



TÉLÉCOM PHYSIQUE STRASBOURG

RAPPORT DE PROTECTION DES DONNÉES

Analyse de données cyberphysiques sur la distribution d'eau

Nathan CERISARA

Clément DESBERG

Ysée JACQUET

Lucas LEVY

21 janvier 2026

Table des matières

1	Données	1
1.1	Visualisation des données	1
1.1.1	Présentation générale des données	1
1.1.2	Analyse des données	2
1.2	Prétraitement	2
2	Application des algorithmes	2
2.1	KNN	2
2.2	Random Forest	2
2.3	XGBoost	2
2.4	MLP	2
2.5	Modèles de transformers	2
2.5.1	Tab Transformer	2
2.5.2	FT Transformer	2
2.5.3	MLP avec attention	2
3	Évaluation et comparaison avec les données de référence	2
3.1	Données physiques	2
	Références	5
	Références	5

Table des figures

1	Matrice de corrélation des erreurs entre les modèles	3
---	--	---

Abréviations employées

- **CART** : Classification And Regression Trees
- **CNN** : Convolutional Neural Network
- **CSV** : Comma Separated Values
- **FT Transformer** : Feature Tokenizer Transformer
- **KNN** : K-Nearest Neighbors
- **MLP** : Multi-Layer Perceptron
- **MCC** : Matthews Correlation Coefficient
- **NLP** : Natural Language Processing

Introduction

Ce projet consiste en l'analyse de données cyberphysiques issues d'un système de distribution d'eau potable [1]. L'objectif principal est de développer et d'évaluer des modèles de machine learning capables de détecter et de classer les anomalies dans le système, telles que les attaques informatiques ou les défaillances physiques. Parmi les modèles demandés, K-Nearest Neighbors (KNN), Random Forest, XGBoost et Multi-Layer Perceptron (MLP) ont été implémentés. De plus, des modèles basés sur des architectures de transformers adaptées aux données tabulaires ont été testés en remplacement de l'algorithme de Classification And Regression Trees (CART). Finalement, les performances des modèles développés ont été comparées avec celles fournies dans l'article de référence [2].

1 Données

L'ensemble des données comprend 5 fichiers Comma Separated Values (CSV) sur les relevés physiques et 5 fichiers CSV sur les relevés réseau. Dans chaque cas, il y a un fichier de mesures prises en période normale – sans anomalie – nommés respectivement `phy_norm.csv` et `normal.csv`, et 4 fichiers de mesures prises en période anormale, avec des anomalies de différents types (attaques informatiques ou défaillances physiques).

1.1 Visualisation des données

1.1.1 Présentation générale des données

Dans chaque fichier CSV, une ligne représente un ensemble de données collectées à un instant donné. Les colonnes renvoient à différents mécanismes sur lesquels une mesure a été faite. Elles sont présentées dans la [Table 1](#).

Nom	Type	Description
feature_1	float	Description de la feature 1
feature_2	int	Description de la feature 2
...
label	string	Type d'anomalie (ou normal)

TABLE 1 – Description des features présentes dans les fichiers CSV

1.1.2 Analyse des données

1.2 Prétraitement

2 Application des algorithmes

2.1 KNN

2.2 Random Forest

2.3 XGBoost

2.4 MLP

2.5 Modèles de transformers

2.5.1 Tab Transformer

2.5.2 FT Transformer

2.5.3 MLP avec attention

3 Évaluation et comparaison avec les données de référence

3.1 Données physiques

D'après les résultats présentés dans la [Table 2](#), le modèle `small_knn` obtient les meilleures performances sur le jeu de données `physical_small`, avec une précision de **97,74%**, un F1-macro de **0,8007**, une précision équilibrée de **0,8235** et un MCC de **0,9402**. Le temps d'entraînement est également très faible (quasi-instantané), ce qui en fait un choix efficace pour ce type de données. À l'inverse, le modèle `small_random_forest` affiche les performances les plus faibles, ce qui peut s'expliquer par une répartition des classes TO COMPLETE.

