant de 1,5 à 5% lorsqu'on compare les prédictions sur l'itinéraire complet à celles obtenues sur l'itinéraire sectionné en segments inter-niveaux (dans cet ordre). En revanche, grâce aux résultats bien souvent exceptionnels de la régression des résidus pour l'itinéraire sectionné avec un segment majeur, les prédictions de TT basées sur cette segmentation sont généralement très significativement meilleurs que les prédictions sur l'itinéraire au complet. On peut alors observer plus de 10% de réduction relative, à la fois pour la MAE et la RMSE. Dans les cas plus rares où le modèle de régression des résidus échouait à donner d'excellents résultats pour l'itinéraire sectionné avec un segment majeur, les prédictions de TT du sous-

modèle pouvaient aller jusqu'à donner des performances très similaires à celles obtenues via le sectionnement de l'itinéraire par segments inter-niveaux.

Enfin, évaluons l'intérêt de notre modèle d'empilement final. Le constat est sans appel : en tout temps, ce modèle permet une amélioration absolument remarquable de nos résultats. En particulier, le modèle d'empilement démontrent les performances les plus impressionnantes lorsque le modèle de prédiction de TT avec un segment majeur donne lui-même d'excellents résultats : comparativement aux prédictions de ce dernier, qui dépassent déjà significativement en qualité celles de nos autres sous-modèles, le modèle d'empilement peut respectivement permettre une réduction relative de plus de 12% pour la MAE et de plus de 9% pour la RMSE. Dans le cas inverse, il peut permettre une amélioration respective de près de 6% et 5%, respectivement. On représente les histogrammes des TT réels et prédits par le modèle d'empilement, dans le cas le plus courant, à la figure 5.12. Il est intéressant de constater que, malgré la supériorité indéniable des prédictions du modèle d'empilement sur tous les autres modèles, la forme des deux histogrammes coïncide remarquablement peu. L'apparition d'un deuxième mode de faible amplitude, considérablement éloigné de la moyenne, est aussi curieuse.

On regroupe les performances observées lors de l'application de notre modèle d'empilement et de chacun des sous-modèles de prédictions initiales de TT dans le tableau 5.17. Les prédictions finales de notre modèle global de prédiction de TT, qui correspondent à celles du modèle d'empilement, sont très satisfaisantes puisqu'elles surpassent très nettement le modèle de référence de l'industrie minière.

TABLEAU 5.17 Performances finales de nos sous-modèles de prédiction initiale de TT, du modèle d'empilement et du modèle de référence

Modèle	Segmentation	MAE	RMSE
Référence	Itinéraire complet	213	307
LSTM	Itinéraire complet	199	312
	Inter-niveaux	208	312
	Majeure	180	278
Empilement (GBR)	Toutes	158	252

Maintenant que tous les résultats concernant notre cas d'étude ont été rapportés, ajoutons quelques évaluations supplémentaires, avant de conclure sur la composition précise finalement

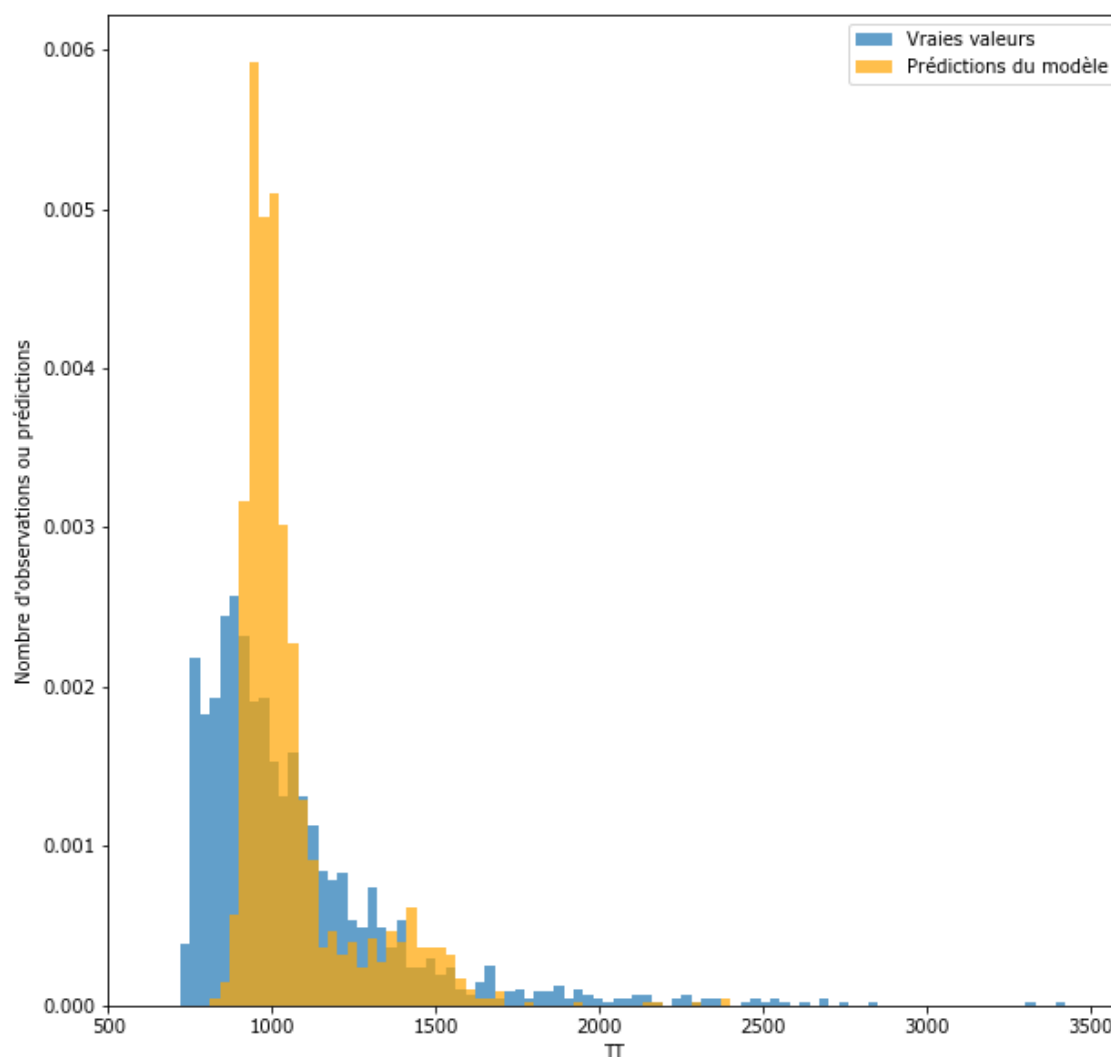


FIGURE 5.12 Histogrammes des TT réels et des TT prédits par notre modèle d'empilement

retenue pour notre modèle de prédiction de TT.

Nous avons eu l'occasion de prédire des TT en montée, uniquement sur l'itinéraire niveau300→surface, à comparer à l'itinéraire inverse. Cette expérimentation a simplement permis de confirmer que, comme nous l'avions supposé, les performances de prédiction de nos modèles de ML sur un itinéraire en montée sont très clairement meilleurs qu'en descente. En effet, les écarts aux prédictions du modèle de référence augmentaient (i.e. amélioration) de plus de 10 points de pourcentage pour la MAE et pour la RMSE sans optimisation préalable. La différence est suffisamment impressionnante pour affirmer que, si notre modèle est capable de prédire convenablement les TT en descente dans une mine souterraine donnée (et que les HT doivent y céder la priorité en descente), il sera capable d'y prédire les TT en mon-

tée plus que convenablement. La distribution des TT en montée sur cet itinéraire, proche d'une simple gaussienne, laissait d'ailleurs penser que les modèles de ML, et en particulier les réseaux de neurones, réussiraient bien à déterminer les relations sous-jacentes entre les conditions opérationnelles fournies et les TT.

Intéressons-nous maintenant à la méthode de filtrage des TT les plus longs, autrement dit quel seuil est employé pour les filtrer. Pour rappel, nous avons éliminé dès le début de notre étude les deux seuils les plus sévères, du fait qu'ils suppriment sûrement une part non négligeable de trajets ordinaires. Pour autant, il est tout de même intéressant de noter que, lors d'expérimentations annexes avec ces seuils, les résultats globaux de notre méthode tendaient à s'améliorer, à comparer à la moyenne de la distribution alors obtenue. Ainsi, la suppression des TT les plus longs paraît effectivement rendre les autres trajets plus prédictibles (bien que l'on sacrifie de nombreux TT ordinaires influents), ce qui corrobore la thèse de l'existence de multiples trajets non ordinaires parmi les TT les plus longs, troublant significativement nos modèles de prédiction.

Concernant la prise en compte de toutes les détections de l'ensemble des HT par chacune des balises, et ce à chaque demi-heure de chaque quart de travail (voir la sous-section 4.5.10 du chapitre 4), notre unique expérimentation a montré une amélioration brutale des indicateurs de performance (après optimisation) pour certains de nos modèles de ML complexes capables de gérer une telle complexité. En particulier, avec une optimisation manuelle sommaire, on a pu observer une impressionnante réduction relative de la MAE du LSTM, qui a atteint près de 20%, tandis que sa RMSE baissait de plus de 9%. La capacité de certains de nos modèles à gérer ces données ne pourra malheureusement pas contribuer directement à l'amélioration des performances concrètes de notre modèle, pour les raisons explicitées dans le chapitre précédent. En revanche, cette expérimentation met remarquablement bien en relief l'impact sur les TT du niveau d'activité des HT sur chaque segment d'un itinéraire. On peut en déduire qu'une excellente évaluation, en amont des quarts de travail, du niveau de congestion dans chaque zone de la mine peut améliorer extraordinairement les performances de notre modèle de prédiction de TT.

Au terme de cette section, en juxtaposant toutes les évaluations précédentes, nous pouvons déterminer une structure optimale de notre modèle global de prédiction de TT, qui le rend capable de résultats de prédiction remarquables pour un itinéraire considérablement long de la mine 1. Ainsi, pour un itinéraire donné, la combinaison optimale de méthodes et de modèles de ML constituant notre modèle global de prédiction de TT est structurée comme suit :

1. Un GMM pour le partitionnement des TT, optimisé par l'AIC ;

2. Un MLP optimisé pour la prédiction des probabilités d'appartenance à chacune des classes issues du partitionnement par le GMM ;
3. Trois sous-modèles de prédiction initiale des TT :
 - * Un LSTM sophistiqué optimisé pour la prédiction sur l'itinéraire complet ;
 - * Un LSTM sophistiqué optimisé sur un segment pour la prédiction sur l'itinéraire sectionné en segments inter-niveaux d'après tous les trajets observés sur chacun de ces derniers, suivi de :
 - Si plus de 100 trajets complets sont disponibles :
 - a. Régression linéaire des résidus (non optimisée) pour ajuster les prédictions par segment du LSTM sur les traversées de segments durant les trajets complets ;
et
 - b. Régression linéaire multiple à fonction de coût personnalisée (adaptée définitivement) pour agréger les prédictions ajustées avec flexibilité.
 - Sinon, sommation des prédictions par segment.
 - * Un LSTM sophistiqué optimisé sur un segment majeur pour la prédiction sur l'itinéraire sectionné avec segment(s) majeur(s) d'après tous les trajets observés sur chacun de ces derniers, suivi de :
 - Si plus de 100 trajets complets sont disponibles :
 - a. Régression linéaire des résidus (non optimisée) pour ajuster les prédictions par segment du LSTM sur les traversées de segments durant les trajets complets ;
et
 - b. Régression linéaire multiple à fonction de coût personnalisée (adaptée définitivement) pour agréger les prédictions ajustées avec flexibilité.
 - Sinon, sommation des prédictions par segment.
4. Un GBR optimisé pour l'empilement final des prédictions des sous-modèles, du tandem GMM-MLP et de la moyenne de la distribution des TT d'entraînement.

5.8 Discussion

La complexité de notre modèle global de prédiction de TT a généré diverses remarques pertinentes, non évoquées directement jusqu'ici, durant sa mise en œuvre. Nous tâcherons de les formuler dans la présente section, en suivant la structure de notre modèle, et nous proposerons également des pistes de recherche prometteuses.

Questionnons-nous d’abord sur l’utilisation directe du GMM pour partitionner les données via les conditions opérationnelles des trajets. Pour rappel, nous avons très rapidement écarté cette utilisation en raison de l’évolution monotone décroissante de l’AIC avec le nombre de classes. Pourtant, elle nous aurait évité de devoir implémenter un MLP et d’optimiser ses hyperparamètres, ajoutant de la complexité à notre modèle et augmentant significativement le temps de calcul nécessaire sur un nouvel itinéraire. Cette utilisation directe avait pourtant été présentée par Fan et al. [18, 21, 22] comme une excellente manière d’utiliser le GMM. Ils avaient ainsi vu augmenter fortement la précision de leurs prédictions finales lorsqu’il était utilisé en amont des autres modèles de ML pour partitionner les trajets directement d’après les données d’entrée. En fait, cette différence majeure provient essentiellement du fait que notre modèle vise à prédire le plus précisément possible le TT à venir sur un itinéraire en particulier (le modèle est entraîné uniquement sur un itinéraire à la fois pour se spécialiser dessus) dans des conditions données, tandis que le leur vise à estimer la productivité horaire pour un itinéraire $A \rightarrow B$ quelconque du réseau minier. De notre côté, nous disposons de données qui peuvent faire varier le TT et qui peuvent permettre à notre modèle très complexe de déterminer les relations latentes entre ces données et le TT, mais nos données ne suffisent pour autant pas à expliquer directement de très grandes variations des TT observables. Au contraire, Fan et al. disposent de trajets en surface, plus réguliers, car nettement moins contraints que les trajets en souterrain. Ils disposent aussi d’informations extrêmement corrélées à la productivité horaire, comme la **distance** de trajet et même la vitesse de déplacement de chaque HT à vide observée sur le cycle en cours d’étude [22], bien que cette dernière donnée ne puisse évidemment pas être connue à l’avance.

Passons maintenant à point de réflexion très significatif pour notre méthode, concernant les prédictions de TT. On a pu observer dans la BDD de la mine 1 que la qualité globale des données de détection des balises tend à augmenter au fil des années. Naturellement, cette tendance nous amène à considérer qu’en ajustant entièrement notre modèle sur les données les plus récentes de la mine 1, il serait possible d’obtenir immédiatement des prédictions de TT plus précises. De surcroît, une autre remarque s’impose : **même sans réajuster** nos différents sous-modèles, ce phénomène devrait nécessairement améliorer la précision de leurs résultats en conditions réelles. En effet, si les données brutes sont progressivement de meilleure qualité, c’est logiquement le cas des données que nous avons préparées. En particulier, elles pourraient contenir une proportion décroissante de trajets non ordinaires. Aussi, ce sont les trajets les plus récents de notre jeu de données, lequel s’étale sur plusieurs années, qui ont certainement dû fournir à nos modèles de ML les données d’entraînement les plus exploitables, influençant ainsi le plus significativement l’ajustement de leurs poids. Par conséquent, lorsque nous avons évalué la qualité des prédictions de tous nos modèles de ML sur l’ensemble des

trajets de test, choisis aléatoirement dans notre jeu de données, nous avons certainement dégradé la précision des prédictions par rapport à celle qui pourrait être obtenue uniquement sur les trajets les plus récents. Par ailleurs, le LSTM est le seul modèle de ML à être exposé simultanément aux caractéristiques de plusieurs trajets, puisqu’il s’appuie sur les données de 10 trajets précédents en plus de celles du trajet en cours. Pour ce modèle, non seulement l’impact des conditions opérationnelles des trajets récents sera plus facilement quantifiable, mais les dynamiques qui lient des trajets rapprochés entre eux seront aussi plus évidentes. Il paraît donc sensé de penser que ses performances s’amélioreront significativement lors de son utilisation en conditions réelles, via cette double amélioration. Nous n’avons pas tâché d’évaluer l’amplitude de cette amélioration, car une problématique apparaît rapidement : nous aurions dû placer une très large part des trajets **récents** dans nos données de **test** pour évaluer convenablement la qualité des prédictions. Le modèle aurait alors été **entraîné** sur une plus grande proportion de données d’apprentissage **anciennes**, qui sont de moins bonne qualité et qui auraient faussé les performances pouvant être obtenues sur les trajets récents.

Concernant la segmentation des itinéraires en courts segments (dans notre cas des segments inter-niveaux) et la baisse de performance observée par rapport à l’itinéraire complet, nos résultats entrent en contradiction avec ceux de la littérature exploitant cette stratégie de segmentation. Pour autant, la littérature parle d’itinéraires de mines à ciel ouvert, et il est donc tout à fait possible que cette approche ne soit pas pertinente dans le contexte des mines souterraines.

D’autre part, le sectionnement d’itinéraires partiellement souterrains en segments majeurs était employé par Li et al. [15], et nos résultats vont dans le même sens que ceux décrits par les auteurs. Ils sectionnaient l’itinéraire étudié en trois segments majeurs, l’un en surface, l’autre en rampe et le dernier en galerie et observaient une belle amélioration de leurs prédictions. De notre côté, notre segment majeur correspond à une rampe simple, et l’association du LSTM et des régressions postérieures semble extrêmement performante pour exploiter cet itinéraire plutôt que l’itinéraire complet, plus complexe. Il nous semble donc indispensable de conserver cette méthode dans toutes les tentatives de prédictions de TT miniers souterrains, en particulier pour analyser spécifiquement chaque milieu de déplacement.

Concernant les résultats de la régression des résidus pour l’itinéraire sectionné avec un segment majeur, nous ne sommes pas capables d’expliquer incontestablement le phénomène à l’origine de l’amplitude des variations observées, mais nous tenons une explication convaincante. En effet, nous avons initialement présumé que ces variations provenaient de performances de prédiction instables du LSTM, en amont, sur le segment majeur surface→niveau300. On suspectait une forme de surapprentissage occasionnel de ce modèle, donc

un phénomène purement négatif, que la régression des résidus corrigerait alors. Pourtant, que cette régression donne d'excellents résultats ou non, les performances du LSTM sont globalement similaires **avant** son utilisation. Même sans cette régression, le LSTM surperforme d'ailleurs largement les autres modèles de ML sur le segment majeur surface→niveau300. Ainsi, lorsque la régression des résidus améliore fortement ses résultats, il semble qu'elle ne comble pas un biais évident du LSTM. On peut d'ailleurs rappeler que quand la régression des résidus a un impact très important sur les performances, ce qui est le cas le plus courant, les prédictions obtenues sont alors bien meilleures que celles des autres méthodes de segmentation. Ajoutons que les très fortes variations des performances suite à la régression peuvent être observées pour deux itérations successives du même programme Python. Nous excluons donc la thèse d'un bogue informatique et favoriserons plutôt celle de l'existence de séparations occasionnellement non représentatives de notre ensemble de données relatives aux trajets sur les itinéraires complets, bien que cet ensemble contienne plus de 5000 observations. Cette dernière thèse sous-entend que ce sont les observations des meilleurs résultats qui sont les plus représentatives des performances véritables de notre régression des résidus sur les itinéraires segmentés avec un segment majeur. Cela correspondrait en effet au cas où ce modèle de régression simple, non optimisé donc très peu robuste, n'a pas été lésé par la séparation des données. Il aura alors pu généraliser correctement à nos trajets de test, soit un peu plus de 1500 observations. La forte amélioration des résultats par la régression des résidus est d'ailleurs l'observation la plus récurrente, ce qui est cohérent avec le fait que les séparations non représentatives sont théoriquement plus rares que l'inverse sur 5000 observations (dont seulement 1500 observations de test et un grand nombre de trajets non ordinaires très influents de par leur TT élevé). Une piste de recherche pertinente pourrait consister à optimiser ce modèle de régression des résidus pour le rendre plus robuste et exploiter son plein potentiel, alors qu'il s'agit déjà d'une des composantes les plus performantes et les plus pertinentes de notre modèle global de prédiction.

Rappelons aussi que notre modèle ne fait que très largement approximer le nombre de HT actifs sur l'itinéraire étudié, et que les planificateurs pourraient lui indiquer bien plus précisément le niveau de congestion attendu sur son itinéraire pour le quart de travail à venir, pour tous véhicules confondus.

Ajoutons quelques pistes de recherche qui nous semblent tout à fait intéressantes à explorer.

La toute première piste est bien simple, elle consiste à ajouter des variables d'entrée pertinentes issues de la télémétrie, notre modèle n'ayant eu accès qu'à des données de détection issues de balises. Ces données ont été maintes fois transformées pour générer de nombreuses autres variables explicatives, mais elles ne nous permettent bien sûr pas tout. Si l'on devait

résumer la liste des variables potentiellement pertinentes mais inutilisées ici, la liste pourrait être impressionnante. Le kilométrage des HT, leur conducteur, le nombre d'autres véhicules actifs sur les itinéraires étudiés, la température et les conditions météorologiques en surface, l'occurrence ou non de ravitaillements en essence au cours d'un trajet et le niveau de chargement du HT en montée sont en particulier des variables influentes que notre modèle n'a jamais eu l'occasion de manipuler. Il est probable qu'il ait en revanche pu inférer certaines d'entre elles indirectement, via les cycles saisonniers et l'identifiant des HT par exemple.

Nous proposons aussi de tester les modèles de ML sur les données les plus récentes de l'échantillon uniquement, sans qu'ils aient été entraînés sur toutes ces dernières, pour comparer les indicateurs de performance. Étant donné le risque de voir apparaître des problématiques liées à l'entraînement sur des données anciennes, il faudrait potentiellement réduire l'échantillon de test pour que les modèles de ML se soient entraînés sur des données récentes de meilleure qualité.

Il serait potentiellement très pertinent de réaliser une analyse en composantes principales de nos variables (en particulier si l'on ajoute d'autres variables de télémétrie) pour fournir les composantes principales à notre modèle. Certaines de nos variables sont d'ailleurs déjà fortement corrélées. Nous pensons en particulier aux quarts de travail (jour/nuit) et au nombre de HT actifs.

Évoquons maintenant une piste de recherche coûteuse en complexité mais potentiellement très intéressante : nous proposons l'utilisation d'un modèle d'empilement dans chacun de nos sous-modèles de prédiction initiale de TT. En effet, les autres modèles de ML étaient parfois capables de surpasser le LSTM, et il y aurait un gain évident de performances pour la prédiction par segments inter-niveaux. L'idée serait donc de combiner les prédictions des six modèles de ML, ou bien plus simplement de combiner celles du LSTM et du XGBoost. Ce dernier est basé sur des arbres de décision, il a donc un fonctionnement radicalement différent de celui des réseaux de neurones. Il était plus régulier que les autres modèles, donnant souvent de bons résultats par rapport au LSTM sur les segments inter-niveaux.

Pour continuer dans cette veine, si une complexité encore supérieure est tolérée, on pourrait même aller jusqu'à utiliser simultanément plusieurs des quatre modèles de ML pour l'empilement proposés initialement, puis à utiliser un méta-modèle d'empilement pour pondérer leurs prédictions de TT respectives. Au vu de l'impressionnante amélioration des résultats qu'a permis le GBR seul, et au vu des performances très similaires de la régression linéaire simple en tant que modèle d'empilement, il y aurait potentiellement une amélioration notable des prédictions à la clé.

