

---

**Proposition de features pour la détection de moteur  
défectueux à partir de l'enregistrement de leurs  
émissions sonores lors de phases de fonctionnement**

---

Octave Jeanne et Elliot Merle

Année de Recherche en Intelligence Artificielle  
Module - Signal  
Ecole Normale Supérieure Paris Saclay

# Table des matières

Introduction . . . . .	2
Présentation du projet . . . . .	2
Présentation de la base de données utilisée . . . . .	2
Pré-traitement des données . . . . .	2
Débruitage . . . . .	2
Fréquence d'échantillonnage . . . . .	4
Données et Démarche . . . . .	5
Largeur de la bande d'émission spectrale . . . . .	11
Profil d'autocorrélation . . . . .	14
Solution choisie . . . . .	17
Sélection des features intéressantes pour la classification . . . . .	17
Résultats . . . . .	18
Visualisation des données dans l'espace des features sélectionnées . . . . .	18
Clustering . . . . .	18
Conclusion . . . . .	21
Annexe : Distributions des features retenues . . . . .	23
Annexe : Distributions des features éliminées . . . . .	25

## Introduction

### Présentation du projet

Le traitement des séries temporelles représente un enjeu fondamental dans de nombreux domaines tels que l'industrie, la recherche scientifique, la finance ou encore le secteur de la santé. En permettant l'analyse et l'interprétation des données brutes recueillies, il facilite l'identification de motifs récurrents, l'extraction de caractéristiques pertinentes et la détection de ruptures ou d'anomalies. Par exemple, dans le domaine médical, le traitement des signaux électrocardiographiques permet de diagnostiquer des troubles cardiaques, tandis que dans la finance, l'analyse des fluctuations des marchés boursiers aide à anticiper les tendances économiques. De même, dans l'industrie où l'analyse des données mesurées par différent capteur lors de phase de fonctionnement d'un produit permet de faire de la maintenance préventive ou encore du contrôle qualité.

Ce mini projet s'ancre dans le cadre de l'analyse de séries temporelles pour le contrôle qualité. Plus précisément, nous chercherons à séparer des sons issus de l'enregistrement des bruits de moteurs en fonctionnement en deux classes : défectueux ou non.

### Présentation de la base de données utilisée

Les bruits de moteurs utilisés lors de ce projet sont issus d'une base de données recueillie par Ford dans le cadre d'une compétition faite au IEEE World Congress on Computational Intelligence en 2008. Cette base de donnée contient deux jeux de données.

- Le jeu d'entraînement contient 3636 données prises dans un environnement calme.
- Le jeu de test contient 810 données prises dans un atelier classique, donc avec potentiellement du bruit parasite venant des autres travaux faits dans l'atelier.

Chaque donnée est constituée d'un signal de 500 points et un numéro de classe, pour utiliser des méthodes supervisées. Cependant, ce mini-projet est réalisé sous la contrainte de ne pas connaître la classe des données utilisées. Par conséquent, la classe de chaque donnée ne sera ni observée, ni utilisée lors de nos travaux.

## Pré-traitement des données

### Débruitage

La description des bases de données indique que le jeu de test était susceptible de présenter des bruits parasites. Nous commençons donc par évaluer ce bruit en visualisant la superposition des périodogrammes des signaux de nos jeux de données [1].

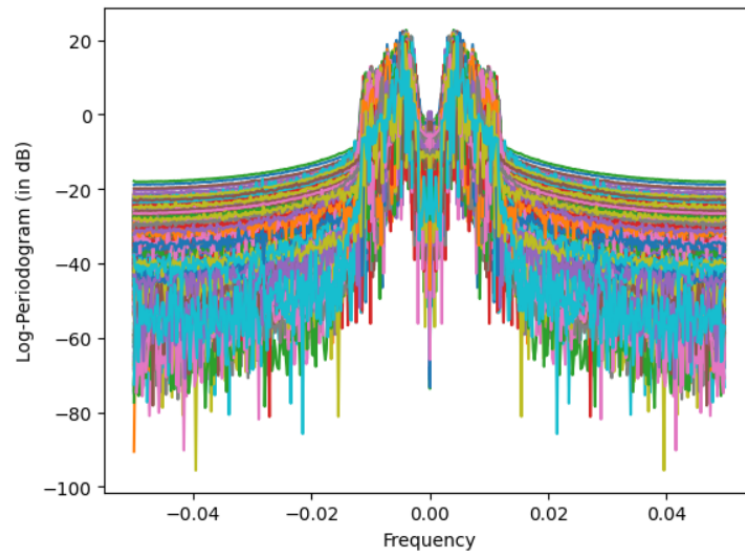


FIGURE 1 – Log-Periodogram de l'ensemble des signaux de la base de données réalisé pour  $f_e = 1\text{Hz}$