Rang	Modèle	Précision	F1 (macro)	Précision équilibrée	MCC	Temps (s)
1	<code>small_knn</code>	0.9774	0.8007	0.8235	0.9402	0.0
2	<code>small_tab_transformer</code>	0.9683	0.8001	0.8250	0.9188	40.7
3	<code>small_attention_mlp</code>	0.9048	0.7259	0.8106	0.7946	35.2
4	<code>small_mlp</code>	0.8786	0.6889	0.7988	0.7515	22.7
5	<code>small_ft_transformer</code>	0.8401	0.6593	0.7970	0.7004	56.1
6	<code>small_xgboost</code>	0.8676	0.6412	0.7906	0.7287	0.6
7	<code>small_random_forest</code>	0.6101	0.4687	0.7242	0.4779	0.2

TABLE 2 – Comparaison des expériences sur le jeu de données physiques `physical_small`

En comparant les résultats obtenus avec ceux de l'article de référence [2], on constate que nos modèles TO COMPLETE. En effet dans la Table 3, on peut voir que les performances des KNN, Random Forest, SVM et Naive Bayes (NB) sont TO COMPLETE.

Algorithme	Exactitude	Rappel	Précision	F1-score
KNN	0,98	0,95	0,95	0,95
RF	0,99	0,98	0,95	0,97
SVM	0,93	0,92	0,64	0,75
NB	0,93	0,92	0,66	0,77

TABLE 3 – Résultats de l'évaluation des algorithmes d'apprentissage automatique sur le jeu de données physique

En s'attardant sur les erreurs commises par les différents modèles, on peut observer certaines corrélation entre celles-ci, comme le montre la Figure 1. Par exemple — et sans surprise — le modèle MLP et le modèle MLP avec attention présentent une forte corrélation dans leurs erreurs.