Une autre piste de recherche consisterait à optimiser les modèles et paramètres que nous

n'avons pas encore optimisés. Deux d'entre eux sont particulièrement évidents : le modèle de régression des résidus, extrêmement efficace pour améliorer les prédictions du LSTM sur un sectionnement avec un segment majeur, et actuellement non optimisé donc non robuste ; mais aussi la fenêtre de trajets du LSTM, i.e. le nombre de trajets précédents à lui fournir pour guider ses résultats. Précisons au passage que, en conditions réelles, ces trajets précédents pourront être automatiquement détectés dans la BDD de la mine 1 via notre modèle de préparation des données qui pourrait être exécuté automatiquement à la fin de chaque quart de travail pour détecter les trajets ayant eu lieu durant ce dernier.

Toujours à propos de la fenêtre de trajets du LSTM, nous proposons de ne prendre en compte que les quelques trajets précédents du même HT que celui dont on essaie de prédire le TT, pour augmenter l'interprétabilité des trajets précédents.

Il pourrait être pertinent de tester l'optimisation des hyperparamètres du LSTM pour chaque segment inter-niveaux, afin de comparer les performances avec une unique optimisation, pour évaluer la rentabilité du temps de calcul additionnel que cela impliquerait.

Il est par ailleurs possible de sectionner les segments majeurs en segments inter-niveaux, et de prédire les TT sur le segment majeur via un modèle d'empilement : on utiliserait l'itinéraire complet et l'itinéraire sectionné par segments inter-niveaux. Les ressources computationnelles nécessaires et la complexité totale du modèle peuvent évidemment être limitantes.

Notons que notre modèle d'empilement actuel pourrait théoriquement aussi utiliser les conditions opérationnelles pour améliorer ses prédictions. Sans relancer une optimisation via **Optuna**, nous avons rapidement testé cette possibilité et les résultats étaient similaires à ceux obtenus **avec** optimisation et sans ces données. Il y aurait donc une piste d'amélioration ici aussi.

Enfin, si les planificateurs attribuent une plus grande valeur au MAE, il serait possible d'adapter les différentes étapes de notre méthodologie de prédiction pour optimiser ce critère de performance plutôt qu'au RMSE. L'impact sur la valeur de ce dernier dans un tel cas serait intéressante à étudier.

5.9 Conclusion

Le bilan de ce chapitre est particulièrement enthousiasmant.

Tout d'abord, le modèle retenu valide tous les requis spécifiés :

- Grâce à toutes les précautions et solutions annexes décrites tout au long de ce chapitre, nous considérons que notre modèle serait capable de généraliser à tous les itinéraires

de la mine 1, avec une précision certes très variable ;

- Il retourne la MAE et la RMSE obtenues sur l'ensemble de test aux utilisateurs ;
- Grâce à sa simplification partielle, il est optimisé et entraîné en moins de trois heures avec notre support informatique ;
- Enfin, une fois entraîné, il retourne plus d'un millier de prédictions de TT en quelques dizaines de secondes, nous sommes donc très loin d'invalider le requis d'une seconde par prédiction.

Pour continuer, les performances finales du modèle que nous avons développé sont évidemment satisfaisantes, puisqu'elles surpassent très largement celles du modèle de référence de l'industrie minière. Rappelons aussi que l'itinéraire étudié est représentatif d'un itinéraire sur lequel notre modèle pourrait être concrètement appliqué dans la mine 1.

De plus, chaque élément de ce modèle est maintenant pleinement justifié, puisque chacun d'entre eux a contribué à faire descendre toujours plus les deux indicateurs de performance, tout en faisant grimper la robustesse de notre modèle global. Pour mettre en relief l'intérêt de la haute complexité de notre modèle, rappelons qu'un simple LSTM serait bien incapable, sans l'aide du tandem GMM-MLP, d'atteindre les performances du modèle de référence sur l'itinéraire étudié, et ce même si on optimisait en amont ses hyperparamètres pendant près d'une heure. Pour notre tâche particulière de prédiction, et en fournissant à tous nos modèles les prédictions du tandem GMM-MLP, nous avons d'ailleurs pu montrer la supériorité globale du LSTM sur les cinq autres modèles de ML qui avaient montré des résultats appréciables dans la littérature (malgré le fait que notre BRNN n'en soit pas un à strictement parler). Nous avons aussi pu affirmer que la plupart de ces modèles ont tout de même été de très bons compétiteurs, qui pourraient assister le LSTM pour les prédictions de TT via un modèle d'empilement.

Au vu de l'immense inventaire de difficultés additionnelles rencontrées successivement par notre modèle de préparation des données puis par notre modèle de prédiction de TT pour la mine 1, il est possible que la partie la plus positive de notre bilan soit la suivante : le modèle que nous avons développé n'est absolument pas arrivé au summum de ses performances pour la mine 1, et est très prometteur pour d'autres sites miniers. En particulier, rappelons certains handicaps exceptionnels auxquels nos deux modèles principaux ont pu être exposés :

- En dehors des lots de données de détection des balises (et du plan de la mine qui s'y rapporte), nous n'avons utilisé **aucune** autre donnée de prédiction. L'ajout des variables évoquées dans la section précédente aux côtés des variables que nous utilisons déjà pourrait apporter un effet d'amélioration quasi-immédiat.
- L'identification des trajets non ordinaires, sans sacrifier de trajets ordinaires, a été si

difficile que nous avons fini par autoriser les HT à prendre jusqu'à une demi-heure pour passer d'une balise de changement de niveau à l'autre. On suppose que le nombre de trajets non ordinaires qui en résulte, avec leurs TT extrêmement longs, est particulièrement gênant pour tous les modèles de ML et que la qualité de nos prédictions en est nécessairement dégradée. Il nous est impossible de dire à quel point. On peut en revanche affirmer qu'il est possible de mieux discriminer ces types de trajets, via des seuils de filtrage choisis par des experts du site, via un meilleur réglage des balises de détection et via d'autres variables de télémétrie.

- La nature particulière de l'itinéraire de prédiction, qui passe au sein même du niveau souterrain qui est naturellement le plus actif toute période confondue, cause des problématiques au niveau de la préparation des données. De très nombreuses balises avoisinent ce trajet et sont capables de détecter le HT à très longue distance dans la mine 1. Bien évidemment, cela crée aussi d'importantes problématiques de prédiction des TT puisque le HT a des chances particulièrement élevées d'y faire une pause ou un petit détour. Nous n'avons aucun moyen de détecter de tels événements. En effet, une pause de moins de trente minutes ne doit pas être exclue pour s'assurer de ne pas surévaluer les performances de notre modèle de prédiction. Pour les détours, nous sommes obligés d'être extrêmement flexibles sur notre choix des balises interdites sur ce niveau pour ne pas réduire à néant l'échantillon des trajets observés (du fait des détections à grande distance des HT par les balises interdites).

Au vu de ces handicaps, rappelons l'objectif que nous nous étions fixés dans la section 5.3 : « Au final, dans le cas où les performances de notre modèle se révéleraient tout de même acceptables, on pourra être certain que son architecture est indubitablement pertinente. ». On constate ainsi mieux à quel point notre modèle a surperformé en tout point.

Au vu de tous ces constats et pour continuer sur cette lancée, nous proposons un plan d'action en trois étapes destiné à améliorer encore les performances de notre modèle dans la mine 1 :

1. Mettre d'abord en place les propositions d'amélioration du chapitre 4 pour obtenir des données initiales de meilleure qualité et prélever d'autres variables pertinentes existantes dans la BDD pour les fournir au modèle ;
2. Préparer minutieusement ces données récentes et consacrer des efforts importants à la discrimination des trajets ordinaires et des trajets non ordinaires, en particulier en fixant des seuils de filtrage des TT aberrants adéquatement sélectionnés par des experts du site ; et
3. Ajuster notre modèle de prédiction de TT sur les données récentes qui auront été bonifiées, pour les itinéraires qui intéressent réellement les planificateurs, en optimisant

ses hyperparamètres.

Enfin, nous allons dès maintenant tâcher de transformer l'essai, en testant la généralisation de notre méthodologie complète à un autre site minier souterrain disposant d'un réseau de balises de détection : la mine 2.

CHAPITRE 6 TEST DE GÉNÉRALISATION DU MODÈLE GLOBAL

Le présent chapitre vise à tester l'application de notre modèle global à un autre site minier, pour évaluer ses capacités de généralisation. Comme le chapitre 4, il débutera par une description du nouveau cas d'étude, présentée à la section 6.1 et d'une analyse des contraintes et des spécificités du cas d'étude, à la section 6.2. Nous consacrerons alors la section 6.4 à la phase de préparation de données et présenterons ensuite la phase de prédiction de TT à la section 6.5. Enfin, nous concluons.

6.1 Description du cas d'étude

La présente section propose une rapide présentation générale de la mine 2.

Tout comme la mine 1, la mine 2 est une mine souterraine canadienne aurifère à accès par rampe, exploitée via la méthode d'abattage par longs trous. Son architecture est plus simple que celle de la mine 1 puisqu'elle ne possède qu'une seule rampe d'accès. Cette dernière relie la zone de déchargement du minerai en surface, au niveau 46 à 460 mètres sous terre, en desservant tous les niveaux intermédiaires (à l'exception du niveau 31, accessible uniquement depuis le niveau 34, mais nous n'aurons pas besoin d'y prêter une plus ample attention). L'écart inter-niveaux n'est pas constant puisqu'il passe de 30 à 40 mètres à partir du niveau 30. Le numéro du niveau indique la profondeur correspondante divisée par 10. On trouve une baie de lavage au niveau 6 et plusieurs refuges le long de la rampe. On trouve aussi des chambres creusées pour y permettre les croisements de véhicules.

Concernant les activités de transport de minerai, un total de 7 HT ont opéré à un moment donné sur ce site sur la période étudiée. Leurs cycles de transport de minerai comportent les mêmes six étapes classiques que ceux de la mine 1. Globalement, les mêmes imprévisibilités peuvent peser sur les trajets des HT.

6.2 Contraintes et spécificités du cas d'étude

Dans cette section, attelons-nous à décrire les particularités de ce second site minier et à identifier les contraintes qui s'y rattachent.

Évoquons tout d'abord le fait que la nature souterraine de ce site minier a des impacts très importants sur les trajets de HT et leur étude, comme pour la mine 1.

Heureusement, concernant l'étude des TT, nous pourrions ici aussi profiter de l'existence d'un

réseau de connectivité souterraine qui inclut cette fois-ci au moins 240 balises de détection fonctionnelles. Pour rappel, la mine 1 ne disposait que de 135 balises malgré ses deux rampes d'accès et un plus grand nombre de niveaux. Les balises de la mine 2 semblent réparties à intervalles réguliers dans la rampe et les niveaux souterrains et sont globalement absentes de la surface. Ces balises sont supposées être bien différentes de celles de la mine 1 puisque ce sont des entreprises différentes qui ont conçu, installé et paramétré les réseaux de balises des deux sites miniers. Naturellement, les HT de la mine 2 sont aussi équipés de puces RFID pour permettre leur reconnaissance. Par la suite, nous tâcherons bien sûr de mettre en lumière les forces et les faiblesses du réseau de balises de détection de la mine 2 relativement au contexte de préparation de données, comparativement au réseau de la mine 1.

La spécificité la plus intéressante de la mine 2 n'a toutefois pas encore été présentée. En effet, cette mine a une particularité exceptionnelle : elle combine des quarts de travail durant lesquels ce sont des conducteurs humains qui pilotent les HT (quarts/trajets dits « conventionnels ») et des périodes additionnelles durant lesquelles un système de conduite autonome de HT prend le relais (quarts/trajets dits « autonomes »). Les quarts de travail conventionnels sont au nombre de 12 par semaine : ils ont lieu à chaque journée, et le soir (et la nuit) de chaque jour de semaine. Leur durée est de huit heures. Deux quarts de travail autonomes complets ont lieu par semaine, l'un dans la nuit de samedi à dimanche et l'autre dans la nuit de dimanche à lundi. Leur durée est d'environ 14 heures. Par ailleurs, des trajets autonomes de HT ont lieu entre les quarts de travail conventionnels pour conserver une certaine productivité durant ces heures creuses. Ajoutons que, d'une part, les conducteurs des HT sont une fois de plus rémunérés à la tonne de minerai extraite, ce qui minimise la quantité potentielle de pauses et de détours durant les trajets conventionnels. D'autre part, le système de conduite autonome permet d'éviter la totalité des pauses typiquement humaines (restauration et autres besoins essentiels) et il n'a pas le libre-arbitre de choisir de faire un détour non ordinaire. Enfin, précisons que le système de conduite autonome est en mesure de gérer les croisements et les conflits de trajectoires en utilisant les chambres creusées dans les parois de la rampe. En lien avec notre objectif de prédiction de TT, nous tâcherons de nous intéresser à la fois aux trajets conventionnels et autonomes pour établir des comparaisons édifiantes et apparemment inédites dans la littérature scientifique entre ces deux modes de pilotage de HT dans les mines souterraines.

6.3 Attributs correspondants aux variables d'intérêt

Dans cette section, nous décrirons précisément les attributs d'intérêt de la BDD de la mine 2 qui devraient nous permettre de retrouver les variables d'intérêt présentées au chapitre 4.

Ces dernières restent naturellement identiques.

Notre exploration de la BDD de la mine 2 nous mène à retenir presque les mêmes attributs que pour la mine 1 issus des lots de données de détection des balises. Nous retiendrons donc ici l'horodatage, l'identifiant du HT détecté et, **au sein d'un seul et même attribut**, l'identifiant de la balise à l'origine de la détection **ou** le niveau nouvellement atteint. Concernant ce dernier attribut, il fonctionne exactement comme l'attribut donnant l'identifiant des balises dans la mine 1 mais, lorsque l'identifiant du niveau associé à la balise nouvellement atteinte est différent de celui de celle qui précède, un lot de données de détection mentionnant l'identifiant du niveau nouvellement atteint est émise en même temps que le lot classique de données de détection mentionnant la balise nouvellement rencontrée. Plus précisément, la détection du nouveau niveau est enregistrée immédiatement avant celle de la nouvelle balise. Nous n'utiliserons aucune autre variable, ce qui permettra une comparaison plus directe entre les performances de notre modèle de prédiction de TT sur les deux sites. Les attributs sélectionnés, leur format et les variables d'intérêt à extraire sont présentés dans la tableau 6.1.

TABLEAU 6.1 Caractérisation des attributs sélectionnés dans la BDD de la mine 2

Source	Attribut	Format	Variables d'intérêt à extraire
Détections des balises	Horodatage des détections	<i>DD/MM/YYYY hh:mm:ss</i>	TT Points de départ et d'arrivée Sens vertical du trajet Type de HT (autonome ou non) Identifiant du HT Détours ordinaires
	Identifiant des balises rencontrées OU niveau	Trois chiffres <u>OU</u> <i>Mine2_0NN</i> avec <i>NN</i> le numéro du niveau (<i>00</i> à <i>46</i>)	Pauses de longue durée Nombre de HT actifs Quart de travail Jour de la semaine Saison
	Identifiant du HT détecté	<i>XXXXXn</i> avec $n \in \{1; 7\}$	Moment dans le quart de travail Moment dans le jeu de données Écart temporel entre les trajets

Enfin, il est indispensable de préciser que nous disposons d'un plan 3D de la mine 2 avec la localisation de chacune des balises de détection.

6.4 Application de notre modèle de préparation de données

Dans la présente section, nous nous intéresserons à la préparation des données extraites de la BDD de la mine 2. Cette préparation s'annonce déjà notablement différente de celle de la mine 1, avec de multiples adaptations des étapes constitutives de notre méthodologie.

6.4.1 Extraction du jeu de données d'intérêt et remarques générales

Pour commencer, nous extrayons l'ensemble des données associées aux trois attributs sélectionnés via une requête écrite en langage SQL, exactement comme nous l'avons décrit pour la mine 1. Au total, plus de deux millions et demi de lots de données de détection sont extraits sur la période s'étalant de juillet 2021 à mars 2023.

Nous remarquons immédiatement qu'une part considérable de ces détections ne fournit aucune information utile. En particulier, sur certaines périodes, certaines balises détectent en boucle le même HT en générant aussi un lot de données de détection de changement de niveau alors que ce HT ne semble pas s'être déplacé. Par ailleurs, certaines détections de HT ne mentionnent pas d'identifiant de balise : elles affichent parfois *None*, *No Data* ou aucun caractère. On supprime toutes ces anomalies, qui représentaient plusieurs centaines de milliers de lignes de données.

On remarque aussi immédiatement que certains HT semblent occuper une part disproportionnée de nos données de détection. D'autres semblent au contraire très peu détectés. On pousse l'investigation en comptabilisant les détections de chacun des HT (on ne compte pas les nombreuses détections correspondant à des changements de niveau, qui augmenteraient artificiellement le nombre de véritables détections de balises). On borne aussi la période de détection de chacun des HT, en comptant le nombre de mois durant lesquels ils ont été actifs. Sur certaines longues périodes incluses entre les bornes, certains HT n'ont jamais été détectés : c'est le cas du HT2 durant quatre mois, du HT3 durant deux mois et du HT4 durant 11 mois successifs. Les résultats, présentés dans le tableau 6.2 montrent un déséquilibre très important entre le nombre de détections respectif de chacun des HT. Cela pourrait handicaper notre modèle de prédiction de TT puisqu'il prend en compte l'identifiant du HT lors de ses prédictions. Le nombre de véritables détections de HT par les balises est finalement proche de deux millions. À titre de comparaison, on avait compté près de trois millions de détections dans la mine 1, mais les conflits de détection entre les balises multipliaient artificiellement le véritable nombre de détections sans qu'il soit possible de déterminer ce dernier. Ici, on ne reconnaît pas un tel phénomène dans les séquences de détection, les successions de détection semblent cohérentes avec l'ordre réel des balises.

TABLEAU 6.2 Bornes de la période d'activité et nombre de détections de chacun des HT

Id. anonymisé du HT	HT1	HT2	HT3	HT4	HT5	HT6	HT7
Bornes période d'acti.	10/21	09/21	07/21	08/21	08/21	06/22	08/21
	02/22	03/23	03/23	03/23	03/23	03/23	03/23
Durée d'acti. (mois)	5	15	18,5	7,5	20	10	20
Nb. détect. (milliers)	111	314	315	66	629	141	443

6.4.2 Identification du niveau de profondeur de chaque balise

Concernant l'identification du niveau de profondeur de chaque balise, il est clair que peu de manipulations sont nécessaires car la coexistence de deux « sous-attributs » dans le même attribut (identifiant de la balise rencontrée et identifiant du niveau nouvellement atteint) nous permet d'associer très rapidement chaque balise à sa profondeur respective. En effet, l'identifiant de chaque niveau contient le numéro du niveau, lequel doit être multiplié par 10 pour obtenir la profondeur correspondante. Dès lors, en lisant la séquence des détections, il suffit de calculer à chaque fois la profondeur du niveau nouvellement atteint par le HT et d'associer cette profondeur à toutes les balises suivantes dans la séquence, jusqu'à ce qu'un autre niveau soit atteint par le HT. Il est donc très simple d'adapter notre méthodologie à la mine 2 sur ce point.

En réalité, ces quelques manipulations ne sont pas nécessaires pour ce site minier. D'une part, on peut directement reconnaître les détours inter-niveaux en lisant les détections correspondant à un niveau nouvellement atteint. D'autre part, notre méthodologie recommande de se servir uniquement du plan de la mine pour sélectionner pertinemment chaque balise de changement de niveau. Nous privilégions donc ces raccourcis bien visibles par la suite mais nous avons pu montrer que notre méthodologie aurait pu être adaptée ici aussi.

6.4.3 Recherche d'incohérences entre les périodes de collecte de données

Tâchons maintenant d'utiliser à bon escient les enseignements de la sous-section 4.5.3.

Pour rappel, cette dernière consistait à mettre en place différentes stratégies de compréhension pour mettre en lumière certaines incohérences dans la période temporelle étendue du jeu de données étudié afin de les corriger au plus tôt. On s'était en particulier concentré sur l'incohérence visible du surnombre de balises dans la BDD de la mine 1, qui était due au renommage de balises et à la conservation des identifiants obsolètes dans les données.

Concernant la mine 2, nous n’identifions aucune problématique similaire avec les mêmes techniques.

Ainsi, notre intérêt va plutôt se concentrer sur l’évocation rapide que nous avons faite d’autres stratégies de compréhension via des représentations graphiques de certaines variables d’intérêt, à la recherche d’incohérences potentielles supplémentaires. Bien que cette stratégie n’avait pas permis d’identifier des problématiques additionnelles, c’est un point crucial de notre méthodologie pour pré-filtrer les problématiques les plus gênantes ensuite. Dans le présent chapitre, nous avons déjà identifié une problématique liée au déséquilibre du nombre de détection de chaque identifiant de HT et nous avons par ailleurs vérifié la cohérence temporelle des identifiants de balises pour la mine 2. Il nous reste donc à nous intéresser à la profondeur des HT. Comme dans le chapitre 4, on trace alors le graphe d’évolution de la profondeur de chacun des HT en fonction du temps sur toute la période étudiée. On ne remarque rien de suspect pour le HT1, comme l’indique la figure 6.1.

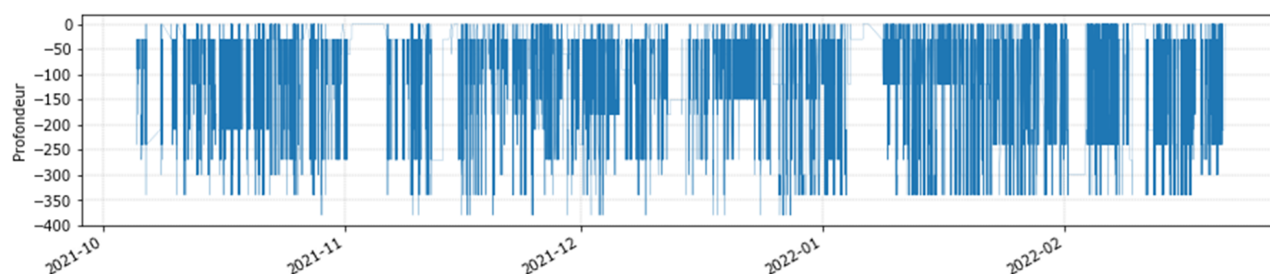


FIGURE 6.1 Tracé continu des profondeurs successives occupées par le HT1 de la mine 2 sur une période de cinq mois

Toutefois, certaines observations de phénomènes indésirables dans notre jeu de données laissent penser qu’une problématique importante reste à éclaircir. En effet, comme le montre le tableau 6.3, on remarque des séquences de détection imprévisibles durant lesquelles certains niveaux n’apparaissent pas du tout. Rappelons que le HT devrait nécessairement y être détecté puisqu’il n’y a qu’une seule rampe et que celle-ci est largement fournie en balises de détection, qui semblent parfois toutes détecter successivement le HT avec une excellente fiabilité. De pures coïncidences ne permettent pas de justifier convenablement nos observations de phénomènes indésirables puisque les HT semblent capables ici de régulièrement parcourir plusieurs niveaux sans jamais passer dans le périmètre de détection de dizaines de balises successives. Ces simples observations, notablement irrégulières, ne nous permettent pas d’estimer la proportion de données contaminées par ce phénomène. En revanche, elles nous font craindre l’existence d’une problématique obscure considérable dans l’acquisition ou la gestion des données de la mine 2.