L'analyse des périodogrammes montre que les signaux sont décalés verticalement les uns par rapport aux autres, caractéristique typique des signaux bruités par des bruits blancs additifs gaussiens. De fait nous avons cherché à réaliser une opération de débruitage. Sur le log-périodogramme de l'ensemble des signaux mais aussi sur les Log-périodogramme individuels, nous avons observé que le bruit était majoritaire pour toutes les fréquences supérieures à 10-15 Hz. Ainsi, nous avons réalisé une opération de débruitage avec un filtre Passe-Bas avec une fréquence de coupure à 15 Hz (En supposant  $f_e = 1\text{Hz}$ ).

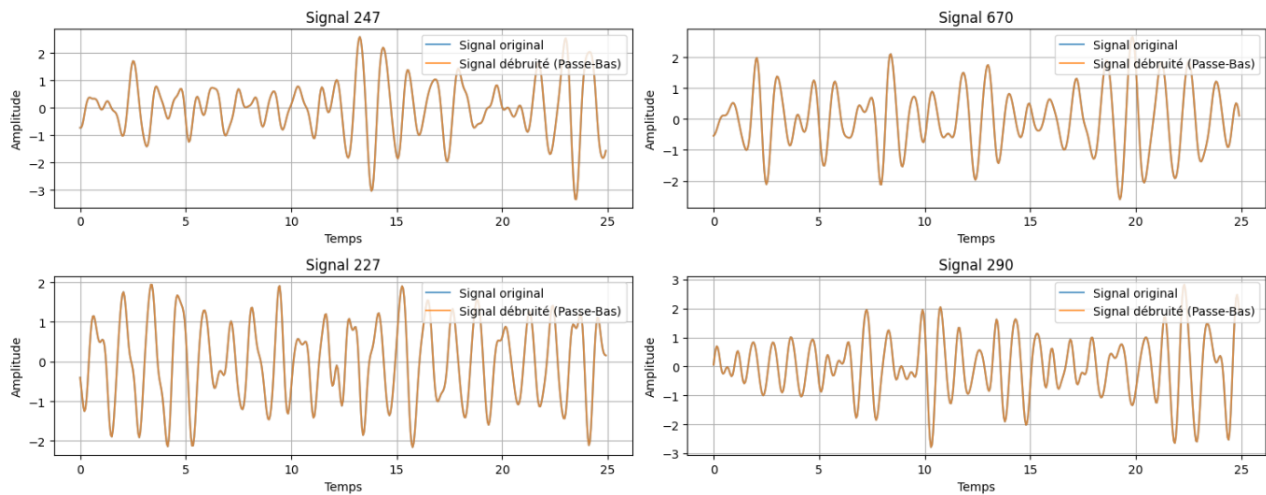


FIGURE 2 – Comparaison des signaux originaux et dé-bruités

La figure 3 montre que les signaux débruités suivent exactement la forme des signaux originaux, ce qui suggère l'absence de bruit dans les signaux. Ce constat est surprenant, car la figure 2 laissait supposer un bruitage, indiqué par un décalage vertical des log-périodogrammes. De plus, la description du jeu de données précise que le fichier "TEST" contient des signaux bruités. Enfin, ce résultat semble s'appliquer à l'ensemble de la base de données, l'expérience ayant été répétée plusieurs fois avec des signaux sélectionnés aléatoirement.

Suite à ces observations, nous avons fait le choix de ne pas appliquer de dé-bruitage pour le reste de notre étude.