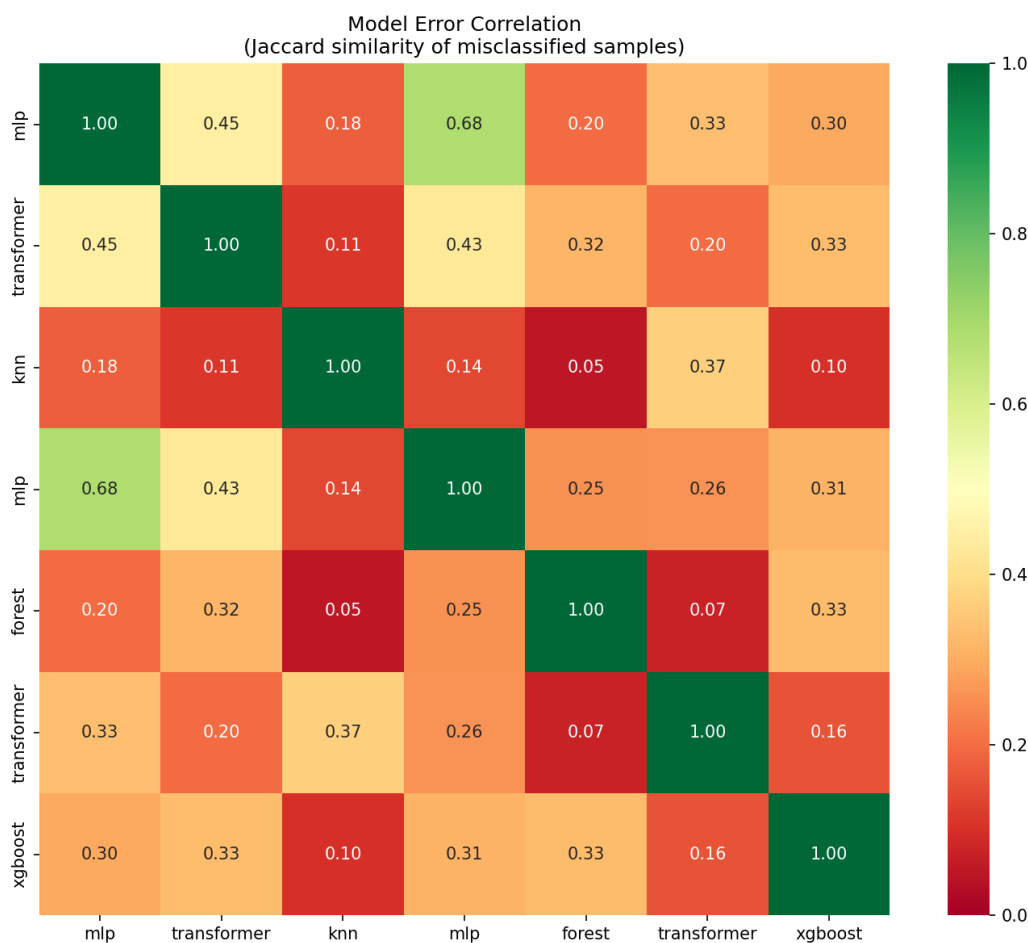


FIGURE 1 – Matrice de corrélation des erreurs entre les modèles

Comme il y a très peu de faible corrélations, il est raisonnable de penser que certains échantillons sont systématiquement mal classés par tous les modèles. La Table 3 liste les

échantillons qui ont été mal classés par l'ensemble des modèles lors de nos différentes exécutions, ainsi que la prédiction majoritaire effectuée par les modèles pour ces échantillons. Heureusement, on constate que leur nombre est très faible — seulement 17 échantillons sur un total de X — donc il n'y a pas beaucoup d'entrées réellement problématiques.

On remarque que la majorité des échantillons mal classés appartiennent à la classe **normal**, ce qui indique que les modèles ont du mal à distinguer les échantillons normaux des anomalies dans certains cas. Cela peut s'expliquer par le fait que les mesures sont effectuées toutes les secondes, et donc la limite entre un état normal et un état anormal peut être très fine. Il y a également des échantillons mal classés appartenant à la classe **scan** qui est sous-représentée, ce qui peut expliquer les difficultés rencontrées par les modèles pour les classer correctement. De même, l'étiquette **scan** est plusieurs fois prédite de manière erronée, ce qui montre bien que les modèles ont du mal à savoir à quoi correspond cette classe.

Entrée	Label réel	Nombre d'erreurs	Prédiction majoritaire
5	normal	7	MITM
35	normal	7	MITM
36	scan	7	normal
210	normal	7	MITM
466	normal	7	DoS
598	normal	7	scan
768	normal	7	physical fault
848	normal	7	physical fault
864	normal	7	scan
892	normal	7	DoS
927	normal	7	MITM
1104	normal	7	MITM
1203	normal	7	scan
1257	normal	7	physical fault
1549	scan	7	normal
1559	normal	7	physical fault
1628	normal	7	MITM

TABLE 4 – Échantillons systématiquement mal classés lors des différentes exécutions

Conclusion

Références

- [1] Simone GUARINO et al. *A hardware-in-the-loop water distribution testbed (WDT) dataset for cyber-physical security testing*. 2021. DOI : [10.21227/rbvf-2h90](https://doi.org/10.21227/rbvf-2h90). URL : <https://dx.doi.org/10.21227/rbvf-2h90>.

- [2] Li TAO et al. « Reduction of Intercarrier Interference Based on Window Shaping in OFDM RoF Systems ». In : *IEEE Photonics Technology Letters* 25.9 (mai 2013), p. 851-854. DOI : [10.1109/LPT.2013.2252335](https://doi.org/10.1109/LPT.2013.2252335). URL : <https://ieeexplore.ieee.org/document/9526562/>.