TABLEAU 6.3 Reproduction simplifiée des séquences indésirables de données de détection observées

Horodatage de la détection	Identifiant de la balise ou du niveau	Identifiant du HT
30/12/2021 3 :22	Mine2__000	HT1
30/12/2021 3 :22	471	HT1
30/12/2021 3 :22	Mine2__003	HT1
30/12/2021 3 :22	234	HT1
30/12/2021 3 :22	235	HT1
30/12/2021 3 :22	236	HT1
30/12/2021 3 :23	435	HT1
30/12/2021 5 :28	Mine2__015	HT1
30/12/2021 5 :28	259	HT1
30/12/2021 5 :29	260	HT1
30/12/2021 7 :24	Mine2__030	HT1
30/12/2021 7 :24	457	HT1
30/12/2021 7 :38	289	HT1
30/12/2021 7 :39	457	HT1
30/12/2021 7 :45	Mine2__015	HT1
30/12/2021 7 :45	260	HT1
30/12/2021 7 :46	259	HT1
30/12/2021 7 :46	Mine2__012	HT1
30/12/2021 7 :46	258	HT1
30/12/2021 7 :47	257	HT1
30/12/2021 7 :47	250	HT1
30/12/2021 7 :47	249	HT1
30/12/2021 7 :51	Mine2__003	HT1

Pour estimer l'impact de l'ensemble de ces phénomènes, nous traçons dans un premier temps le graphe d'évolution de la profondeur de chacun des HT sur toute la période étudiée mais cette fois-ci **sans lier les points de profondeur** entre eux. Quel que soit le HT étudié, la figure obtenue nous permet alors toujours d'effectuer les mêmes observations. On décide de représenter ici les profondeurs ponctuelles du HT1 sur la figure 6.2, pour la même période temporelle de détection que celle de la figure 6.1. Rappelons que cette dernière figure nous avait semblé parfaitement cohérente, tandis que nous observons ici une problématique évidente d'absence de données de détection sur plusieurs niveaux, et ce durant parfois plus d'un mois.



FIGURE 6.2 Graphe des profondeurs ponctuelles successives occupées par l'un des HT de la mine 2 sur une période de cinq mois

Bien que ce graphe mette clairement en lumière l'existence de cette problématique considérable, il ne semble pas avoir livré toute l'étendue des incohérences observées dans les séquences de détection. En effet, alors que la totalité des niveaux disposent de trois à cinq balises dans la rampe, il semble souvent arriver qu'une seule de ces balises ne génère des détections. Lorsque l'on s'intéresse plus spécifiquement à cette balise, on remarque généralement qu'elle n'est pas infallible pour autant puisqu'elle cesse parfois à son tour de se manifester pendant plusieurs mois. La problématique semble donc plus vaste encore que ce que laisse penser la figure 6.2.

Nous avons alors tâché de trouver un mode de représentation des absences de données de balises intempestives dans la mine 2. Représenter cet ensemble de phénomènes de durée variable sur une longue période, pour toutes les balises et en une seule figure facilement interprétable sans avoir manipulé directement les données est un véritable défi. Notre représentation la plus aboutie, explicitée ci-après, peut être visualisée à la figure 6.3. Au contraire des graphes précédents, la profondeur est en abscisse et le temps en ordonnée. Chaque ligne de la matrice correspond à un quart de travail au complet. Il y en a ici 612 (soit près d'un an de détections) triés dans l'ordre chronologique, de haut en bas. Chaque colonne correspond à une certaine balise de la mine 2. Il y en a ici 211, triées en ordre de profondeur croissante. Chaque case de la matrice a été colorée en gris foncé si au moins une détection de HT a été enregistrée pour la

balise correspondante durant le quart de travail correspondant, et laissée en blanc sinon. On a aussi affiché les séparations entre les différents niveaux de la mine 2 par des lignes segmentées bleues verticales. Entre deux de ces lignes, on trouve donc à la fois les balises de la rampe et celles des galeries (bien plus rarement détectées). Ce dernier point permet de grandement faciliter l'interprétation, mais aussi de rendre compte de l'absence parfois quasi-totale de détections de HT sur une part très conséquente des niveaux. En effet, sur la seconde moitié des quarts, les HT sont détectés par quelques balises très en profondeur sur la majorité des quarts de travail (niveaux 38, 42 ou 46 en particulier) et par les balises les plus proches de la surface (les balises du niveau 3 semblent totalement épargnées par le phénomène), souvent sans être détectés sur près de la moitié des niveaux. De manière générale, cette visualisation montre le caractère artificiel et imprévisible de ce phénomène, avec des blocs complets de détections de balises pouvant disparaître brutalement pendant plusieurs mois, puis réapparaître du jour au lendemain. Cela rend aussi très bien compte de l'impossibilité de détecter systématiquement les détours des HT avec notre algorithme puisque, au cours d'un trajet donné, il n'est souvent même pas possible de dire à quel niveau se trouve le HT. On peut aussi remarquer qu'il est parfois difficile de choisir la meilleure balise de changement de niveau. Par exemple, au niveau 30, la cinquième balise est extrêmement fiable sur la première moitié des quarts (près de six mois) par rapport aux quatre premières, mais aucune nouvelle détection de cette balise n'a été enregistrée dans la BDD sur les derniers mois. On peut d'ailleurs remarquer qu'aucune nouvelle balise du niveau 30 ne l'a remplacée. Elle figure toujours sur le plan 3D, sans qu'aucune autre balise n'ait pris sa position. Les quatre premières balises sont quant à elles devenues les plus fiables du niveau. Au niveau 18, le phénomène est encore bien plus visible, avec plusieurs alternances évidentes des balises fiables.

Nous avons d'abord présenté ces résultats à l'administrateur de la BDD opérationnelle de notre partenaire industriel mais nous n'avons pas pu aboutir à une explication parfaitement convaincante bien que nous supposions à ce stade une désactivation temporaire de lots de balises voisines. Nous sommes aussi allés jusqu'à contacter le PDG de l'entreprise qui installe ces balises, qui n'avait alors jamais entendu parler d'un tel phénomène. Son ingénieur logiciel principal a pu nous confirmer qu'aucune des balises installées dans la mine 2 n'avait été victime d'un tel dysfonctionnement dans le dernier mois (le seul auquel l'entreprise peut directement accéder) qui était alors le mois d'août 2023. Pourtant, il semble d'après notre graphe qu'aucune période précédente n'ait été totalement épargnée par les absences de données de détection de balises.

Au terme de ces investigations, nous proposons les raisonnements logiques suivants, dont la vraisemblance est variable, mais sur lesquels on ne peut pas définitivement se prononcer.

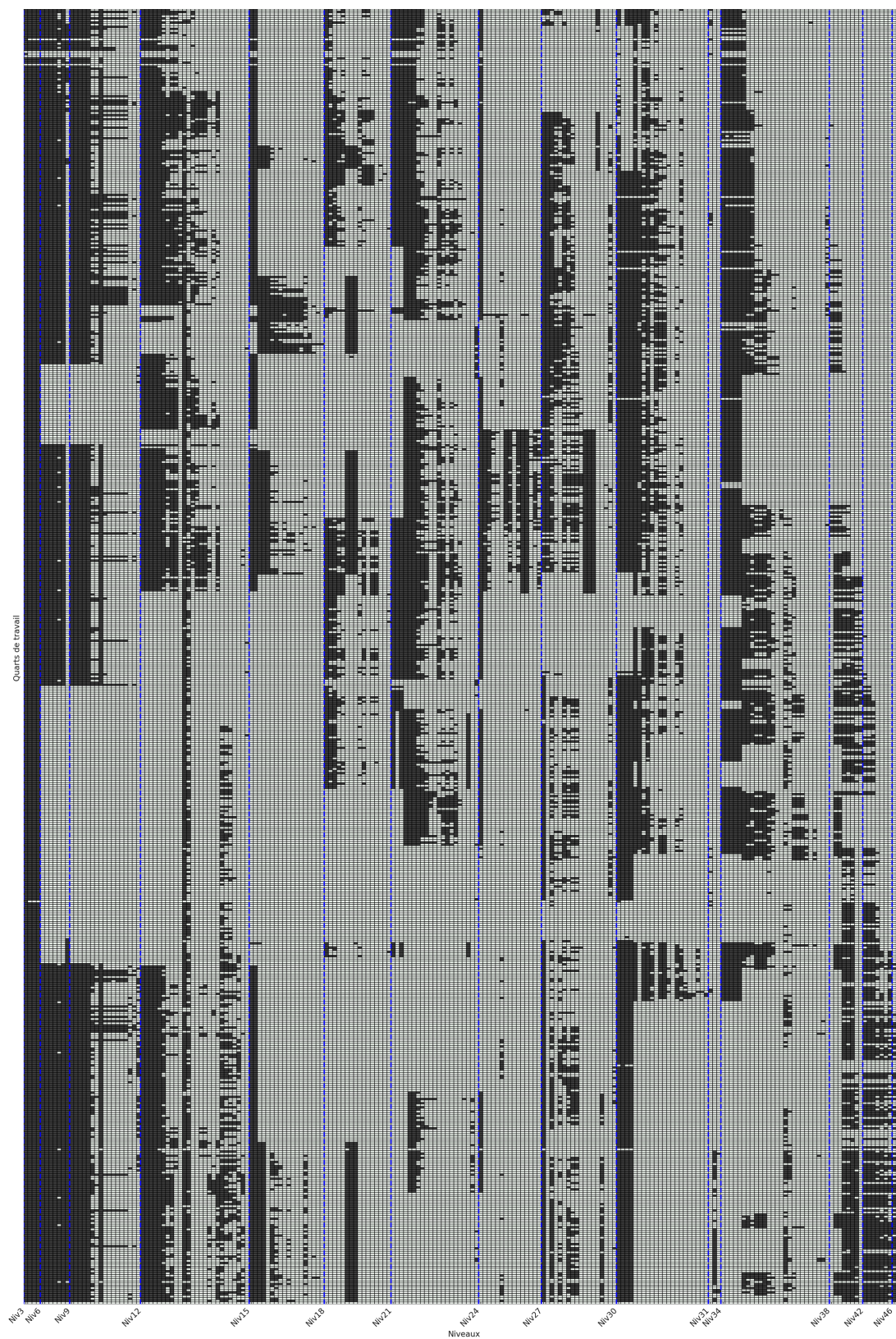


FIGURE 6.3 Visualisation matricielle des absences de données de balises dans la mine 2

Scénario 1 - Dysfonctionnement du réseau de balises puis résolution obscure :

- Nous savons que les balises sont branchées en circuit mixte et qu'il y a donc des balises branchées en série ;
- L'environnement minier souterrain est très contraint donc, selon certains industriels du secteur, les câbles reliant les balises pourraient être régulièrement débranchés accidentellement par les véhicules qui y circulent sans que personne ne le remarque (ou n'y prête attention) pendant quelques semaines ou quelques mois ;
- Par ailleurs, diverses maintenances peuvent nécessiter de débrancher des câbles de balises, sans assurance que ceux-ci soient ensuite rebranchés pour permettre de nouvelles détections ;
- Les observations de disparitions de blocs complets de balises seraient donc totalement justifiées par les affirmations précédentes ;
- Après avoir été victimes de ces dysfonctionnements hasardeux pendant près de deux ans a minima, toutes les balises de la mine 2 ont été rebranchées systématiquement et n'ont subitement plus rencontré une seule des problématiques précédentes seulement quelques mois après la période temporelle que nous avons prélevé ; et
- Les personnes les plus qualifiées des deux entreprises concernées (notre partenaire industriel et l'entreprise ayant implanté les balises) n'ont jamais été informés de cette bonne nouvelle.

Scénario 2 - Problématique de transfert des données post-acquisition :

- Les balises génèrent un volume impressionnant de données lors des quarts de travail ;
- Toutes les autres données opérationnelles (de télémétrie) sont collectées dans le même temps ;
- Ces dernières sont nettement plus utilisées par le personnel de la mine 2 car plus anciennes donc bien mieux maîtrisées ;
- Elles sont critiques pour continuer à disposer des indicateurs existants ;
- En cas de saturation de la bande passante dédiée au transfert de données opérationnelles dans la BDD de la mine 2, elles seraient potentiellement affectées d'une priorité de transfert absolue sur les données issues des détections de balises ;
- Pour une raison obscure, au lieu de désactiver temporairement l'enregistrement de données correspondant aux attributs des balises, seules certaines balises géographiquement groupées seraient désactivées temporairement (mais pas nécessairement toutes celles d'un même niveau) ;

- Lors d’une ré-initialisation du système de transfert de données, la plupart des exclusions de balises pourraient être elles-mêmes réinitialisées collectivement ; et
- L’entreprise ayant implanté les balises aurait uniquement accès aux données brutes située en amont de la BDD opérationnelle de la mine 2 dans un espace de stockage directement relié aux balises, lequel n’est aucunement victime des disparitions de données liées à la saturation de la bande passante de transfert vers la BDD opérationnelle de la mine 2, ce qui expliquerait que l’entreprise n’ait jamais eu connaissance de ces dysfonctionnements.

Quoi qu’il en soit, en fin de compte, les enseignements tirés de la sous-section 4.5.3 se sont révélés cruciaux. Ils nous ont fait adopter ici une démarche proactive pour véritablement approfondir notre analyse des quelques incohérences initialement observées, ce qui a facilité la détection d’une anomalie particulièrement problématique dans la base de données de la mine 2. En l’absence de cette démarche, nous aurions potentiellement laissé passer les quelques anomalies initiales, en espérant qu’elles seraient noyées dans la masse et filtrées correctement par notre algorithme d’identification de trajets. Ce dernier aurait pourtant été bien incapable de trouver suffisamment de trajets valides dans de telles circonstances et nos pistes de compréhension du phénomène auraient été bien maigres. En effet, nous aurions dû revérifier si chacune des étapes précédentes de notre méthodologie est à l’origine du phénomène alors que ce dernier est déjà extrêmement difficile à cerner à ce stade sans la vision globale offerte par la visualisation matricielle présentée ci-avant. Précisons enfin que le phénomène affecte à la fois les quarts de travail conventionnels et autonomes, toutes les données ayant été utilisées sans distinction. Il va sans dire que les conséquences de ce phénomène sur les prochaines étapes de notre méthodologie vont assurément être considérables.

6.4.4 Identification de l’itinéraire prédominant à étudier

La problématique que nous venons d’évoquer pourrait rendre sous-optimale la méthode que nous avons proposée pour identifier l’itinéraire prédominant à étudier. En effet, si l’on considère la représentation matricielle précédente, on constate que le niveau 30 n’aurait pas paru intéressant avec notre méthode présentée au chapitre 4 puisque les niveaux 34, 38, 42 et 46 auraient quasiment toujours enregistré la détection la plus profonde du HT. On aurait alors retenu l’itinéraire Niv3→Niv34 comme étant l’itinéraire prédominant le plus long et donc le plus pertinent de la mine 2. Pour autant, bien que ce choix n’aurait pas été catastrophique, il reste moins pertinent que l’itinéraire Niv3→Niv30. En effet, les quatre premières balises du niveau 30 sont visiblement plus souvent fonctionnelles que les balises les plus fiables du niveau 34. Elles étaient par ailleurs visiblement fiables durant les derniers mois de la période

temporelle étudiée, auxquels on accorde plus de valeur que les toutes premières semaines de la période étudiée, bien couvertes par les balises du niveau 34. L’itinéraire Niv3→Niv30 sera finalement notre itinéraire d’étude pour tout le reste du présent chapitre, au vu de la fiabilité des balises situées sur la rampe du niveau 30 par rapport aux autres niveaux.

Ainsi, grâce à la visualisation matricielle que les incohérences précédentes nous ont poussé à concevoir spécialement pour la mine 2, il a été possible de choisir avec une pertinence accrue son itinéraire prédominant à étudier. Nous avons par ailleurs ici mis en évidence un léger défaut de cette étape de notre méthodologie, qui reste néanmoins fonctionnelle. Elle nous avait pourtant parue optimale dans le chapitre 4 pour la sélection de l’itinéraire prédominant à étudier lorsque la très grande majorité des balises fonctionnaient correctement sur toute la période étudiée.

6.4.5 Listage des balises de changement de niveau

Intéressons-nous maintenant aux balises de changement de niveau, qui sont ensuite omniprésentes dans notre méthodologie globale. Évidemment, la problématique d’absence intempes- tive de données de balises va bousculer là aussi quelque peu notre méthodologie.

Dans la mine 2, il n’y a de toute façon pas d’identifiant récurrent permettant de reconnaître les balises de changement de niveau avec certitude. Notre méthodologie nous demanderait donc d’essayer de trouver une autre forme de récurrence dans les identifiants de balises ou, à défaut, d’utiliser le plan du site minier pour sélectionner manuellement et individuellement chaque balise de changement de niveau. Ici, nous devons en plus tenir compte du fait que toutes les balises de la rampe ne sont pas aussi fiables, pour un niveau donné. Par ailleurs, nous remarquons que le niveau associé à chaque balise de la rampe correspond au niveau qui la suit immédiatement en descente. Ainsi, pour détecter les changements de niveau en descente, il faut d’abord s’intéresser à la balise de la rampe située la plus en aval de chaque niveau, puis remonter progressivement les balises précédentes de la rampe s’il existe des balises plus fiables en amont, dans le même segment inter-niveaux. En montée, c’est le raisonnement inverse qu’il faut adopter : lorsque l’on lit les séquences de détection, on utilisera idéalement la dernière balise de la rampe du niveau immédiatement précédent (i.e. un peu plus profond) pour s’assurer de détecter la changement de niveau. On applique ces stratégies en s’appuyant sur notre visualisation matricielle pour chaque niveau. L’opération est donc manuelle et laborieuse mais on obtient quasi-certainement les balises de changement de niveau les plus fiables pour chaque niveau de l’itinéraire Niv3→Niv30, conformément à l’une des alternatives que nous décrivions dans le chapitre 4. Précisons aussi que l’on choisit une balise du niveau 3 qui sera constamment protégée du phénomène des distributions à deux modes : une autre

balise du niveau 3 (épargnée elle aussi par les disparitions de blocs de données) est en effet située plus proche de la surface. Comme explicité au chapitre 4, elle servira, pour simplifier, de « fusible » pour la balise sélectionnée. De même pour la balise sélectionnée au niveau 30, la balise située à peine plus en profondeur sur le même segment de rampe génère toujours des détections lorsque la balise sélectionnée en génère.

6.4.6 Utilisation de notre algorithme de reconnaissance de trajets

Nous cherchons ici à utiliser notre algorithme de reconnaissance de trajets sur l'itinéraire étudié, en utilisant toutes les améliorations déjà mises en place pour cet algorithme. Selon le motif régulier qui semble discrètement émerger, l'étape correspondante de notre méthodologie devra être là aussi adaptée.

L'adaptation principale consiste à gérer l'absence régulière des détections issues des balises de changement de niveau ayant été retenues pour jalonner l'itinéraire. Pour cela, on ne peut plus forcer les HT à être détectés par chacune de ces balises de changement de niveau sur leurs trajets et on adapte donc cette condition. Maintenant, lorsqu'une balise de changement de niveau ne détecte pas le HT, celui-ci n'est évidemment plus forcé de parcourir les segments inter-niveaux précédent et suivant en moins de 30 minutes chacun puisqu'on ne peut pas le localiser. Lorsqu'il sera détecté par la balise de changement de niveau suivante (elle peut être située plusieurs niveaux en dessous ou au dessus), on le force en revanche à avoir parcouru les multiples segments inter-niveaux successifs en moins d'une heure pour filtrer tout de même des trajets non ordinaires extrêmement aberrants.

De même, pour les détours inter-niveaux, nous excluons uniquement les trajets pour lesquels nous sommes certains que le HT en a effectivement réalisé un. Nous utilisons bien sûr les lots de données mentionnant un changement de niveau pour détecter ces détours inter-niveaux, dès que l'ordre des profondeurs atteintes est incohérent. Nous n'aurons aucune garantie d'avoir effectivement filtré tous les détours de ce type puisque de nombreux niveaux ne sont plus balisés sur certaines périodes.

Enfin, exactement comme nous l'avons décrit au chapitre 4, nous interdirons aux HT d'être détectés par les balises intra-niveaux, mais aussi en amont de la balise initiale du niveau 3 et en aval de la balise finale du niveau 30.

6.4.7 Analyse d'histogrammes de TT sur l'itinéraire prédominant et amélioration résultante

Les histogrammes des TT issus de notre algorithme de reconnaissance de trajets sur l'itinéraire Niv3→Niv30 sont présentés via les figures 6.4 et 6.5, correspond respectivement aux TT conventionnels et autonomes. Le second histogramme correspond aux quarts de travail autonomes complets, car on préfère éviter de risquer de s'intéresser aux trajets autonomes ayant lieu entre les quarts de travail conventionnels. Ce second histogramme paraît évidemment très suspect puisqu'il affiche fièrement une distribution à deux modes, alors que l'autre histogramme semble globalement crédible.

Étant donné que nous avons déjà éliminé le phénomène connu générant des distributions à deux modes, on suspecte très fortement une confusion entre les trajets autonomes et des trajets conventionnels malgré le fait que les deux types de quarts de travail sont séparés entre eux de plusieurs heures. Remarquons par ailleurs que le premier mode de la figure 6.5 correspond plutôt bien à l'unique pic de la figure 6.4. Remarquons enfin, bien plus subtilement, que l'histogramme de cette figure semble présenter un second mode extrêmement discret qui correspondrait précisément au second mode de la figure 6.5. Toutes nos observations pointent donc clairement vers la même explication.

On entreprend alors de s'intéresser aux TT issus de notre algorithme de reconnaissance de trajets pour tâcher de reconnaître un motif dans les conditions opérationnelles associées. Très vite, on remarque que les quarts de travail réputés autonomes ne sont pas du tout homogènes d'une semaine à l'autre puisqu'on observe parfois uniquement des TT correspondant au premier mode, et parfois uniquement des TT correspondant au second mode. Cela signifierait que les quarts de travail autonomes ne le sont réellement que la moitié du temps d'après la proportion des détections incluses dans chacun des deux modes. Étant donné que le premier mode correspond remarquablement bien au pic de la figure 6.4, on en déduit que l'on doit supprimer de notre étude des TT autonomes tous les TT issus de quarts de travail durant lequel la très large majorité des trajets observés se composait de trajets rapides. Pour identifier ces quarts, on calcule simplement la moyenne des TT sur chacun des quarts de travail (en excluant les trajets de plus de 40 minutes du calcul), et l'on supprime de nos données tous les quarts dont la moyenne des TT est inférieure à 20 minutes. Notons que les quarts exclus seront identiques quel que soit l'itinéraire étudié. On représente l'histogramme obtenu à la figure 6.6 (le pas de temps des barres de l'histogramme a été allongé).