## Fréquence d'échantillonnage

La fréquence d'échantillonnage n'est pas fournie avec la base de donnée. Cependant, comme les signaux considérés sont des mesures de son, le critère de Shannon-Nyquist appliqué à la fréquence maximale audible par l'humain permet d'affirmer que cette fréquence d'échantillonnage est d'au moins à 40kHz. De plus, les mesures ont été prises dans le cadre d'une compétition lors d'un congrès, ce qui n'est pas un enjeu majeur pour Ford. On peut donc s'attendre à ce que les micros ayant été utilisé ne sont pas des micros coûteux permettant d'échantillonner à haute fréquence mais des micros standard dont la fréquence d'échantillonnage est comprise entre 44 kHz (disque de musique) et 48kHz (audio DVD). (cf Wikipédia). Nous ferons donc le choix de supposer que la fréquence d'échantillonnage est de 46kHz pour la suite du projet.

On peut alors raisonnablement affirmer que les données ont été faite en découpant des enregistrement plus long. En effet, compte tenu de notre estimation de la fréquence d'échantillonnage et de la longueur des signaux fournit chaque donnée dure quelques centièmes de seconde, ce qui est une durée extrêmement faible pour faire des mesures en une prise. Du fait de l'obtention des différentes données, on s'attend à une certaine cohérence entre celles-ci, en particulier on peut supposer que les phénomènes caractéristiques des bruits traduisant des défauts se ressembleront.