Il est limpide que ce filtrage a parfaitement fonctionné. Point important, il permet de supprimer d'éventuels TT longs issus de ces quarts. Cela n'aurait pas été possible en ajoutant simplement un seuil de filtrage aux alentours des 15 minutes pour supprimer le premier mode.

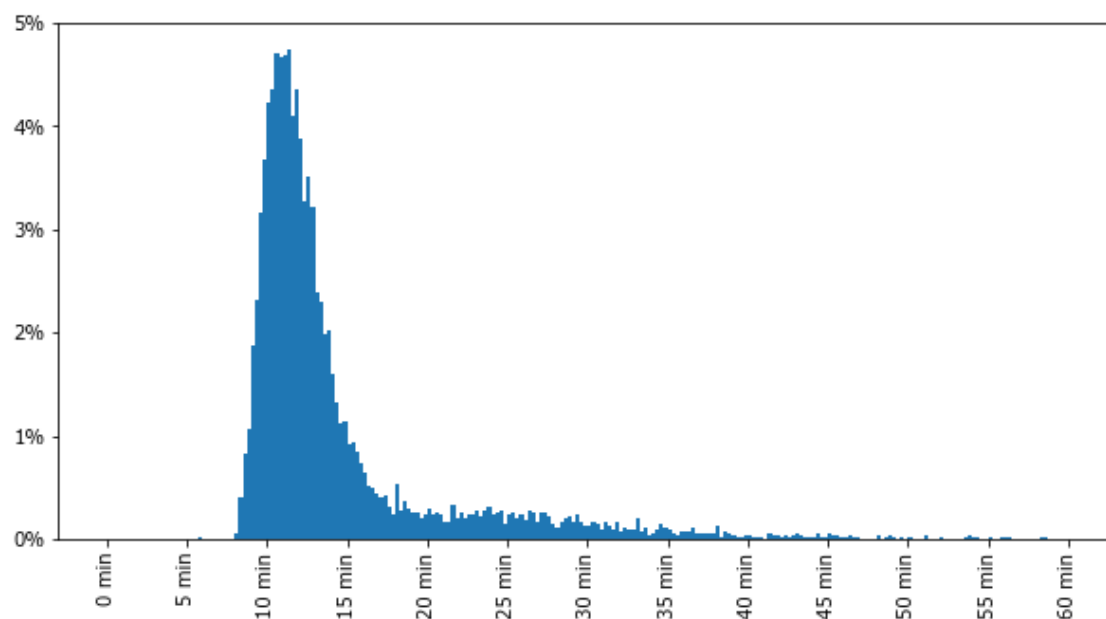


FIGURE 6.4 Histogramme des TT conventionnels sur Niv3→Niv30 dans la mine 2

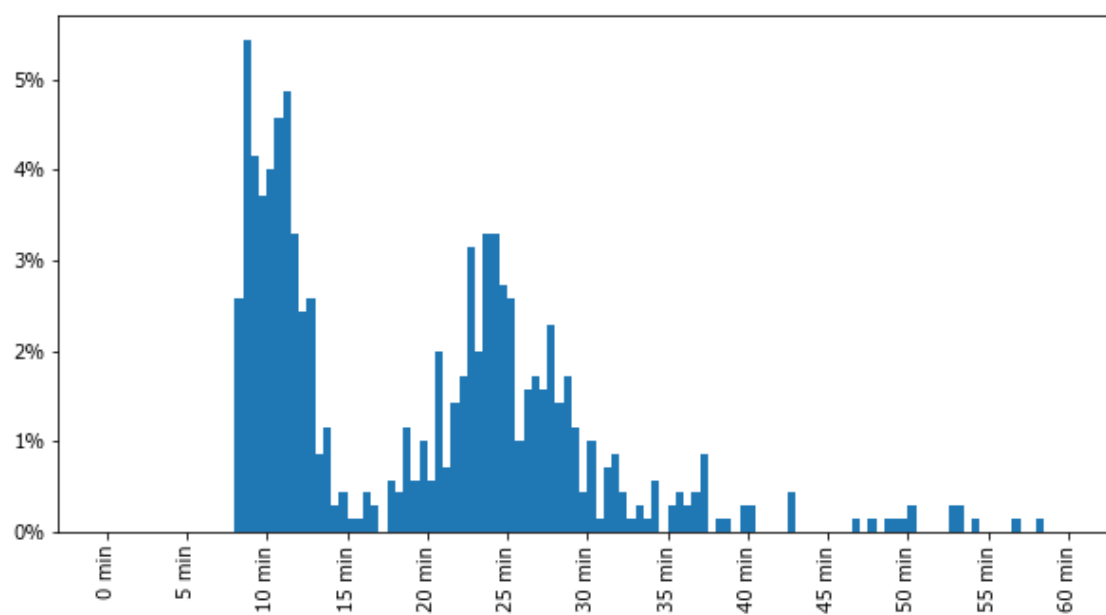


FIGURE 6.5 Histogramme à deux modes suspect des TT théoriquement autonomes sur Niv3→Niv30 dans la mine 2

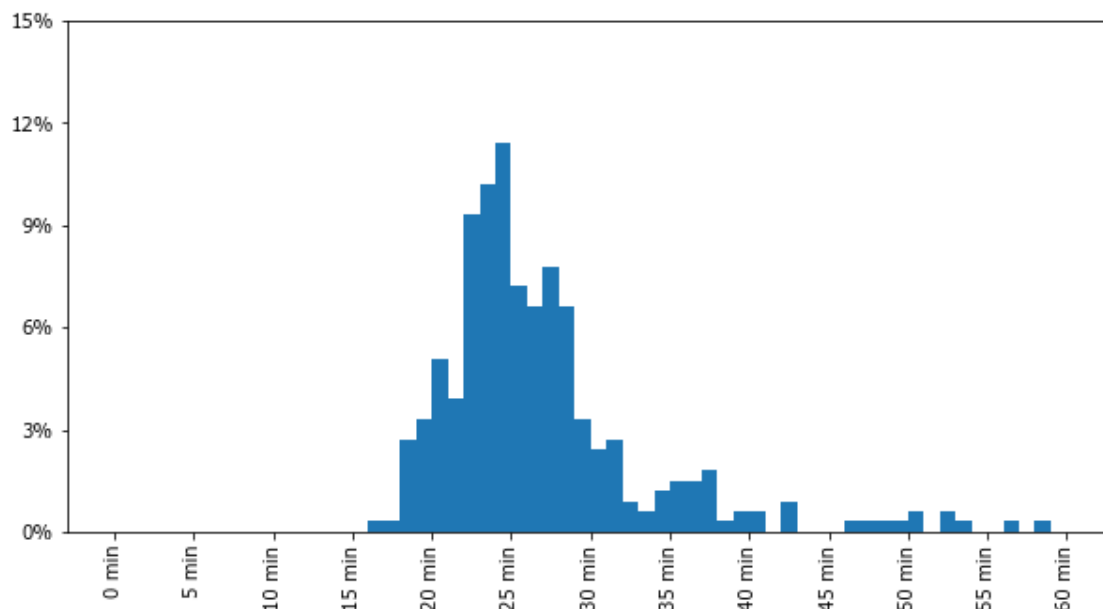


FIGURE 6.6 Histogramme des TT autonomes sur Niv3→Niv30 dans la mine 2

Précisons bien par ailleurs que notre modèle de prédiction de TT n'aurait jamais pu prédire lesquels de ces quarts de nuit de la fin de semaine sont réellement purement autonomes et inversement. Les planificateurs l'auraient en revanche su à l'avance, donc il est parfaitement cohérent de bien filtrer les TT non autonomes ici, en amont de l'utilisation du modèle de prédiction de TT.

Forts de ce filtrage couronné de succès, nous tentons alors de filtrer le très discret second mode de la figure 6.4. Une longue analyse nous convint que ce nouveau filtrage est en fait bien trop laborieux à réaliser a posteriori. En effet, la lecture des séquences de détection ne fait apparaître aucun motif régulier dans les conditions opérationnelles associées. Tous les HT semblent être pilotés habituellement de manière conventionnelle sur les quarts réputés conventionnels, mais l'un d'entre eux semble parfois tout à coup être piloté par le système de conduite autonome car ses TT deviennent tout à coup beaucoup plus longs sur l'itinéraire étudié. En revanche, il est exceptionnellement complexe de généraliser ces observations de manière rigoureuse, car elles pourraient parfois simplement résulter de répétitions de trajets ordinaires longs effectués par l'un des HT sur l'itinéraire Niv3→Niv30. Nous décidons finalement d'ignorer cette problématique, car il n'est pas absolument nécessaire de s'en occuper : la proportion de trajets réputés conventionnels dont les durées sont aberrantes est déjà plus faible que ce que nous observions pour la mine 1, ce qui devrait faciliter l'obtention de bonnes performances de prédiction.

Concernant les trajets en montée, notre filtrage précédent sur les quarts de travail autonomes a parfaitement fonctionné, comme le montre la figure 6.7. Pour les TT conventionnels en revanche, nous n'avons aucune explication solide pour justifier le fait que le second mode soit bien plus marqué qu'en descente (voir fig. 6.8), comme si les HT remontaient régulièrement en toute autonomie. Nous n'accorderons toutefois qu'un faible intérêt aux trajets en montée puisqu'ils se sont avérés plus faciles à prédire que les trajets en descente. Par ailleurs, la variance des trajets conventionnels en montée est visiblement bien plus faible qu'en descente (si on ignore le second mode) ce qui réduit notablement l'intérêt de notre modèle de prédiction de TT comparativement à l'utilisation de la moyenne des TT contenus dans le premier mode.

D'un point de vue plus général, l'observation la plus surprenante est évidemment ici la très longue durée des trajets autonomes comparativement aux trajets conventionnels. Habituellement, les HT autonomes semblent mettre près de deux fois la durée de trajet des HT pilotés manuellement. Il est difficile de trouver une unique explication imparable à ce phénomène. En revanche, nous pouvons émettre quelques hypothèses complémentaires.

- Tout d'abord, les quarts de travail entièrement autonomes ont théoriquement lieu sans qu'aucun humain ne soit présent dans la mine. On peut alors supposer qu'une collision d'un HT serait extrêmement fâcheuse sur ces périodes. Par conséquent, le système de conduite autonome devrait être configuré pour respecter des critères de sécurité absolument draconiens. Ces dernières ralentiraient exceptionnellement le HT, ainsi que les manœuvres d'évitement ou de croisement des HT comparativement aux capacités d'opérateurs aguerris. Cela expliquerait à la fois la durée moyenne très élevée des TT autonomes et leur très forte variance ;
- Par ailleurs, le système de conduite autonome de HT en souterrain est une technologie évidemment récente. Il serait cohérent de penser que cette technologie est encore immature. De multiples facteurs pourraient alors obliger le HT à adopter une faible vitesse de croisière. Citons par exemple des normes de sécurité très strictes et une puissance de calcul insuffisante (et/ou une optimisation insuffisante de la stratégie de calcul) pour déterminer en temps réel la trajectoire à adopter si le HT se déplace à la vitesse maximale autorisée dans la mine 2. Idem pour les manœuvres de croisement, qui constituent certainement un remarquable défi à gérer et, plus encore, à optimiser.

Une autre observation rendue possible par nos histogrammes est la très faible quantité de trajets autonomes détectables durant les quarts de travail de la fin de semaine comparativement au nombre total de trajets conventionnels détectés : on en recense respectivement 321 et 7742 sur l'itinéraire Niv3→Niv30. Mentionnons au passage que nous disposons de plus de trajets conventionnels complets sur cet itinéraire que sur l'itinéraire Surface→Niveau350 de

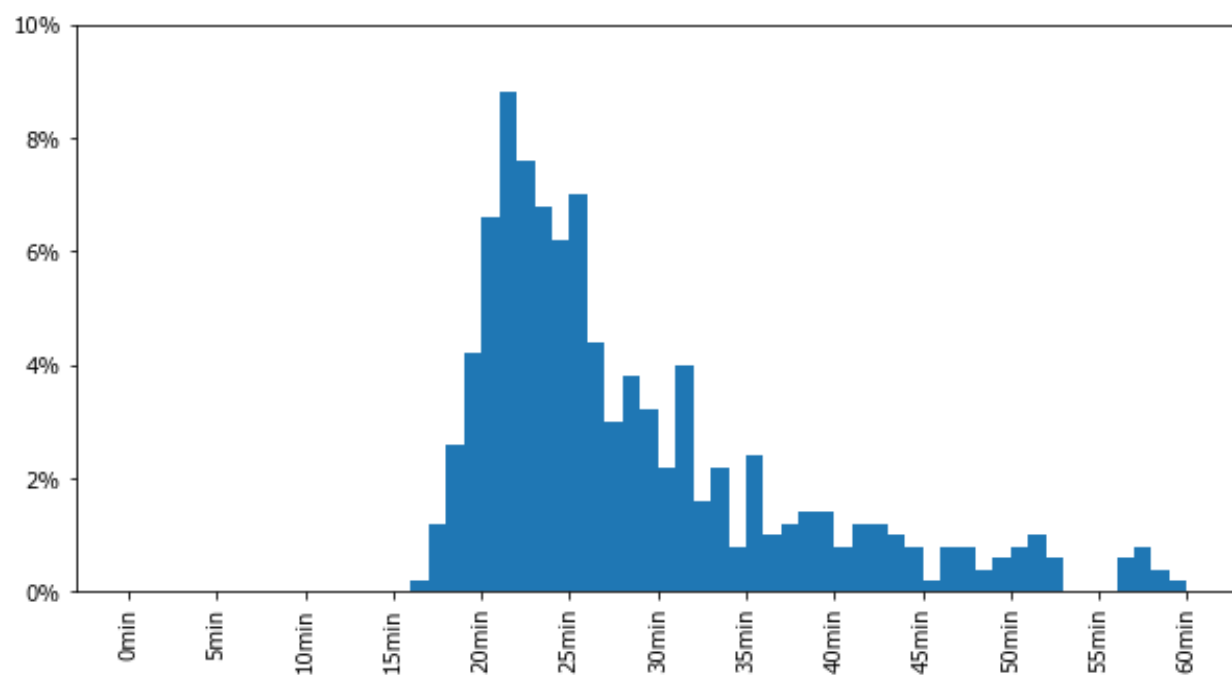


FIGURE 6.7 Histogramme des TT autonomes sur Niv30→Niv3 dans la mine 2

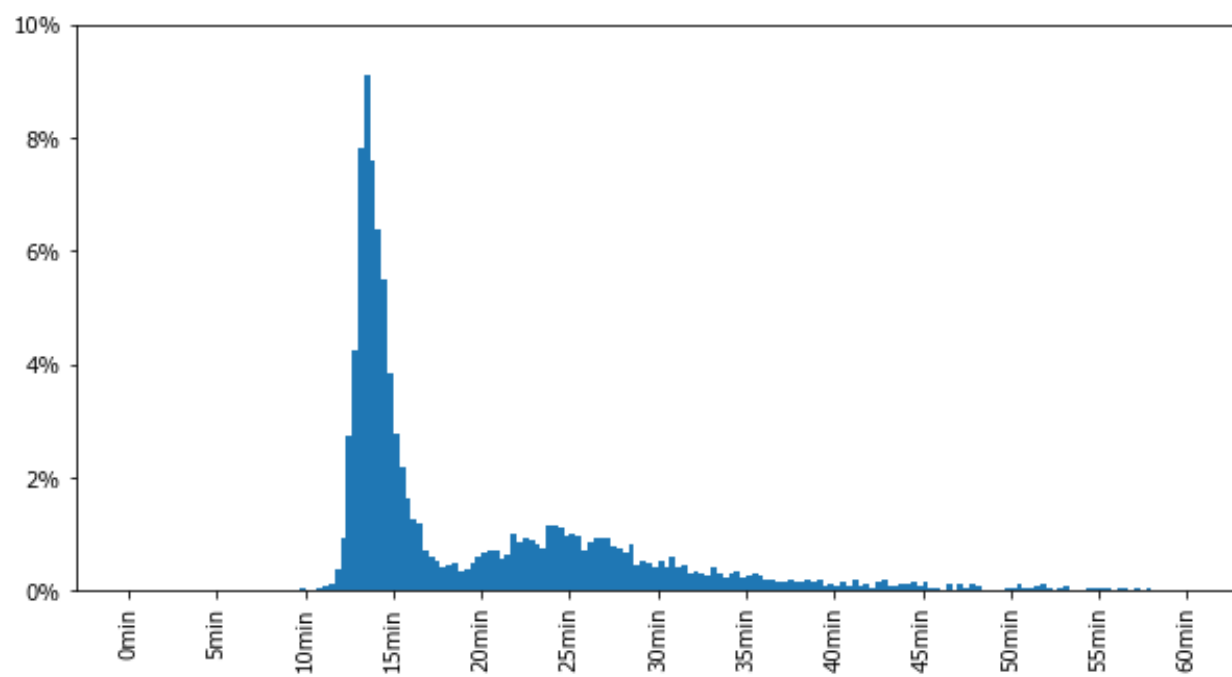


FIGURE 6.8 Histogramme des TT conventionnels sur Niv30→Niv3 dans la mine 2

la mine 1 (un peu plus de 5000 observations). Étant donné que des responsables de la mine 2 ont annoncé que le système de conduite autonome leur permettait d'augmenter d'environ 10% la productivité de la mine 2, on peut imaginer que l'on devrait trouver un total de près de 800 ($7742 * 10/100$) trajets autonomes sur l'itinéraire Niv3→Niv30, soit près de 450 ($774 - 321$) trajets entre les quarts de travail conventionnels, tout au long de la semaine. Cela semble cohérent avec le volume horaire que représentent ces durées chaque semaine par rapport aux quarts de travail purement autonomes. Nous ne pousserons pas plus loin notre investigation, car ces trajets s'annoncent très difficiles à cerner.

6.4.8 Génération de variables d'intérêt additionnelles

Dans la présente sous-section, nous allons nous atteler à générer les dernières variables d'intérêt manquantes à notre modèle de prédiction de TT.

Nous ne consacrons pas de sous-section à la génération de variables d'intérêt temporelles, car cette étape est parfaitement identique à la description que nous en avons faite dans le chapitre 4.

En revanche, concernant l'estimation du nombre de HT se déplaçant activement sur l'itinéraire étudié durant chaque quart de travail, force est de constater que l'on ne peut pas négligemment compter le nombre total de détections d'un HT par chaque balise de changement de niveau de l'itinéraire Niv3→Niv30 sur chaque quart de travail. En effet, les absences de données de détection de balises ne nous permettent absolument pas de fixer un seuil qui serait valable quelle que soit la période étudiée. Plutôt que d'implémenter un seuil dynamique pour adapter cette étape de notre méthodologie, nous décidons de simplifier le problème en établissant les constats suivants :

- Nous savons qu'à cette étape seuls les trajets de l'itinéraire A→B (Niv3→Niv30) nous intéressent ;
- Nous nous intéressons donc uniquement aux périodes durant lesquelles la balise A et la balise B fonctionnent simultanément ;
- Nous savons par ailleurs que les HT remontent régulièrement à la surface lors du transport de minerai et qu'ils sont donc détectés par la balise A ;
- Il n'y a qu'une seule rampe, tous les HT qui remontent à la surface peuvent donc croiser les HT se déplaçant sur l'itinéraire Niv3→Niv30, ce qui ralentit ces derniers ; et
- Les HT qui passent aussi par la balise B sont les plus gênants puisqu'ils ont parcouru l'intégralité de l'itinéraire même s'ils sont moins souvent détectés par la balise A.

Aussi, nous finissons par arriver à la conclusion qu'en additionnant pour chaque HT ses

détections respectives par les balises A et B uniquement pendant chaque quart de travail, on aura un bon aperçu de son « score » d'activité sur cet itinéraire (et donc de sa capacité de perturbation des HT en descente) à chaque quart de travail. Après quelques itérations manuelles, on fixe à **huit** détections des balises A et B le seuil permettant de considérer qu'un HT est actif. Cela équivaut théoriquement à deux allers-retours sur un itinéraire incluant l'itinéraire A→B (chacune des balises A et B est détectée à l'aller et au retour) ou à quatre allers-retours sur un itinéraire plus court (la balise A est détectée à l'aller et au retour).

6.4.9 Évaluation de l'influence de variables d'intérêt sur les TT

Passons à présent à la visualisation graphique de l'influence de différentes variables d'intérêt sur les TT.

Le seul seuil de filtrage utilisé sera le seuil des 30 minutes inter-niveaux (et une heure pour plusieurs segments inter-niveaux successifs). Nous utiliserons la médiane et la moyenne pour trouver nos remarques concernant la significativité statistique des différentes variables. La sous-représentation de certains HT dans nos données nous obligera à augmenter la taille maximale autorisée pour les barres d'erreur et à réduire le nombre minimal d'échantillons par barre bien que cela rende parfois la lecture des données un peu plus difficile. Toutes nos analyses se feront pour les quarts conventionnels en raison du manque criant d'observations de trajets autonomes pour ces analyses.

Concernant l'influence du quart de travail, nous pouvons affirmer à un niveau de confiance de 99,7% que la moyenne et a fortiori la médiane des TT de jour dépassent très largement la moyenne et la médiane des TT de nuit sur l'itinéraire étudié, comme le montrent respectivement les figures 6.9 et 6.10. Ces observations, contraires à celles faites lors de l'étude de la mine 1, pourraient parfaitement s'expliquer par une simple baisse du nombre de HT actifs durant les quarts de nuit.

En l'occurrence, concernant l'influence du nombre de HT actifs par quart de travail sur les TT, la tendance est éminemment positive pour la moyenne et plus encore pour la médiane (voir fig. 6.11 et 6.12). Rappelons que, lors de l'étude de la mine 1, la significativité statistique de cette relation était validée de peu à un niveau de confiance de 95% avec le seuil des 30 minutes, et que nous ne pouvions pas conclure à la significativité statistique du nombre de HT sur les TT observés lorsque nous appliquions le seuil de Tukey et le seuil des 5%. Ici, nous pouvons clairement affirmer à un niveau de confiance de 95% que lorsque cinq HT sont actifs (et dans une moindre mesure six) les TT observés sont significativement supérieurs aux trajets ayant eu lieu lorsque seulement deux, trois voire quatre HT étaient actifs. À vrai dire, nous pouvons même affirmer à un niveau de confiance de 99,7% que les TT médians

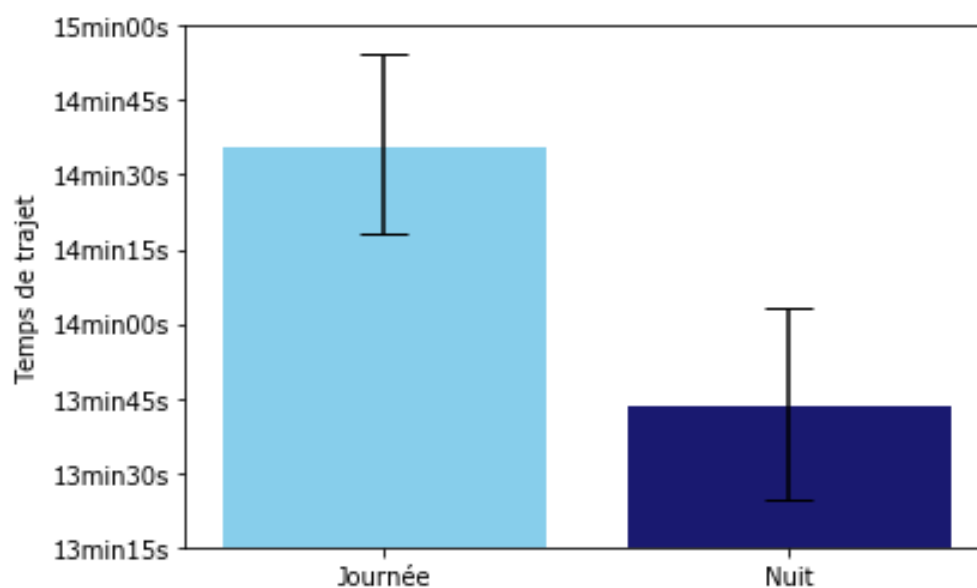


FIGURE 6.9 Diagramme à barres du TT conventionnel moyen selon le quart de travail sur Niv3→Niv30 dans la mine 2 et barres d'erreur pour un niveau de confiance de 99,7%

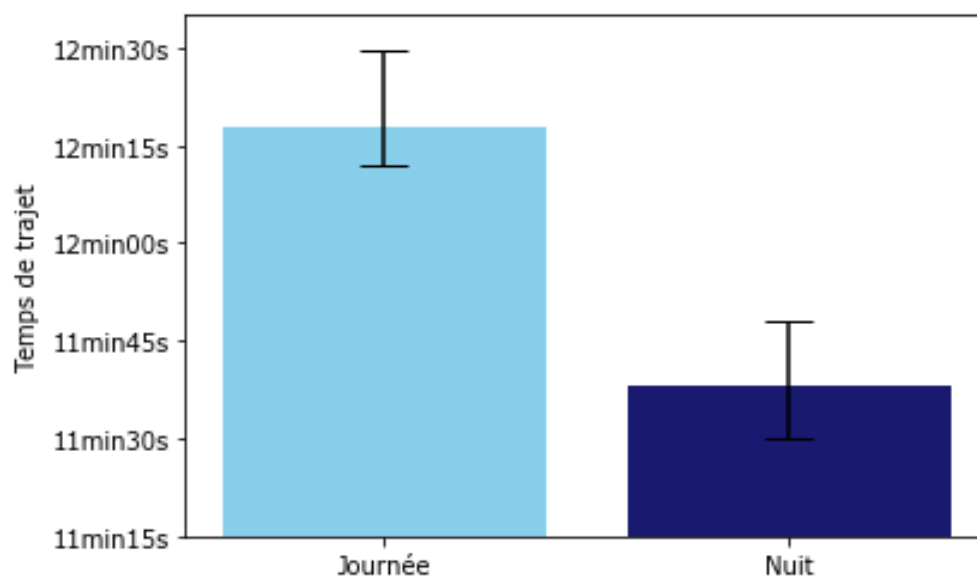


FIGURE 6.10 Diagramme à barres du TT conventionnel médian selon le quart de travail sur Niv3→Niv30 dans la mine 2 et barres d'erreur pour un niveau de confiance de 99,7%

pour cinq HT actifs sont significativement supérieurs aux TT médians pour un à quatre HT actifs. Il apparaît ainsi évident que l'on ne peut pas augmenter le nombre de HT affectés à une mine souterraine donnée sans faire significativement décroître leur productivité respective (en termes de TT uniquement). Ici, les TT autonomes inclus dans les quarts de travail réputés conventionnels ne biaisent pas les résultats (puisque notre remarque est aussi valable pour les trajets autonomes), au contraire des très nombreux TT non ordinaires que présentaient certainement les histogrammes de la mine 1. Selon nous, c'est la raison pour laquelle la significativité statistique est si marquée ici, en comparaison avec nos observations pour la mine 1.

Concernant les TT moyens et médians des différents HT de la mine 2, certains HT se démarquent tout particulièrement (voir fig. 6.13 et 6.14). Le HT1, en particulier, a été peu actif sur la période étudiée mais concernant la médiane de ses TT nous pouvons tout de même affirmer à un niveau de confiance de 99,7% qu'il sous-performe largement voire très largement cinq des six autres HT sur l'itinéraire Niv3→Niv30. Nous parlons tout de même de près de 10% de différence du TT médian de ce HT par rapport aux TT médians des HT les plus performants. Concernant la moyenne de ses TT, on peut aussi affirmer à un niveau de confiance de 99,7% qu'il sous-performait trois des six autres HT. À une moindre échelle, le HT2, le HT3 et le HT4 sous-performent significativement le HT5 et le HT7 lorsqu'on regarde la médiane de leurs TT, là encore à un niveau de confiance de 99,7%.

Enfin, pour la mine 2 seulement, nous avons pu prouver aussi la significativité statistique du jour de la semaine sur le TT moyen des HT, à un niveau de confiance de 95% (voir fig. 6.15). Les trajets du dimanche sont ainsi significativement plus longs que ceux du mardi et du vendredi. Plus généralement, il semblerait que les TT observés durant les deux quarts de travail conventionnels de la fin de semaine soient plus longs que ceux des jours de semaine. De prime abord, nous aurions pourtant été amenés à penser qu'un moindre nombre de HT seraient actifs durant ces quarts de travail, générant des TT plus courts la fin de semaine. Rappelons que les quarts de nuit s'étalent sur deux jours à la fois, et que la comparaison des TT selon le jour de la semaine n'est donc pas parfaite. Étrangement, les observations précédentes ne sont pas valables pour les TT médians, qui permettraient pourtant de démontrer plus facilement la significativité des autres variables d'intérêt sur les TT. Nous n'avons pas creusé davantage cette piste.

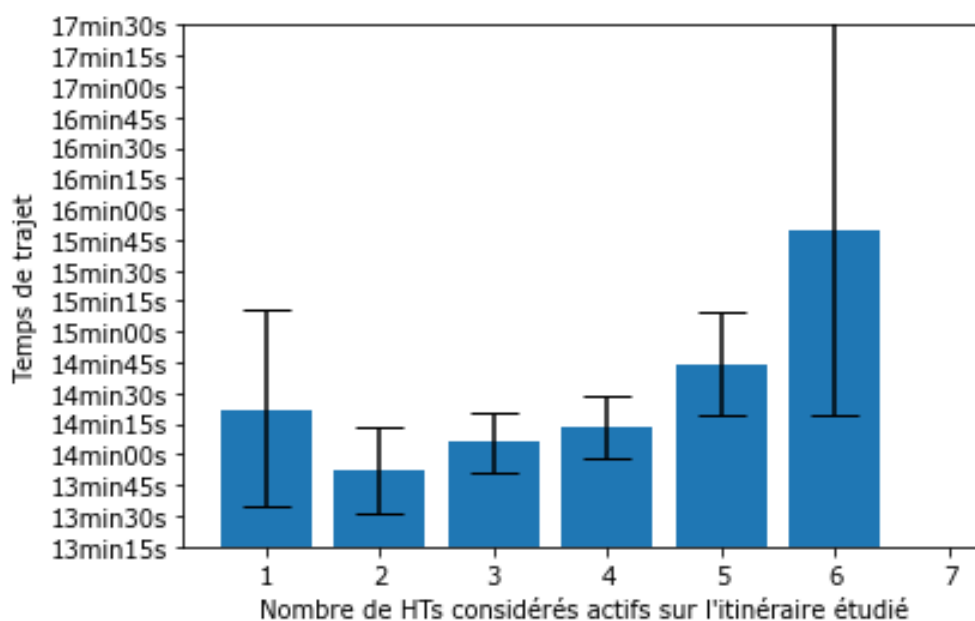


FIGURE 6.11 Diagramme à barres du TT conventionnel moyen selon le nombre de HT actifs sur Niv3→Niv30 dans la mine 2 et barres d'erreur pour un niveau de confiance de 95%

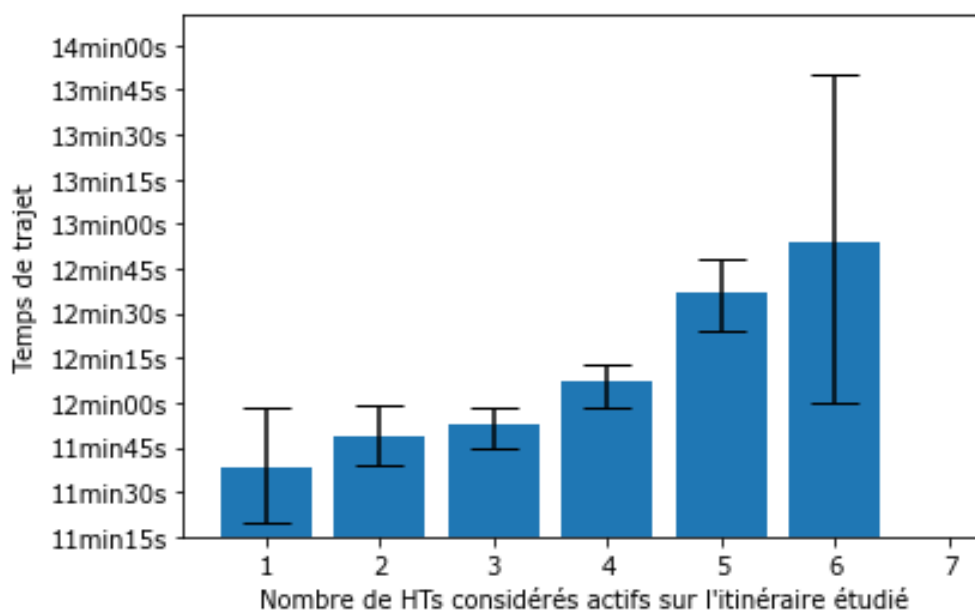


FIGURE 6.12 Diagramme à barres du TT conventionnel médian selon le nombre de HT actifs sur Niv3→Niv30 dans la mine 2 et barres d'erreur pour un niveau de confiance de 95%

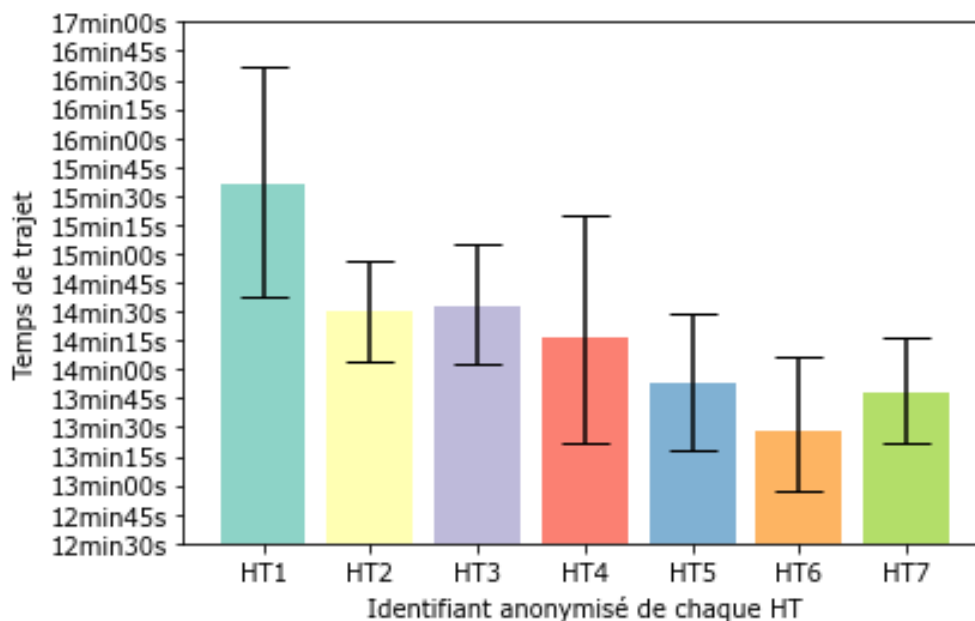


FIGURE 6.13 Diagramme à barres du TT conventionnel moyen selon le numéro du HT sur Niv3→Niv30 dans la mine 2 et barres d'erreur pour un niveau de confiance de 99,7%

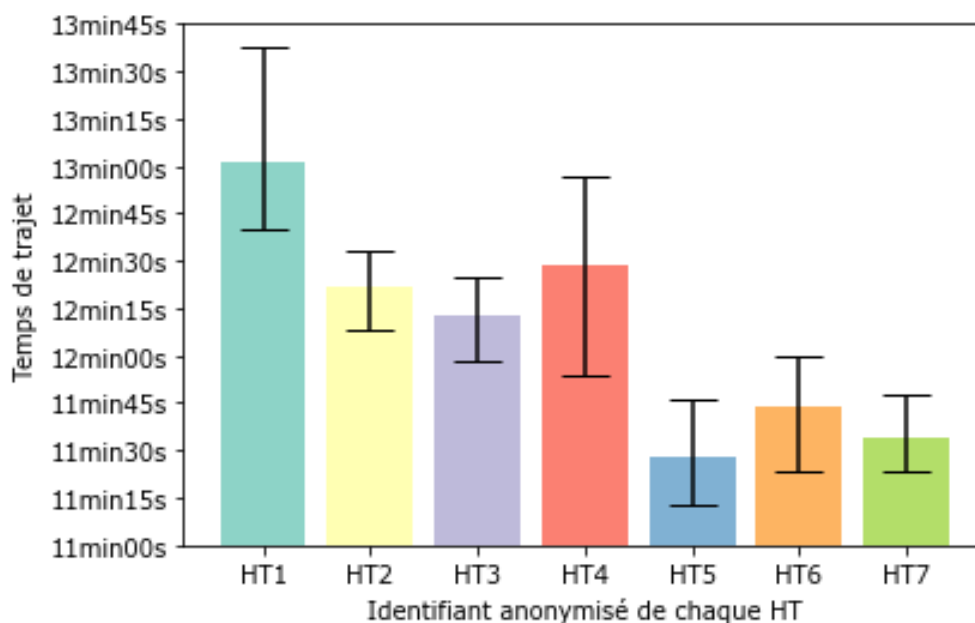


FIGURE 6.14 Diagramme à barres du TT conventionnel médian selon le numéro du HT sur Niv3→Niv30 dans la mine 2 et barres d'erreur pour un niveau de confiance de 99,7%

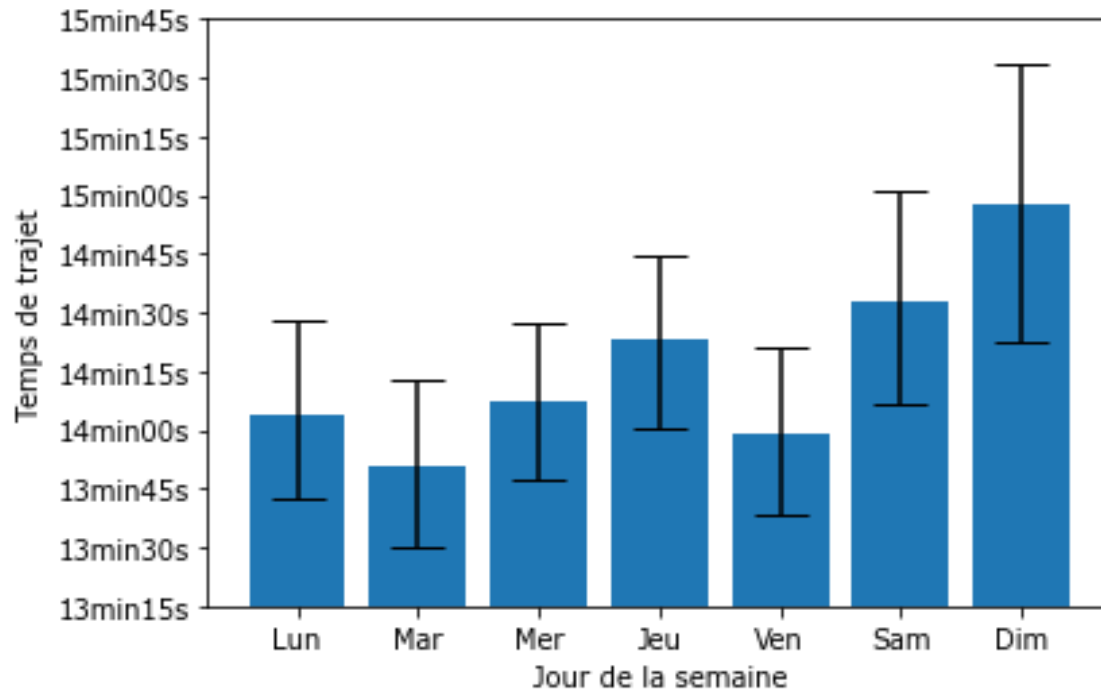


FIGURE 6.15 Diagramme à barres du TT conventionnel moyen selon le jour de la semaine sur Niv3→Niv30 dans la mine 2 et barres d'erreur pour un niveau de confiance de 95%

6.4.10 Inventaire des pistes d'amélioration relatives aux capteurs, à l'acquisition des données et à leur gestion

Étant donné qu'une majeure part des problématiques rencontrées sont issues des mêmes causes racines et que les détectations de balises semblent extrêmement fiables hors des périodes de disparition de données, l'inventaire des pistes d'amélioration est réduit pour la mine 2.

Concernant les absences de données de balises par blocs complets, qui est clairement la cause racine centrale, les raisonnements logiques présentés auparavant nous mènent à recommander :

- De mener une surveillance systématique de la « santé » des balises, en implémentant des indicateurs visuels automatiques à l'image de la visualisation matricielle précédente, pour agir au plus tôt lors des disparitions de données ; et
- D'augmenter la bande passante de transfert de données vers la BDD de la mine 2 pour éviter absolument la disparition de données cruciales pour la recherche, bien que celles-ci soient aujourd'hui non exploitées par les industriels du site pour la reconnaissance de trajets de véhicules.

Il serait particulièrement important de surveiller très attentivement les balises déjà existantes

à l'entrée des différents niveaux (balises de changement de niveau valables dans les deux sens de circulation) et d'en installer là où aucune balise ne correspond parfaitement à cette description.

Par ailleurs, nous proposons d'ajouter un nouvel attribut dans les lots de données de détection de balises. Il indiquerait si le HT en cours de déplacement est piloté manuellement ou bien de manière autonome, ce qui est évidemment connu au moment du trajet. Cet attribut permettrait selon nous d'améliorer très considérablement les performances de notre modèle de prédiction de TT dans le cas des trajets conventionnels. En effet, les trajets autonomes qui se sont mélangés à ces derniers risquent d'avoir des effets similaires à ceux des trajets non ordinaires lorsque leurs TT seront fournis à notre modèle de prédiction. L'effet sera tout de même moins intense puisque ces TT sont plus fortement corrélés aux conditions opérationnelles d'entrée que les TT non ordinaires.

6.4.11 Prétraitement des données

Le prétraitement des données de la mine 2 est globalement identique à celui de la mine 1. Seule la normalisation des quarts de travail (jour/nuit) diffère légèrement : malgré le fait que le nombre de HT actifs y soit sûrement pour quelque chose, les quarts de nuit font visiblement réduire le TT des HT comparativement aux quarts de jour. On inverse donc leurs valeurs normalisées entre elles, associant ainsi une valeur de 1 aux quarts de jour et de -1 aux quarts de nuit.

6.4.12 Conclusion de la préparation de données

Au terme d'adaptations conséquentes de notre méthodologie, comparativement à l'application que nous en avons faite pour la mine 1, nous avons visiblement pu préparer avec succès les données de détection de balises de la mine 2 pour fournir les variables d'intérêt nécessaires au bon fonctionnement de notre modèle de prédiction de TT.

Les adaptations réalisées ont été rendues incontournables par une problématique absolument majeure d'absence de données de balises, que nous avons globalement pu cerner et qui ne nous a pas définitivement bloqué lors de l'application de notre méthodologie. Au contraire, malgré cette problématique, un concours de circonstances partiellement expliquées nous a permis d'obtenir un histogramme de TT conventionnel en descente (sur Niv3→Niv30) contenant finalement nettement moins de TT aberrants que l'un des histogrammes obtenu pour la mine 1 (sur Surface→Niv300 via la rampe 1) malgré la longueur similaire des itinéraires étudiés et leurs propriétés globalement similaires.

Ainsi, en admettant que ces deux itinéraires sont totalement comparables, force est de constater que la stratégie d’installation des balises (incluant leur position, leur nombre et leur paramétrage) semble être très clairement le critère le plus critique pour exploiter les détections de ces dernières aux fins de reconnaissance précise des trajets et des TT associés. Cette stratégie d’installation des balises est en effet nettement plus aboutie et précise dans la mine 2 que dans la mine 1. Toujours aux mêmes fins, cette stratégie semble bien plus importante que la capacité des balises à rester toujours fonctionnelles ou à pouvoir faire sauvegarder leurs détections l’essentiel du temps. Ce constat nous semble particulièrement déstabilisant étant donné les lourdes adaptations auxquelles nous avons été poussés pour limiter les pertes de précision et l’incapacité finale de notre algorithme à détecter, dans la mine 2, la même proportion de détours évidents que dans la mine 1 (boucles inter-niveaux et intra-niveaux).

L’adaptabilité de notre modèle de préparation de données a finalement été parfaitement mise en évidence par l’intense mise à l’épreuve à laquelle il vient de se soumettre avec succès. L’adaptation nécessaire peut en revanche être notablement exigeante à mettre en œuvre dépendamment des problématiques propres à chaque site minier étudié.

6.5 Application de notre modèle intégré de prédiction de TT

Nous allons pouvoir nous attaquer maintenant à l’application de notre modèle intégré de prédiction de TT à la mine 2. Gageons qu’il sera moins significativement perturbé par les spécificités de cette dernière que notre modèle de préparation des données. Comparativement à ce dont il avait été capable pour la mine 1, ce dernier a en effet fourni des ensembles très similaires de TT et de conditions opérationnelles associées pour la mine 2, mais des différences persistent.

6.5.1 Sélection de l’itinéraire de test

Comme annoncé, nous continuerons de nous intéresser à l’itinéraire Niv3→Niv30. Nous avons déjà entièrement défini les segments inter-niveaux ainsi que les balises de changement de niveau correspondantes. Concernant le sectionnement de l’itinéraire en segments majeurs, aucun sectionnement ne nous semble totalement évident, en raison de l’architecture très répétitive de la mine 2. Finalement, étant donné l’importance de ce sectionnement d’après nos conclusions du chapitre 5, nous décidons de retenir l’une des balises du niveau 12 pour sectionner notre itinéraire. Cette dernière, située dans la rampe, est remarquable d’après notre visualisation matricielle. Située environ au tiers de l’itinéraire Niv3→Niv30, c’est presque la seule balise de l’itinéraire ayant été globalement active à chacun des quarts de travail durant

lesquels notre balise B du niveau 30 a été elle-même active.