## Données et Démarche

### Visualisation des signaux temporels

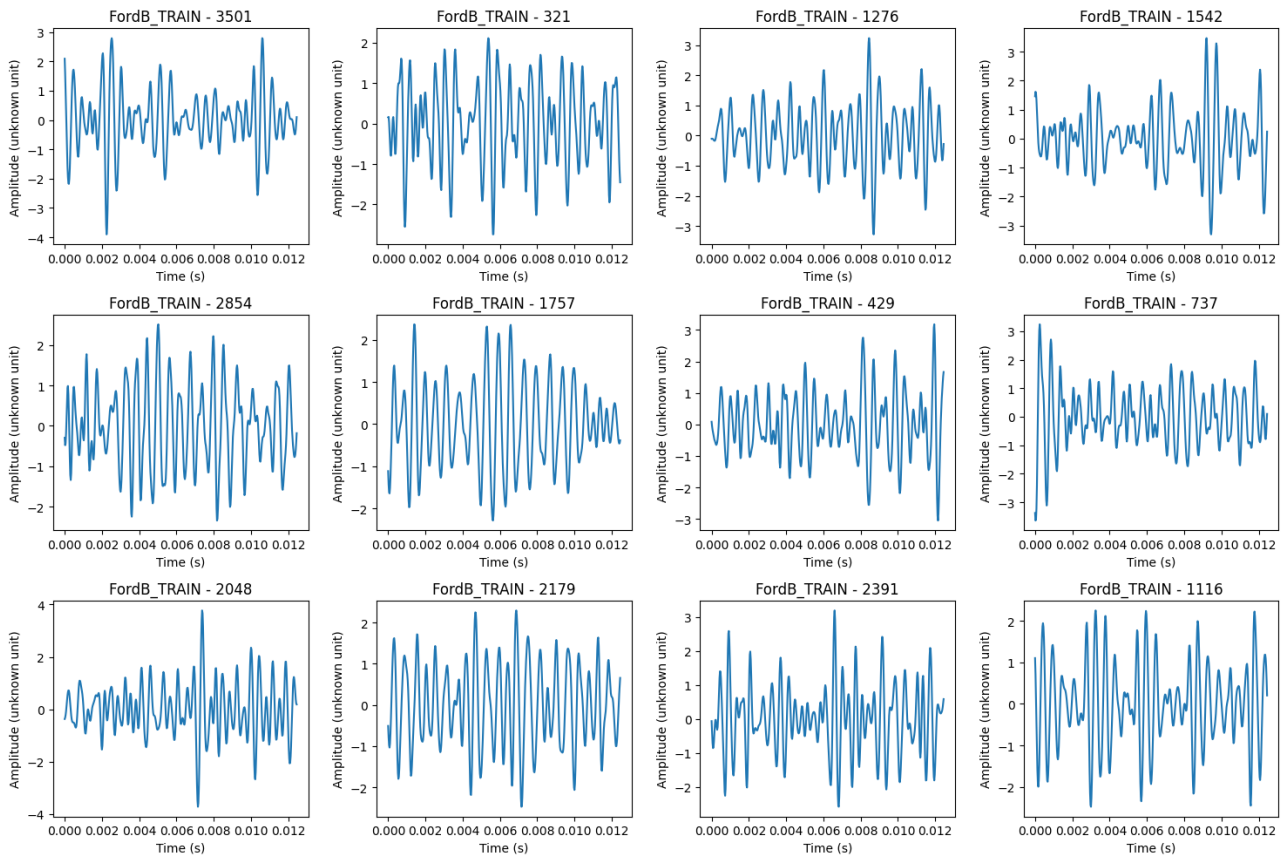


FIGURE 3 – Echantillons de signaux du jeu d'entraînement

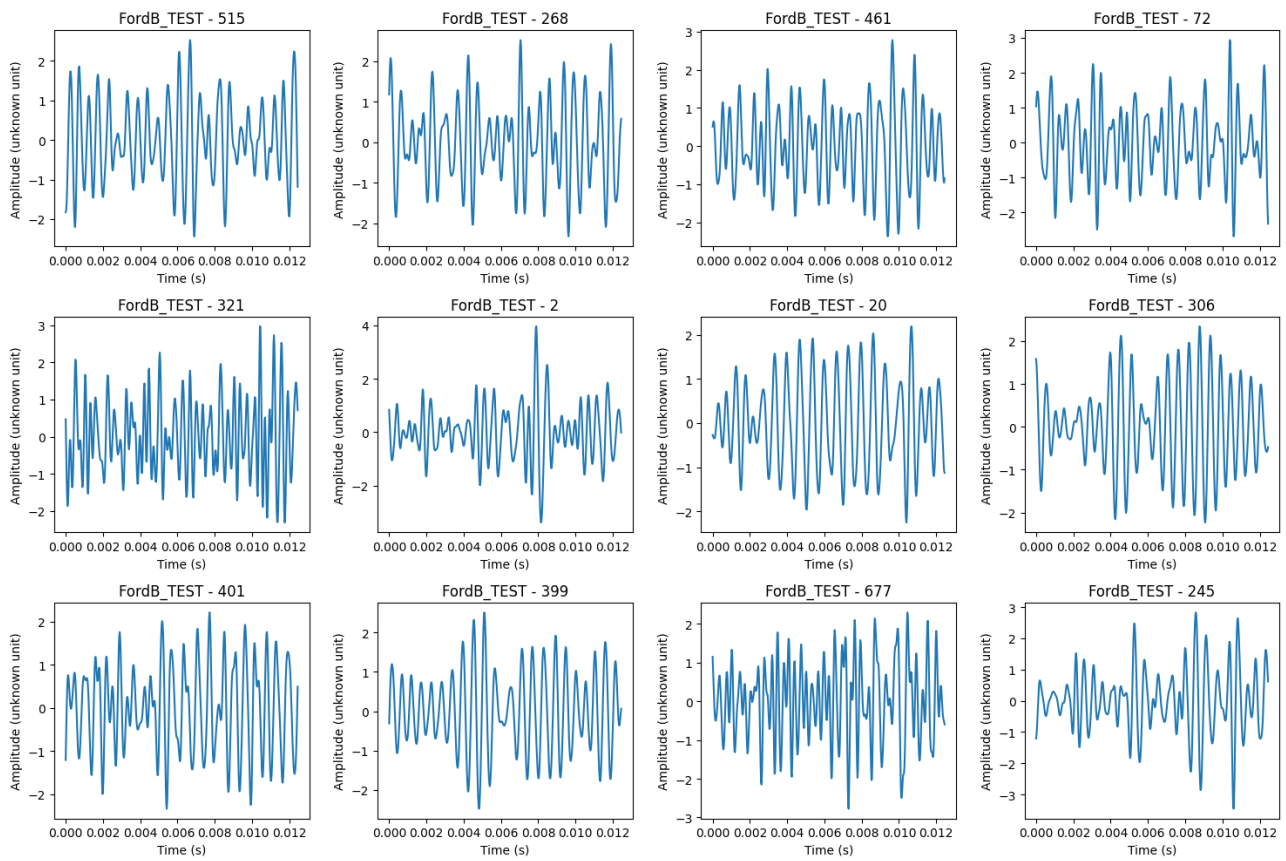