Naturellement, nous continuerons à utiliser simultanément le seuil de filtrage des 30 minutes inter-niveaux et le seuil de filtrage d’une heure pour plusieurs segments inter-niveaux successifs.

6.5.2 Partitionnement des données

Toute la théorie ayant déjà été présentée, tâchons dès à présent de partitionner les données d’entrée récoltées.

Pour vérifier que nos observations du chapitre 5 se confirment ici, on tente dans un premier temps d’appliquer le GMM à X_{train} conventionnel (i.e. l’ensemble d’apprentissage des conditions opérationnelles des TT conventionnels sur Niv3→Niv30). L’observation de l’évolution de l’AIC pour un partitionnement allant de un à 20 groupe(s) est effectivement similaire à ce que nous pouvions observer au chapitre 5, comme le montre la figure 6.16. Ainsi, avec les données opérationnelles dont nous disposons, nous considérons que cette stratégie est vouée à l’échec quel que soit le site minier étudié (pour prédire les TT sur un itinéraire précis donné).

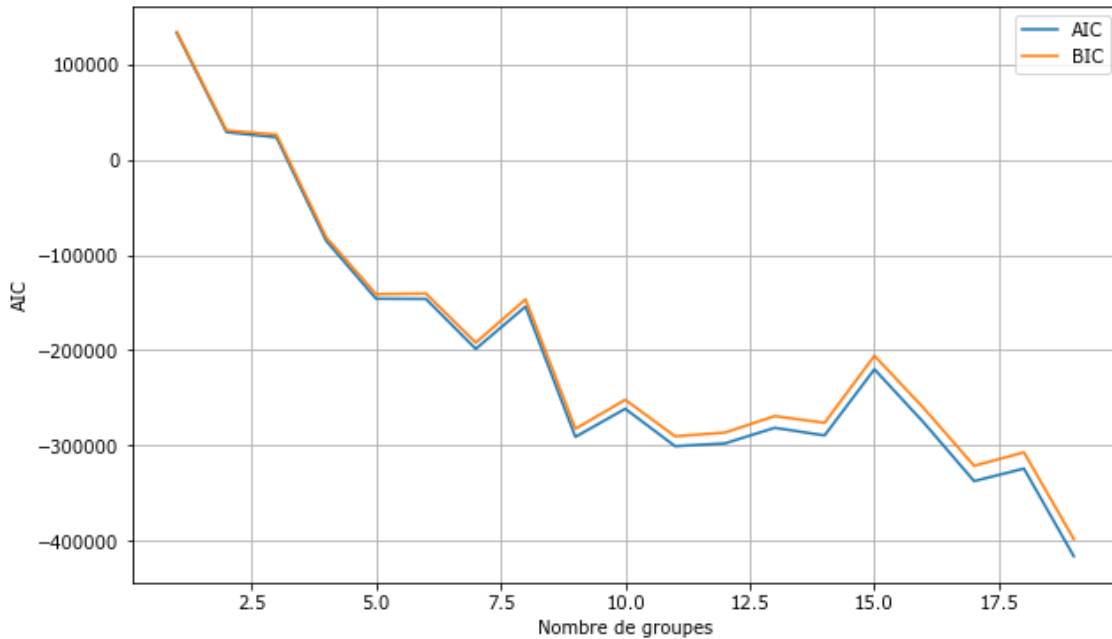


FIGURE 6.16 Évolution de l’AIC et du BIC lorsque le GMM est directement appliqué à X_{train} conventionnel sur Niv3→Niv30 dans la mine 2

Lorsque l’on répète les mêmes opérations sur y_{train} , on trouve qu’un partitionnement en quatre groupes est optimal d’après l’AIC. On représente l’évolution de cette dernière à la

figure 6.17. On représente par ailleurs le partitionnement en quatre groupes de y conventionnel tout entier à la figure 6.18, une fois que le GMM a été ajusté sur y_{train} .

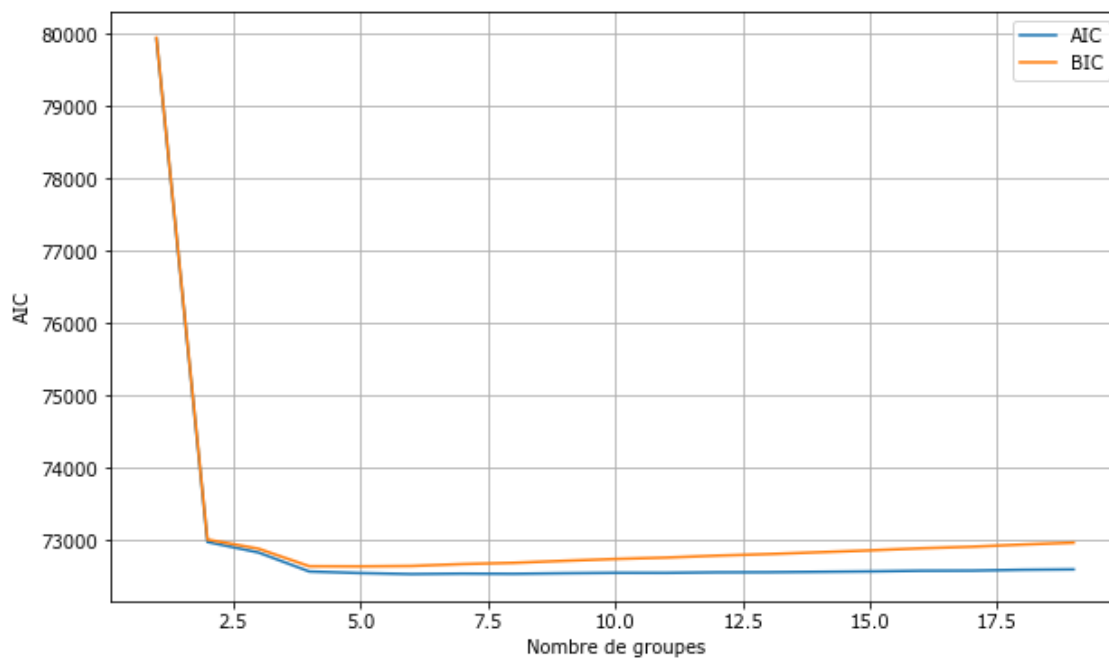


FIGURE 6.17 Évolution de l'AIC et du BIC lorsque le GMM est appliqué à y_{train} conventionnel sur Niv3→Niv30 dans la mine 2

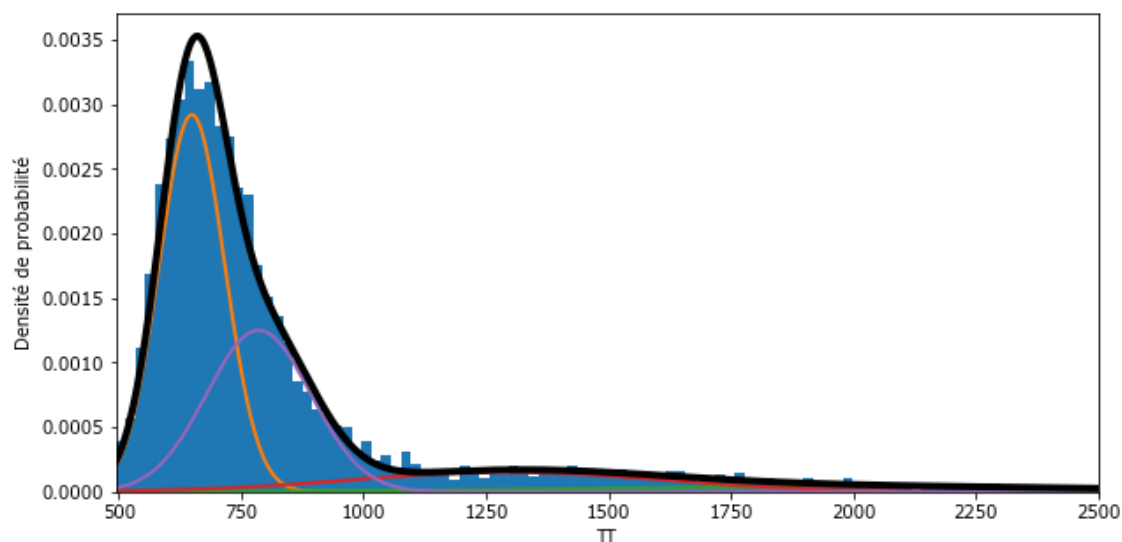


FIGURE 6.18 Partitionnement en quatre groupes de y conventionnel sur Niv3→Niv30 dans la mine 2 via le GMM

Pour les TT conventionnels, on optimise alors les hyperparamètres de notre MLP pour prédire

les groupes établis par le GMM via la succession d'étapes expliquées au chapitre 5. Les hyperparamètres optimaux du MLP sont rassemblés dans le tableau 6.4.

TABLEAU 6.4 Combinaison optimale des hyperparamètres du MLP pour les trajets conventionnels

Hyperparamètre	Valeur optimale
$nbUnits1MLP$	64
$nbUnits2MLP$	8
$tailleLotsMLP$	128
$tauxApprMLP$	1×10^{-2}

Avec ces hyperparamètres, on ajuste ensuite le MLP sur X_{train} conventionnel puis on l'applique à X_{test} conventionnel. Nous obtenons des prédictions apparemment correctes des groupes proposés par le GMM, comme le montre la figure 6.19.

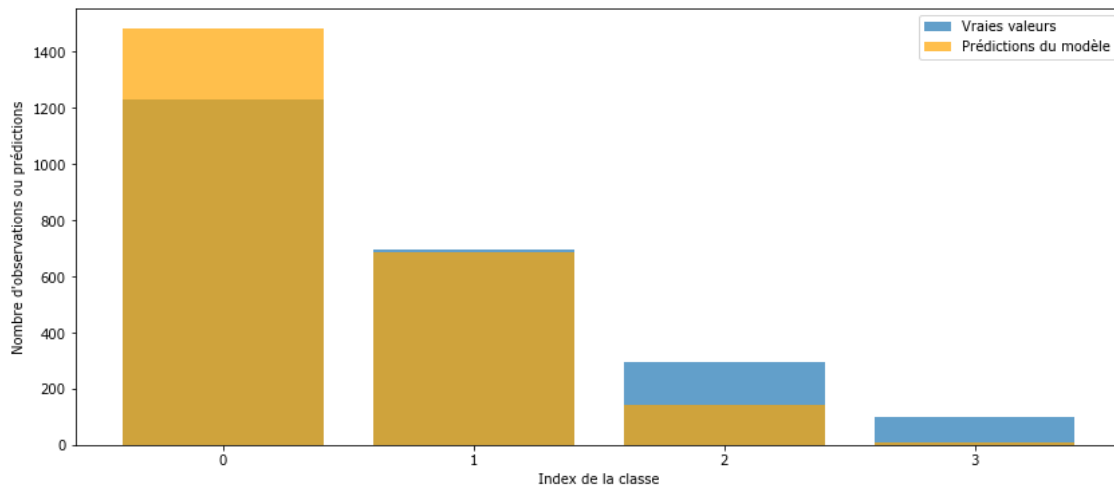


FIGURE 6.19 Diagramme à barres des groupes prédits par le MLP par rapport à ceux prévus par le GMM pour les TT conventionnels sur Niv3→Niv30 dans la mine 2

Concernant les trajets autonomes, étant donné que notre échantillon de trajets est notablement faible, on juge qu'il faut minimiser coûte que coûte le nombre de gaussiennes pour éviter au maximum les écueils du surajustement. D'après la figure 6.20, nous jugeons qu'un partitionnement en deux groupes serait alors optimal.

On obtient alors le partitionnement représenté à la figure 6.21.

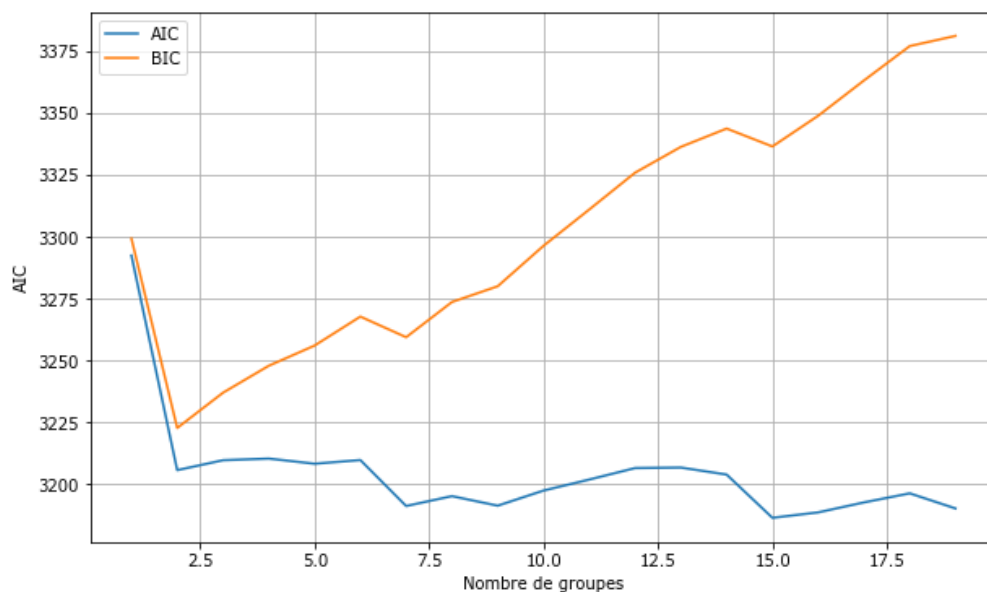


FIGURE 6.20 Évolution de l'AIC et du BIC lorsque le GMM est appliqué à y_{train} autonome sur Niv3→Niv30 dans la mine 2

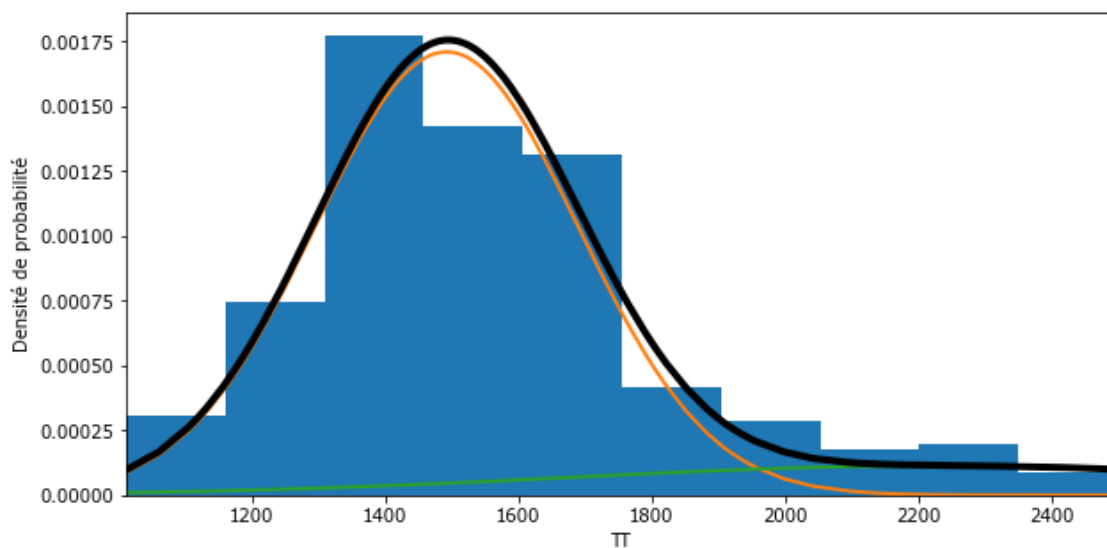


FIGURE 6.21 Partitionnement en quatre groupes de y autonome sur Niv3→Niv30 dans la mine 2 via le GMM

On optimise alors le MLP sur ces trajets autonomes, ce qui nous donne la combinaison optimale d'hyperparamètres du tableau 6.5

TABEAU 6.5 Combinaison optimale des hyperparamètres du MLP pour les trajets autonomes

Hyperparamètre	Valeur optimale
<i>nbUnits1MLP</i>	224
<i>nbUnits2MLP</i>	16
<i>tailleLotsMLP</i>	16
<i>tauxApprMLP</i>	1×10^{-2}

Une fois optimisé et ajusté, le MLP pourrait être particulièrement doué pour prédire les groupes du partitionnement correspondant aux TT autonomes, comme l'indique la figure 6.22. Nous avons fait de nombreux essais et le diagramme à barres obtenu est toujours très similaire à celui-ci.

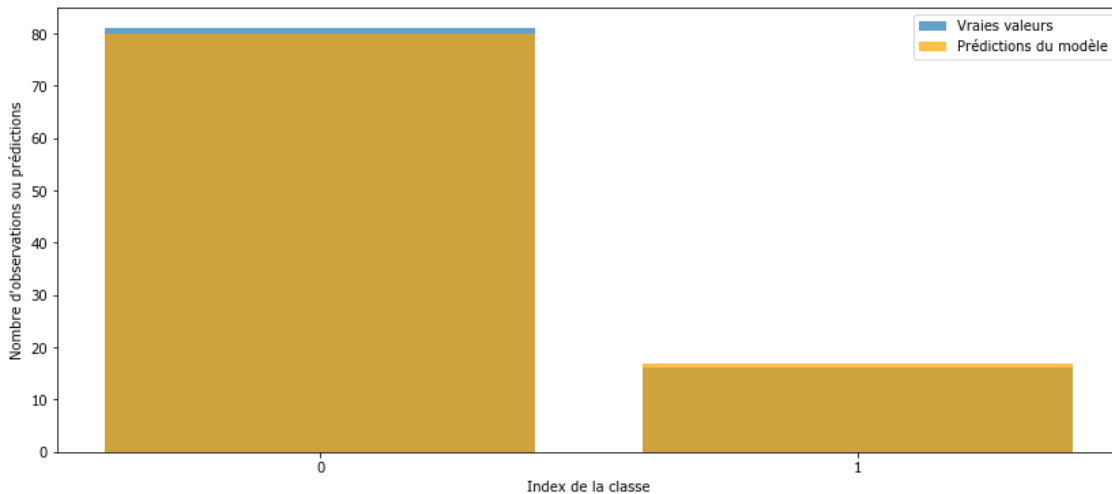


FIGURE 6.22 Diagramme à barres des groupes prédits par le MLP par rapport à ceux prévus par le GMM pour y autonome sur Niv3→Niv30 dans la mine 2

Pour les trajets conventionnels et les trajets autonomes, on enrichit alors l'ensemble de conditions opérationnelles avec les prédictions effectuées par le MLP sur X au complet, comme l'indique notre méthodologie. Tout est alors en place pour s'attaquer réellement à la prédiction de TT.

6.5.3 Prédiction de TT

Dans la présente sous-section, nous allons présenter d’abord l’application de notre modèle de prédiction de TT aux trajets complets puis aux trajets segmentés correspondant à chacun des sectionnements d’itinéraires. On ne s’intéressera alors qu’au LSTM. Nous terminerons par le modèle d’empilement. Il n’y a aucune différence entre la prédiction de TT conventionnels et la prédiction de TT autonomes, il nous suffit de sélectionner leur ensemble respectif de données préparées et prétraitées.

Comme au chapitre 5, nous analyserons cet histogramme et les performances de prédiction associées uniquement après avoir obtenu toutes les prédictions, à la sous-sous-section 6.5.3.4.

6.5.3.1 Prédiction initiale de TT sur l’itinéraire complet

Commençons par générer les prédictions initiales de TT sur l’itinéraire complet. Dans notre contexte, nous n’avons remarqué aucune différence notable avec la succession d’étapes présentée pour la mine 1.

On optimise donc directement les hyperparamètres du LSTM sur les données d’apprentissage. Pour les trajets conventionnels et les trajets autonomes, on aboutit respectivement aux combinaisons optimales d’hyperparamètres figurant dans les tableaux 6.6 et 6.7.

TABLEAU 6.6 Combinaison optimale des hyperparamètres de notre LSTM pour les trajets conventionnels sur l’itinéraire complet Niv3→Niv30

Hyperparamètre	Valeur optimale
$nbUnitsLSTM$	244
$L2_1SLTM$	9×10^{-3}
$L2_2LSTM$	2×10^{-1}
$tailleLotsLSTM$	109
$tauxApprLSTM$	2×10^{-3}

En appliquant ensuite notre LSTM optimisé à X_{test} , on obtient notre premier histogramme de prédiction initiale de TT pour chacun des types de pilotage de HT.

TABLEAU 6.7 Combinaison optimale des hyperparamètres de notre LSTM pour les trajets autonomes sur l’itinéraire complet Niv3→Niv30

Hyperparamètre	Valeur optimale
$nbUnitsLSTM$	174
$L2_1LSTM$	4×10^{-4}
$L2_2LSTM$	6×10^{-2}
$tailleLotsLSTM$	33
$tauxApprLSTM$	2×10^{-3}

6.5.3.2 Prédiction initiale de TT sur les itinéraires sectionnés

Concernant les itinéraires sectionnés, l’approche est la même quel que soit le sectionnement adopté. Nous présenterons ici le cas du sectionnement par segments inter-niveaux. Nous ne nous intéresserons à partir d’ici qu’aux trajets conventionnels, car nous jugeons suffisant de comparer les prédictions de TT conventionnels et autonomes sur l’itinéraire complet uniquement.

Commençons par optimiser les hyperparamètres de notre LSTM. Pour cela, rappelons que l’on applique d’abord directement le modèle de préparation de données à chacun des itinéraires inter-niveaux pour maximiser la quantité de données d’apprentissage. On choisit ensuite l’itinéraire inter-niveaux qui semble le plus représentatif de l’ensemble de ces itinéraires. Ici, on ajoute simultanément un autre critère : il doit s’agir de l’un des itinéraires inter-niveaux les plus fréquemment reconnu. En effet, certains itinéraires inter-niveaux, bien que très représentatifs, sont plus rarement observés en raison de dysfonctionnements successifs des deux balises de changement de niveau qui les délimitent. On optimise puis l’on ajuste immédiatement notre LSTM sur ces trajets inter-niveaux, en s’assurant qu’aucun d’eux ne corresponde aux trajets de test. Pour les itinéraires sectionnés en segments majeurs, précisons qu’un LSTM a été optimisé pour chacun des deux segments majeurs. Nous avons rassemblé les combinaisons optimales des hyperparamètres obtenus pour les itinéraires sectionnés dans le tableau 6.8.

Notre méthodologie indique ensuite, pour chacun des segments, de demander au LSTM optimisé de prédire le TT de l’intégralité des trajets inter-niveaux ayant eu lieu durant chacun des itinéraires complets observés.

TABLEAU 6.8 Combinaison optimale des hyperparamètres de notre LSTM pour les trajets conventionnels sur l'itinéraire sectionné Niv3→Niv30

Segment ou segmentation	Hyperparamètre	Valeur optimale
Inter-niveaux	<i>nbUnitsLSTM</i>	333
	<i>L2_1LSTM</i>	7×10^{-3}
	<i>L2_2LSTM</i>	1×10^{-1}
	<i>tailleLotsLSTM</i>	218
	<i>tauxApprLSTM</i>	1×10^{-3}
Segment majeur 1	<i>nbUnitsLSTM</i>	142
	<i>L2_1LSTM</i>	4×10^{-4}
	<i>L2_2LSTM</i>	6×10^{-2}
	<i>tailleLotsLSTM</i>	480
	<i>tauxApprLSTM</i>	8×10^{-4}
Segment majeur 2	<i>nbUnitsLSTM</i>	410
	<i>L2_1LSTM</i>	4×10^{-3}
	<i>L2_2LSTM</i>	$1,4 \times 10^{-1}$
	<i>tailleLotsLSTM</i>	25
	<i>tauxApprLSTM</i>	1×10^{-3}

C'est ici que rentre en jeu l'importante problématique explicitée au chapitre 5 : étant donné que nous utilisons un LSTM, nous sommes affectés par les données manquantes directement dans nos données d'entrée, les X_i . Nous appliquons alors la stratégie que nous avons présenté, consistant à remplir chronologiquement les valeurs manquantes en prédisant successivement chaque TT. Le temps de calcul de notre modèle de prédiction de TT en est fortement affecté puisque cette tâche complexe requiert plusieurs dizaines de minutes à notre support informatique pour être réalisée entièrement.

Une fois cela fait, il faut ajuster le modèle de régression des résidus sur les prédictions obtenues pour l'ensemble d'apprentissage en les comparant aux $y_{i,appr}$ des trajets complets. Nous disposons de plus de cent observations pour chaque segment inter-niveaux, cette approche est donc valide. Une fois que ce modèle est ajusté, on l'applique aux prédictions de TT inter-niveaux des données de test pour les corriger. On l'applique aussi aux prédictions des données

d'apprentissage issues du LSTM pour corriger toutes les anciennes valeurs vides des $y_{i,appr}$, afin de permettre un meilleur ajustement du modèle suivant.

Passons enfin à l'utilisation du modèle de régression linéaire multiple. Son utilisation implique de disposer de l'intégralité des prédictions de TT inter-niveaux d'un itinéraire complet pour tâcher d'obtenir la prédiction la plus cohérente possible du TT complet. Nous venons justement de les générer et de les corriger de surcroît. Son ajustement sur les données d'apprentissage puis ses prédictions des TT complets de l'ensemble de test sont donc immédiates.

Il est remarquable que, grâce à notre remplissage artificiel des TT inter-niveaux, nous puissions si facilement utiliser le modèle de régression linéaire multiple. En effet, il n'existe que de rares trajets complets pour lesquels nous disposons des détections de la totalité des balises de changement de niveau. Seuls 206 trajets sur 7743 présentant une telle particularité sur l'itinéraire Niv3→Niv30, soit 2,66% des trajets. Il s'agit par ailleurs des seuls trajets qui auraient été détectés par notre modèle de préparation de données si l'on avait gardé la configuration sélectionnée pour la mine 1. La criticité du travail d'adaptation de notre méthodologie est donc d'autant plus mise en exergue par cette dernière remarque.

Concernant les itinéraires sectionnés en deux segments majeurs, une particularité intéressante peut être observée : pour chaque trajet complet, soit les deux TT majeurs sont connus, soit aucun des deux. Nous disposons en effet nécessairement des données de détection des balises A et B lorsqu'un trajet complet est détecté. Seule la balise du niveau 12 peut donc faire défaut, anéantissant alors l'intérêt de sectionner le trajet observé.

6.5.3.3 Application du modèle d'empilement aux prédictions initiales obtenues

Passons enfin à l'application de notre modèle d'empilement. Aucune différence notable ne transparaît comparativement aux étapes de son implémentation pour la mine 1.

On optimise donc simplement le modèle à sélectionner et ses hyperparamètres associés sur les prédictions de TT d'apprentissage de nos modèles de prédiction initiale. Nous n'aurons pas besoin d'afficher les résultats obtenus dans un tableau puisque le modèle finalement retenu est simplement une régression linéaire sans terme constant (i.e. `LinearRegression` configurée avec `fit_intercept = False`).

En appliquant enfin ce modèle d'empilement optimisé à X_{test} conventionnel, on obtient la toute dernière liste de prédiction de TT.

6.5.3.4 Évaluation de la qualité des prédictions de TT et discussion

Toutes les prédictions ont été obtenues, ce qui nous mène directement à débiter notre analyse des résultats obtenus. Mentionnons dès à présent que, concernant les TT conventionnels, nous rassemblerons les valeurs des critères de performance successivement obtenus dans un seul et même tableau à la fin de la présente sous-section. Il s'agit du tableau 6.9. Dans les prochains paragraphes, nous préférons uniquement calculer des réductions/augmentations relatives de ces valeurs pour donner des éléments de compréhension et de comparaison plus explicites.

Intéressons-nous dans un premier temps aux prédictions du LSTM seul sur l'itinéraire complet Niv3→Niv30.

Pour les TT conventionnels, la comparaison entre les histogrammes des TT réels et prédits semble indiquer une bonne cohérence entre les prédictions obtenues et les vraies valeurs (voir fig. 6.23).

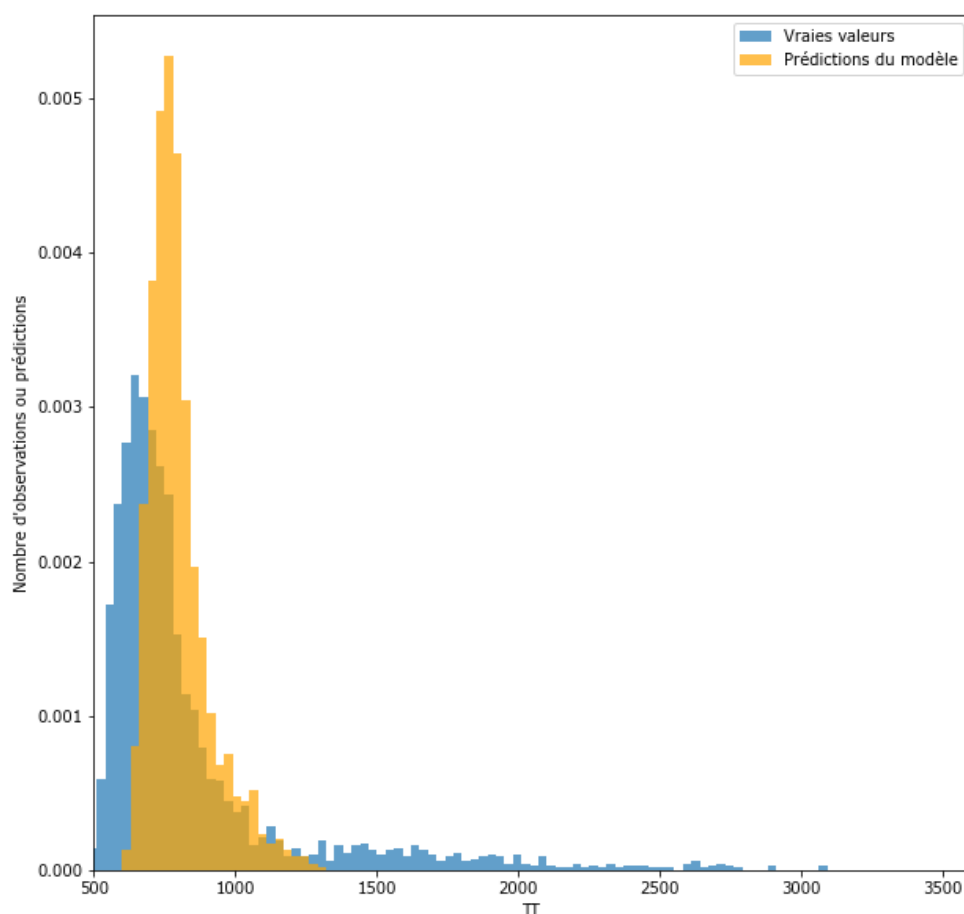


FIGURE 6.23 Comparaison des histogrammes des TT conventionnels réels et prédits par le LSTM seul sur Niv3→Niv30 dans la mine 2

Concernant les TT autonomes, la même remarque pourrait être formulée bien qu'on ne dispose que d'une centaine de trajets de test (voir fig. 6.24).

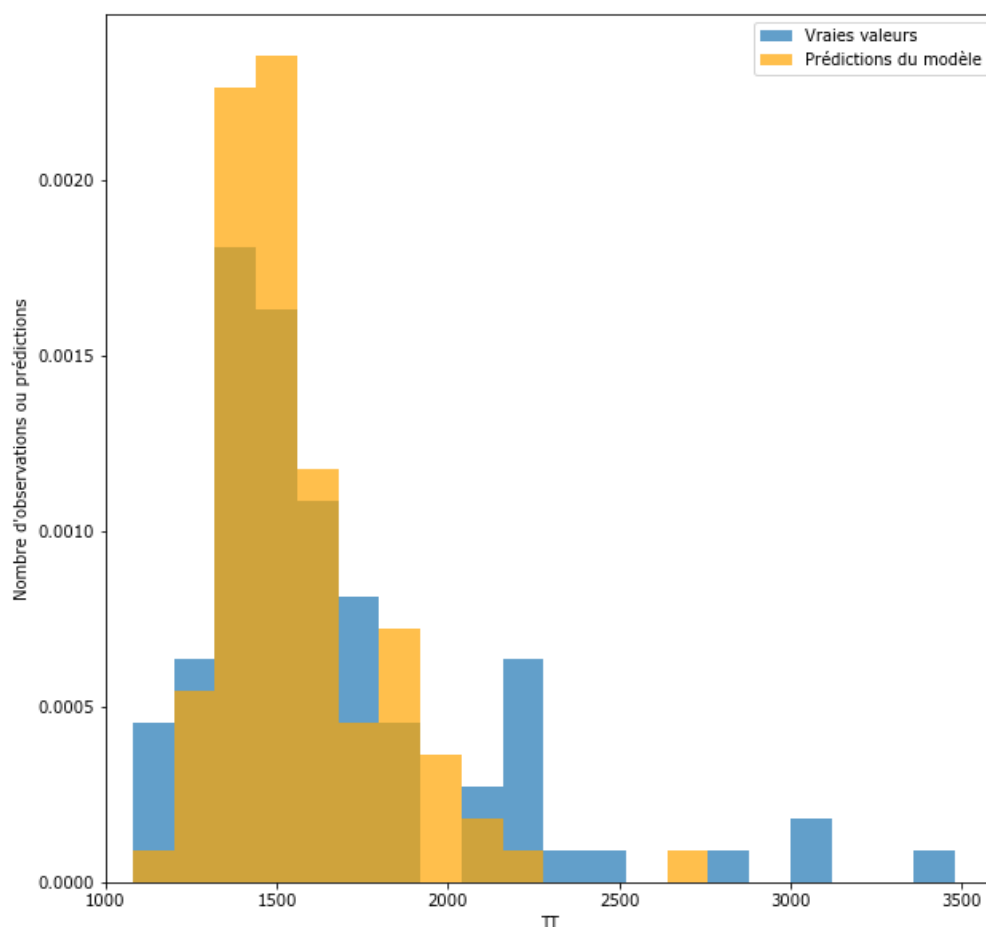


FIGURE 6.24 Comparaison des histogrammes de TT autonomes réels et prédits par le LSTM seul sur Niv3→Niv30 dans la mine 2

Passons maintenant aux critères de performance correspondants.

En s'intéressant uniquement aux critères de performance correspondant aux prédictions de TT conventionnels, un constat s'impose : notre modèle affiche ici de bien meilleures performances que lors de son application aux trajets de la mine 1. En effet, dans cette dernière, le LSTM seul permettait une réduction relative de la MAE d'environ 6,6% (par rapport au modèle de référence) mais une augmentation relative d'environ 1,6% pour la RMSE, qui est pourtant son critère d'optimisation. Dans le cas de la mine 2, la réduction relative de la MAE dépasse les 20% et la réduction relative de la RMSE avoisine les 10,5%. Le calcul de la différence entre ces réductions relatives nous donne respectivement 12 et 14 points de pourcentage pour la MAE et la RMSE. L'écart de précision est donc pour le moins impressionnant.

Concernant les performances de notre modèle pour les TT autonomes, elles sont encore bien meilleures dans l'ensemble. En revanche, elles varient très considérablement en fonction de la sélection aléatoire des ensembles d'apprentissage et de test. Sur de multiples essais, nous avons pu observer une réduction relative de la MAE allant de 12,8% à 38,6%, et une réduction relative de la RMSE allant de 10,7% à 38,7%. Respectivement, ces performances correspondaient à une distribution de test comportant une quantité disproportionnée de très longs trajets et à une distribution de test comportant très peu de très longs trajets. Quelques remarques intéressantes peuvent alors être formulées :

- Au vu des performances exceptionnelles que nous observons ici (en particulier pour la RMSE), les TT autonomes pourraient être plus faciles à prédire que les TT conventionnels grâce à une très grande régularité du système de conduite autonome quel que soit le HT (contrairement aux opérateurs humains, naturellement plus imprévisibles) ;
- Une réduction du nombre de TT très longs, supposés généralement non ordinaires, serait à même de faire exploser la qualité des prédictions de notre modèle. Cette remarque nous permet de réaffirmer la capacité de notre modèle à prédire les « véritables » trajets avec une précision encore accrue. En effet, rappelons encore une fois que les trajets non ordinaires ne sont pas de véritables trajets du point de vue des planificateurs, bien que notre algorithme de reconnaissance de trajets les considère comme tels ;
- Lorsque les performances de notre modèle sont maximales comparativement au modèle de référence, la distribution de test (certes plus facile à prédire en théorie) est considérablement différente de la distribution d'apprentissage puisque cette dernière contient à l'inverse une part disproportionnée de très longs TT. Notre modèle de prédiction gère alors excellemment bien la nouvelle distribution de test, ce qui l'oppose radicalement au modèle de référence qui reste biaisé par les très longs TT d'apprentissage. C'est une très belle démonstration de la robustesse accrue de notre modèle comparativement au modèle de référence dans une telle situation ; et
- Notre LSTM seul vient de démontrer d'excellentes performances sur un itinéraire pour lequel seules quelques centaines de détections sont disponibles. Là aussi, c'est un excellent signe de robustesse. Notre modèle nous convint ainsi qu'il pourrait aussi atteindre de très solides performances s'il était appliqué aux véritables trajets se déroulant sur les itinéraires les plus récents et les plus longs.

Intéressons-nous maintenant aux prédictions de TT conventionnels par itinéraires sectionnés.

Quel que soit le sectionnement adopté (resp. inter-niveaux et segments majeurs), le LSTM seul permet une importante réduction relative de la MAE (resp. 16,1% et 14,2%) mais une bien maigre réduction de la RMSE comparativement aux prédictions sur l'itinéraire au complet

(resp. 2,8% et 3,6%).

Lors de l'application de la régression des résidus aux résultats obtenus, un phénomène intrigant est observé : bien que la qualité des performances augmente nettement dans le cas du sectionnement en segments majeurs, elles vont jusqu'à empirer dans le cas du sectionnement inter-niveaux. Respectivement et comparativement aux prédictions du LSTM seul, la réduction relative de la MAE est alors de 6,4% et -5,2%. La réduction relative de la RMSE atteint quant à elle 5,8% et -1,8%.

Nous supposons alors que certains segments inter-niveaux sont trop peu observés durant les trajets complets. Notons que nous savions d'avance que ce nombre d'observations est toujours supérieur à 321 puisqu'il s'agit du nombre de trajets complets pour lesquels toutes les balises ont détecté le HT correspondant, ce qui nous avait conforté à penser que la régression des résidus fonctionnerait normalement. Après investigation, le nombre de trajets inter-niveaux détectés durant les 7743 trajets complets atteint respectivement 4735, 4648, 4365, 5114, 674, 1122, 4767, 1246 et 1253 pour chacun des neuf segments (de Niv3 vers Niv30). Nous sommes donc amenés à penser qu'entre un (le cinquième segment) et quatre (les segments 5, 6, 8 et 9) de ces segments manquent de détections de trajets pour le modèle de régression des résidus. Celui-ci est alors incapable de généraliser et d'afficher de bonnes performances. Les segments majeurs ne sont pas affectés par cette problématique car la balise sélectionnée au niveau 12 est bien plus fiable. La régression des résidus dispose donc d'un grand ensemble d'apprentissage (5178 trajets ont été observés au total sur chacun des deux segments majeurs durant les 7743 trajets complets, soit 3625 trajets d'apprentissage) qui lui permet de corriger pertinemment les résultats du LSTM sur les trajets de test.

Rappelons que le LSTM est moins exposé aux faibles échantillons évoqués puisqu'il est entraîné sur tous les trajets ayant été reconnus sur chacun des segments, excepté ceux de l'échantillon de test créé pour les trajets complets. Les trajets inter-niveaux d'apprentissage sont respectivement au nombre de 16000, 16200, 12236, 11462, 3680, 4303, 10593, 1602 et 1185 (ce dernier nombre est inférieur à 1253 car on lui a soustrait les trajets de test correspondant aux trajets complets). Ainsi, étant donné que notre LSTM a démontré des performances exceptionnelles lorsqu'il était entraîné sur seulement quelques centaines de trajets autonomes, nous admettons qu'il ne devrait pas être confronté ici à des difficultés notables de prédiction.

Pour les itinéraires sectionnés, il nous reste à présenter l'amélioration des performances permise par le modèle de régression linéaire multiple. Pour respecter la méthodologie présentée au chapitre précédent, nous appliquerons ce modèle aux résultats du modèle de régression des résidus comme si nous ne savions pas que les résultats obtenus ont empiré suite à l'utilisation de cette régression sur les segments inter-niveaux. Naturellement, nous recommandons dès

maintenant aux futurs utilisateurs de notre méthodologie de vérifier aussi pour chaque segment que la régression des résidus améliore les résultats de prédiction et, dans le cas contraire, de ne pas l'appliquer au segment correspondant. Il s'agit là d'une correction mineure tout à fait pertinente.

Fidèle à ses habitudes, la régression linéaire multiple fait quant à elle empirer la MAE pour favoriser une amélioration de la RMSE. La perte de précision sur la MAE est toutefois bien plus importante qu'au chapitre précédent. Relativement aux résultats de prédiction de la régression des résidus, pour le sectionnement inter-niveaux et en segments majeurs, la MAE subit une augmentation relative de 10,7% et de 13,2% respectivement. La RMSE, quant à elle, profite d'une réduction relative bien plus faible : l'amélioration est de 3,9% et de 1,1%, respectivement. Il est difficile de trouver une explication à cette différence importante entre les deux chapitres. Dans le contexte où l'on cherche à minimiser la RMSE (comme cherchent à le faire tous nos sous-modèles), il est difficile de qualifier ces résultats autrement que de « bons » puisqu'il y a eu une amélioration effective du critère de performance central.

Une remarque tout à fait intéressante doit être soulignée, en lien avec les remarques que nous avons formulées précédemment pour le modèle de régression des résidus. Les k_i proposés par la régression linéaire multiple pour les segments inter-niveaux sont en effet les suivants : 1,017, 1,017, 1,023, 1,017, 0,986, 0,574, 1,011, 0,983 et 0,988 (en allant de Niv3 vers Niv30). On remarque immédiatement que tous les k_i prennent des valeurs cohérentes excepté le sixième segment. Il ne s'agit étrangement pas du segment le plus suspect évoqué précédemment (c'était alors le cinquième segment, qui comporte deux fois moins de détections que le sixième). En revanche, c'est une preuve irréfutable que, malgré un coefficient λ_{reg} égal à 10×10^4 pour les trajets inter-niveaux, le modèle de régression linéaire multiple considère qu'il faut absolument minimiser le poids des TT prédits pour ce segment dans les prédictions de TT complets. On peut donc supposer que les prédictions de TT formulées par le modèle de régression des résidus pour ce segment sont effroyablement mauvaises. Le modèle de régression linéaire multiple ne peut en revanche pas réduire ce k_i à zéro car la pénalité serait alors bien trop forte.

Par ailleurs, notre modèle de régression linéaire multiple fixe le terme d'interception à 167,42 afin de compenser l'importante proportion du TT du sixième segment non prise en compte (i.e. le TT moyen du sixième segment multiplié par la différence de 0,574 et 1, soit 0,426).

D'après ces observations, une modification mineure de notre méthode semblerait pertinente : pour le i ème segment, si k_i prend une valeur inférieure à 0,90 on pourrait par exemple remplacer les prédictions de TT du segment par la moyenne des TT observés. Ainsi, les résultats du modèle de régression linéaire multiple ne seront plus affectés par des prédictions

aberrantes massives sur l'un des segments.

Mentionnons aussi que les valeurs prises par les autres k_i indique qu'un coefficient λ_{reg} égal à 10×10^4 est bien adapté à différents sites miniers et qu'il est donc pertinent de l'avoir fixé lors de l'étude de la mine 1.

Enfin, passons à l'analyse des résultats de prédiction du modèle d'empilement. Tout d'abord, on peut dire que ce modèle a l'effet escompté puisque le RMSE de ses prédictions est inférieur à celui de chacun des modèles de prédiction initiale de TT. De plus, il ne semble pas affecté par les mauvaises performances du sous-modèle de prédiction de TT par sectionnement de l'itinéraire en segments inter-niveaux. En revanche, la réduction relative du RMSE est extrêmement faible : elle n'atteint que 0,6% lorsque l'on compare les performances du modèle d'empilement à celles du LSTM seul appliqué à l'itinéraire complet. Dans le même temps, la précision chute au regard du MAE : l'augmentation relative observée est de 9,9%. On représente la comparaison des histogrammes des TT réels et prédits par le modèle d'empilement à la figure 6.25.

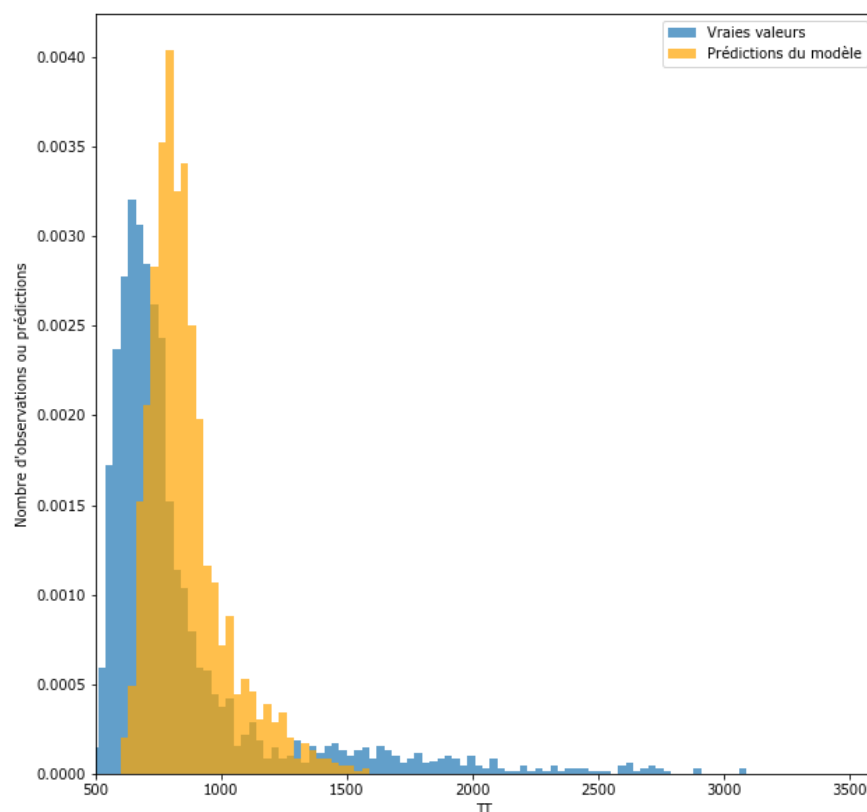


FIGURE 6.25 Comparaison des histogrammes des TT conventionnels réels et prédits par notre modèle global de prédiction sur Niv3→Niv30 dans la mine 2

TABLEAU 6.9 Performances successives de nos sous-modèles de prédiction initiale, du modèle d’empilement et du modèle de référence appliqués aux TT conventionnels

Modèle	Segmentation	Phase	MAE	RMSE
De référence (moyenne)	Itinéraire complet	/	254	392
De prédiction initiale (LSTM)	Itinéraire complet	/	203	351
	Inter-niveaux	LSTM seul	213	381
		Rég. résidus	224	388
		Finale	248	373
	Segments majeurs	LSTM seul	218	378
		Rég. résidus	204	356
		Finale	231	352
D’empilement (rég. linéaire)	Toutes	/	223	349

Au regard de la complexité nécessaire pour obtenir ces ultimes résultats, comparativement aux résultats pouvant être obtenus grâce au LSTM seul appliqué à l’itinéraire complet, il est clair que les performances finalement observées sont décevantes.

Pour autant, la pertinence de notre modèle de prédiction de TT n’est pas remise en question : une amélioration des résultats de prédiction a été observée sur le critère de performance principal.

De plus, les propositions d’amélioration précédentes pourraient être à même d’améliorer la qualité des prédictions initiales fournies au modèle d’empilement et d’ainsi lui permettre de démontrer de meilleures performances.

Nous ne pouvons pas non plus exclure la possibilité que nos modèles aient approché les meilleures prédictions de TT possibles au regard des conditions opérationnelles qui leur ont été rendues disponibles. En effet, le LSTM seul appliqué à l’itinéraire complet permettait déjà d’atteindre quasiment la même RMSE que notre sous-modèle de prédiction initiale de TT par sectionnement de l’itinéraire en segments majeurs (les prédictions de TT inter-niveaux n’auraient jamais pu atteindre un tel score au vu du dysfonctionnement évident observé). Le fait que le modèle d’empilement arrive très difficilement à faire mieux, malgré le fait qu’il

dispose des prédictions des deux sous-modèles les plus performants, irait en ce sens.

Quoi qu'il en soit, nous voilà prêt à conclure ce sixième et dernier chapitre.

6.6 Conclusion

En conclusion, le test de généralisation de notre modèle global sur un second site minier a révélé plusieurs points clés.

Naturellement, pour commencer avec le point le plus critique, il faut dire que notre modèle a su démontrer sans ambiguïté sa capacité à généraliser lors de la préparation des données puis lors de la prédiction de TT conventionnels et autonomes. L'architecture de la mine 2, sa technologie de balises, ses HT autonomes et les contraintes exceptionnelles rencontrées sont autant d'éléments radicalement différents des spécificités de la mine 1. Quelques adaptations laborieuses, décrites dans les chapitres précédents et faisant ainsi partie intégrante de notre méthodologie, ont été nécessaires en raison de données manquantes de balises inter-niveaux. Aussi, en veillant à intégrer minutieusement les différents éléments que nous avons évoqués aux chapitres 4 et 5 pour adapter notre méthodologie aux spécificités de chaque site minier, nous considérons cette dernière comme étant robuste et flexible.

Quelques modifications mineures de notre modèle de prédiction de TT ont été jugées pertinentes à la lumière des difficultés rencontrées par l'un de nos sous-modèles de prédiction initiale de TT. Toutefois, l'existence des autres sous-modèles de prédiction initiale a brillamment permis à notre modèle de prédiction de TT de conserver sa stabilité et d'offrir de très bonnes performances finales comparativement au modèle de référence.

Concernant l'influence des différentes variables opérationnelles, nos analyses statistiques ont confirmé dans le présent chapitre que les variables telles que le nombre de HT actifs, le quart de travail et même le jour de la semaine ont un impact très significatif sur les TT de HT en souterrain. Il est selon nous indispensable de prendre en compte ces facteurs dans un tel contexte afin d'améliorer la précision des prédictions.

Des pistes d'amélioration ont été identifiées, telles que la surveillance systématique des balises pour prévenir les absences de données, l'augmentation de la bande passante pour le transfert des données et l'ajout d'attributs supplémentaires pour mieux distinguer les trajets ordinaires des trajets non-ordinaires.

Les résultats obtenus montrent que la stratégie d'installation et de paramétrage des balises est cruciale pour la reconnaissance précise des trajets et la prédiction des TT. Des implications majeures en découlent sur les meilleures pratiques d'installation et de paramétrage des balises. Nous jugeons par ailleurs que la gestion des données ne doit jamais être négligée.

En somme, ce chapitre a démontré la capacité de notre modèle global à s'adapter et à performer sur un nouveau site minier significativement dissemblable du précédent. Des améliorations supplémentaires ont été soulignées, ouvrant la voie à l'obtention de prédictions de plus haute précision et à une amplification supplémentaire de la stabilité de notre modèle.

CHAPITRE 7 CONCLUSION

7.1 Synthèse des travaux

Ce mémoire a présenté une avancée significative dans le domaine de la prédiction de TT de HT dans les mines souterraines. Il s'agit en particulier de la toute première étude ayant réussi à prédire avec succès les TT de HT en s'appuyant exclusivement sur les lots de données émis par des balises de détection de HT. Il ne fait aucun doute que le modèle de référence couramment utilisé dans l'industrie minière est désormais dépassé par notre méthodologie innovante.

7.1.1 Préparation des données

Nous avons développé une méthodologie rigoureuse pour la préparation des données, qui est essentielle pour garantir la qualité et la pertinence des données fournies aux modèles de prédiction. Elle se décompose comme suit :

- **Stratégie robuste de reconnaissance de trajets ordinaires** : elle permet théoriquement d'identifier les trajets ayant eu lieu sur un itinéraire quelconque d'un site minier quelconque, en se basant sur les détections de balises et en fixant de nombreux critères de filtrage. Cette stratégie a permis de créer de larges ensembles d'observations contenant une proportion importante de véritables trajets ordinaires ;
- **Extraction de conditions opérationnelles pertinentes** : nous avons réussi à extraire des variables qui sont toujours disponibles avant chaque quart de travail. Nous nous sommes ainsi assurés que, lors de l'utilisation réelle de notre méthodologie, les planificateurs puissent disposer d'au moins autant d'informations avec une précision généralement supérieure ;
- **Prétraitement des données** : nous avons détaillé le processus permettant de transformer et de normaliser pertinemment nos données pour nos modèles de prédiction.

7.1.2 Modèles de prédiction

Pour la prédiction de TT de HT dans les mines souterraines, nous avons prouvé successivement, et pour la première fois d'après la littérature existante, la pertinence de plusieurs techniques de ML sophistiquées :

- **Modèle de partitionnement (tandem GMM-MLP)** : nous avons prouvé la pertinence d'utiliser un tel modèle en amont de nos modèles de prédiction de TT. La combinaison d'un GMM et d'un MLP est une technique de partitionnement complexe mais nécessaire pour améliorer les résultats du modèle global de prédiction ;
- **RNN (LSTM)** : La pertinence d'un LSTM pour formuler les prédictions de TT a été démontrée. Le LSTM s'est révélé capable de capturer les dynamiques temporelles complexes des trajets dans les mines souterraines, surpassant généralement les performances des modèles traditionnels et montrant une robustesse bien supérieure ;
- **Modèle d'empilement (GBR ou régression linéaire)** : Nous avons montré l'intérêt de combiner plusieurs prédictions via différentes segmentations en utilisant un modèle d'empilement final optimal. Cette approche a permis de renforcer la robustesse et la précision des prédictions.

7.1.3 Application à deux sites miniers

Notre méthodologie a été appliquée avec succès à deux sites miniers distincts, démontrant sa robustesse et sa flexibilité :

- **Mine 1** : ce site nous a permis de mettre au point notre modèle global. Nous avons mené à bien une préparation de données minutieuse et une prédiction de TT remarquable malgré de nombreux défis posés par des réglages inadaptées des balises de détection de HT. Sur ce site minier, les résultats obtenus ont démontré la supériorité de notre modèle par rapport au modèle de référence de l'industrie. Aucune absence importante de données n'était à déclarer ; et
- **Mine 2** : sur ce second site, nous avons pu mener notre test de généralisation du modèle global. La capacité de notre modèle à s'adapter à des paramètres et un contexte différent a alors été mise en évidence. Des adaptations préalablement envisagées ont été nécessaires en raison d'absences massives de données de détection de balises. Malgré tout, notre modèle y a clairement redémontré sa robustesse en affichant des résultats très satisfaisants tant pour des prédiction de TT conventionnels que des TT autonomes, auxquels il n'avait jamais été exposé jusqu'alors.

7.2 Limitations de la solution proposée

Bien que notre modèle ait démontré des performances remarquables, certaines limitations subsistent :

- **Qualité des données** : la basse qualité des données de détection de balises, ainsi que certaines contraintes spécifiques à chaque site minier, sont à même d'introduire de redoutables défis additionnels parfois insolubles durant la phase de préparation des données ;
- **Trajets indésirables** : les trajets non ordinaires ou un mélange de trajets conventionnels et autonomes, très difficilement filtrables sans sacrifier une part considérable des trajets ordinaires, défavorise assurément notre modèle de prédiction ; et
- **Données manquantes** : les absences de données de balises inter-niveaux entraînent des adaptations laborieuses des LSTM et des modèles de régression des résidus et dégradent nécessairement la précision des prédictions de notre modèle de prédiction de TT.

7.3 Améliorations futures

Pour surmonter ces limitations et améliorer encore notre méthodologie, plusieurs axes d'amélioration ont été identifiés :

1. **Surveillance de la santé des données de détection** : il s'agirait de mettre en place une surveillance systématique des balises pour assurer une collecte de données de haute qualité et prévenir du même coup les absences massives de données ;
2. **Ajout d'attributs** : nous proposons d'intégrer des attributs supplémentaires des BDD opérationnelles pour mieux distinguer les trajets ordinaires des trajets non ordinaires, améliorant ainsi la précision des prédictions ;
3. **Discrimination des trajets ordinaires et non ordinaires** : il est certain que des efforts majeurs sont requis pour pleinement y arriver. Toutefois, une bien meilleure exclusion de ces trajets pourrait déjà être atteinte en demandant à des experts des activités de HT de chaque site de fixer des seuils de filtrage adéquats pour notre modèle de préparation de données ;
4. **Augmentation de la bande passante** : cette amélioration, à première vue mineure, pourrait être critique pour éviter des pertes massives de données de détection lors de leur transfert ; et
5. **Pistes de recherche concernant le modèle de prédiction** : les nombreuses pistes de recherche que nous avons évoquées concernant la prédiction de TT pourront permettre de continuer à développer et à optimiser notre modèle. Il est nécessaire de bien prendre en compte les spécificités de chaque site minier et d'évaluer l'impact de l'introduction de nouvelles variables de prédiction sur ses performances.

En résumé, le présent mémoire a établi une base solide pour la prédiction de TT de HT dans les mines souterraines via l'utilisation de systèmes de balises de détection de HT. Grâce à une méthodologie rigoureuse et à des techniques avancées, nous avons démontré la faisabilité et l'efficacité d'une approche innovante capable de surpasser le modèle de référence actuel. Les défis complexes du contexte opérationnel minier souterrain ont été pleinement adressés et une solution spécifique et convaincante a ainsi pu émerger. Elle ouvre la voie à une liste substantielle d'améliorations majeures qui, combinées, pourraient produire des résultats de prédiction d'une précision exceptionnelle.

RÉFÉRENCES

- [1] La langue française, “Définition de secteur primaire,” 2023, [Page disponible le 20-septembre-2023]. [En ligne]. Disponible : <https://www.lalanguefrancaise.com/dictionnaire/definition/secteur-primaire>
- [2] C. Mion et M. Ajmi, “Les 10 principaux risques et opportunités business pour le secteur minier en 2023,” 2022. [En ligne]. Disponible : https://www.ey.com/fr_fr/mining-metals/principaux-risques-e-opportunités-du-secteur-minier-en-2023
- [3] V. Bellehumeur *et al.*, “TRANSFORMATION NUMÉRIQUE ET COMPÉTENCES DU 21^e SIÈCLE POUR LA PROSPÉRITÉ DU QUÉBEC,” 2018. [En ligne]. Disponible : <https://numerique.banq.qc.ca/patrimoine/details/52327/4027329>
- [4] M. Normand *et al.*, *PORTRAIT NUMÉRIQUE DE L’INDUSTRIE MINIÈRE AU QUÉBEC*. Bibliothèque et Archives nationales du Québec, 2019. [En ligne]. Disponible : <https://numerique.banq.qc.ca/patrimoine/details/52327/4027329>
- [5] R. Leite-Corthésy, “Vérification d’un plan de production minier à court terme par simulation,” Thèse de doctorat, École Polytechnique de Montréal, 2016.
- [6] G. L’Heureux, “Modèle d’optimisation pour la planification à moyen terme des mines à ciel ouvert,” Thèse de doctorat, École Polytechnique de Montréal, 2011.
- [7] L. St-Amant Dyotte, “Développement d’une solution opérationnelle 4.0 pour le transport minier souterrain,” Thèse de doctorat, École de technologie supérieure, 2022.
- [8] Office québécois de la langue française, “mine,” 1990. [En ligne]. Disponible : <https://vitrinelinguistique.oqlf.gouv.qc.ca/fiche-gdt/fiche/17027685/mine>
- [9] Merriam-Webster, “Hauler,” [En ligne; Page disponible le 10 octobre 2023]. [En ligne]. Disponible : <https://www.merriam-webster.com/dictionary/hauler>
- [10] BELAZ, “Mining dump trucks,” [En ligne; Page disponible le 10 octobre 2023]. [En ligne]. Disponible : <http://belaz.ca/products/dumptrucks>
- [11] SKF, “Solutions for dump trucks in mining,” [En ligne; Page disponible le 10 octobre 2023]. [En ligne]. Disponible : <https://www.skf.com/my/industries/mining-mineral-processing-cement/mining-quarrying/haul-trucks>
- [12] CAT, “Underground mining load haul dump (lhd) loaders,” [En ligne; Page disponible le 11 octobre 2023]. [En ligne]. Disponible : https://www.cat.com/en_US/products/new/equipment/underground-hard-rock/underground-mining-load-haul-dump-lhd-loaders.html
- [13] Queen’s Mine Design Wiki, “Equipment selection,” 2015, [En ligne; Page disponible le 11 octobre 2023]. [En ligne]. Disponible : http://minewiki.engineering.queensu.ca/mediawiki/index.php/Equipment_selection#:~:text=Haul%20trucks%20tend%20to%20have,for%20both%20mucking%20and%20hauling.
- [14] Mining Technology, “Mining trucks, loaders and haulage equipment for the mining industry,” [En ligne; Page disponible le 10 octobre 2023]. [En ligne]. Disponible : <https://www.mining-technology.com/buyers-guide/mining-loaders-trucks-haulage/>

- [15] N. Li *et al.*, “Underground mine truck travel time prediction based on stacking integrated learning,” *Engineering Applications of Artificial Intelligence*, vol. 120, p. 105873, 2023. [En ligne]. Disponible : <https://www.sciencedirect.com/science/article/pii/S095219762300057X>
- [16] K. Hlophe, “GPS-deprived localisation for underground mines,” 2010. [En ligne]. Disponible : <https://researchspace.csir.co.za/dspace/handle/10204/4225>
- [17] E. K. Chanda et S. Gardiner, “A comparative study of truck cycle time prediction methods in open-pit mining,” *Engineering, Construction and Architectural Management*, vol. 17, n^o. 5, p. 446 – 460, 2010. [En ligne]. Disponible : <https://www.scopus.com/inward/record.uri?eid=2-s2.0-77956620853&doi=10.1108%2f09699981011074556&partnerID=40&md5=f318284b4cab6b0abf69acef338819e7>
- [18] C. Fan *et al.*, “Weighted ensembles of artificial neural networks based on gaussian mixture modeling for truck productivity prediction at open-pit mines,” *Mining, Metallurgy and Exploration*, vol. 40, n^o. 2, p. 583 – 598, 2023. [En ligne]. Disponible : <https://www.scopus.com/inward/record.uri?eid=2-s2.0-85149141065&doi=10.1007%2fs42461-023-00747-9&partnerID=40&md5=6176e039bd81bb060376a8ff29afb6fa>
- [19] M. Ao, C. Li et S. Yang, “Prediction method of truck travel time in open pit mines based on lstm model,” 07 2023, p. 8651–8656.
- [20] P. Gun, A. J. Hill et R. Vujanic, “Coordinating multiple cooperative vehicle trajectories on shared road networks,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 24, n^o. 1, p. 274–290, 2023.
- [21] C. Fan *et al.*, “Prediction of truck productivity at mine sites using tree-based ensemble models combined with gaussian mixture modelling,” *International Journal of Mining, Reclamation and Environment*, vol. 37, n^o. 1, p. 66 – 86, 2022. [En ligne]. Disponible : <https://www.scopus.com/inward/record.uri?eid=2-s2.0-85141373250&doi=10.1080%2f17480930.2022.2142425&partnerID=40&md5=98593edf4c26d48f7acbc2927cbbc1ae>
- [22] —, “Preprocessing large datasets using gaussian mixture modelling to improve prediction accuracy of truck productivity at mine sites,” *Archives of Mining Sciences*, vol. 67, n^o. 4, p. 661 – 680, 2022. [En ligne]. Disponible : <https://www.scopus.com/inward/record.uri?eid=2-s2.0-85141413665&doi=10.24425%2fams.2022.143680&partnerID=40&md5=e93223b136028a7ff9b145b3a5a42ed8>
- [23] S. Choudhury et H. Naik, “Use of machine learning algorithm models to optimize the fleet management system in opencast mines,” 2022. [En ligne]. Disponible : <https://www.scopus.com/inward/record.uri?eid=2-s2.0-85135618019&doi=10.1109%2f12CT54291.2022.9825450&partnerID=40&md5=f0c3e760e1c26c9f04886011d6eae48>
- [24] S. Upadhyay *et al.*, “A simulation model for estimation of mine haulage fleet productivity,” 01 2020, p. 42–50.
- [25] J. Baek et Y. Choi, “Simulation of truck haulage operations in an underground mine using big data from an ICT-based mine safety management system,” *Appl. Sci. (Basel)*, vol. 9, n^o. 13, p. 2639, juin 2019.
- [26] X. Sun *et al.*, “The use of a machine learning method to predict the real-time link travel time of open-pit trucks,” *Mathematical Problems in Engineering*, 2018. [En ligne]. Dis-

- ponible : <https://www.scopus.com/inward/record.uri?eid=2-s2.0-85048666422&doi=10.1155%2f2018%2f4368045&partnerID=40&md5=e0c9457a517caeae0dfd69c34ce447f2>
- [27] K. Ristovski *et al.*, “Dispatch with confidence : Integration of machine learning, optimization and simulation for open pit mines,” dans *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. New York, NY, USA : Association for Computing Machinery, 2017, p. 1981–1989. [En ligne]. Disponible : <https://doi.org/10.1145/3097983.3098178>
 - [28] S. Upadhyay *et al.*, “Simulation and optimization in open pit mining,” 05 2015.
 - [29] K. Erarslan, “Modelling performance and retarder chart of off-highway trucks by cubic splines for cycle time estimation,” *Transactions of the Institutions of Mining and Metallurgy, Section A : Mining Technology*, vol. 114, n°. 3, p. A161–A166, 2005. [En ligne]. Disponible : <https://www.scopus.com/inward/record.uri?eid=2-s2.0-27844597136&doi=10.1179%2f037178405X54006&partnerID=40&md5=31cee47d1efd097f8542e323c0e1baf2>
 - [30] V. Temeng, *A Computerized Model for Truck Dispatching in Open Pit Mines*. Michigan Technological University, 1997. [En ligne]. Disponible : <https://books.google.ca/books?id=D0lQNwAACAAJ>
 - [31] R. Grosse et N. Srivastava, “Lecture 16 : Mixture models,” [Page disponible le 7-novembre-2023]. [En ligne]. Disponible : https://www.cs.toronto.edu/~rgrosse/csc321/mixture_models.pdf
 - [32] N. B. M. A. Julnasir, “Full simulation of cycle time using talpac software at lafarge kanthan quarry,” 2018. [En ligne]. Disponible : http://eprints.usm.my/53135/1/Full%20Simulation%20Of%20Cycle%20Time%20Using%20Talpac%20Software%20At%20Lafarge%20Kanthan%20Quarry_Nurwahidah%20Meor%20Akil%20Julnasir_B1_2018.pdf
 - [33] S. Dyer et J. Dyer, “Cubic-spline interpolation. 1,” *IEEE Instrumentation & Measurement Magazine*, vol. 4, n°. 1, p. 44–46, 2001.
 - [34] “RPMGlobal adds to simulation suite,” mars 2020. [En ligne]. Disponible : <https://www.miningmagazine.com/simulation-optimisation/news/1382423/rpmglobal-adds-to-simulation-suite>
 - [35] L. Blessing et A. Chakrabarti, *DRM, a Design Research Methodology*. Springer, 01 2009.
 - [36] M. L. Guen, “La boîte à moustaches de TUKEY, un outil pour initier à la Statistique,” *Statistiquement Votre - SFDS*, 2001, [Page disponible le 16-février-2024]. [En ligne]. Disponible : <https://shs.hal.science/halshs-00287697>
 - [37] J.-H. Guay, “Utiliser le bootstrap pour estimer des intervalles de confiance pour différentes statistiques,” 2017. [En ligne]. Disponible : <https://dimension.usherbrooke.ca/pages/76>
 - [38] R. Hyndman et G. Athanasopoulos, *Forecasting : principles and practice*, 3^e éd. OTexts : Melbourne, Australia, 2021. [En ligne]. Disponible : <https://otexts.com/fpp3>
 - [39] T. Akiba *et al.*, “Optuna : A next-generation hyperparameter optimization framework,” dans *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2019.