FIGURE 4 – Echantillons de signaux du jeu de test

## Visualisation des spectrogrammes

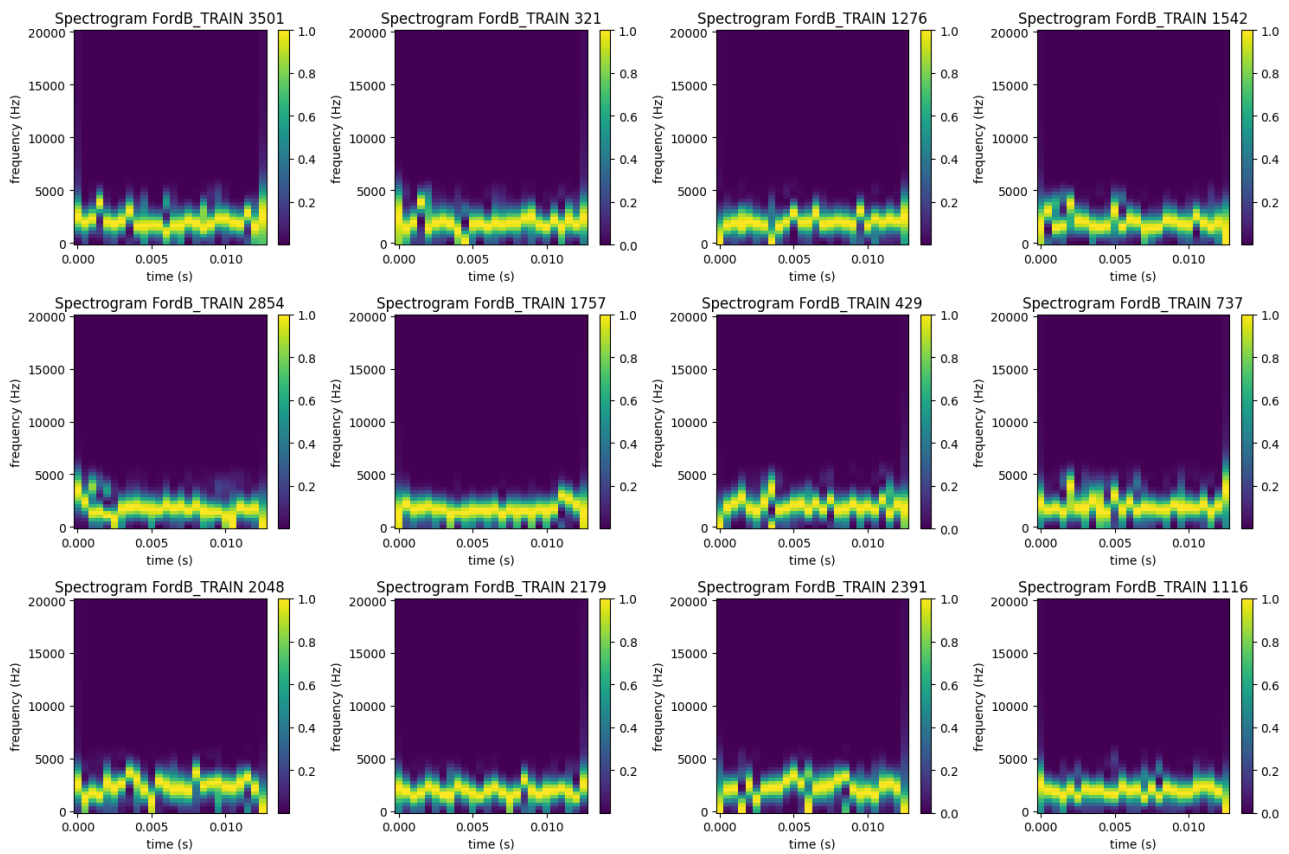


FIGURE 5 – Echantillons de spectrogrammes normalisés du jeu d'entraînement



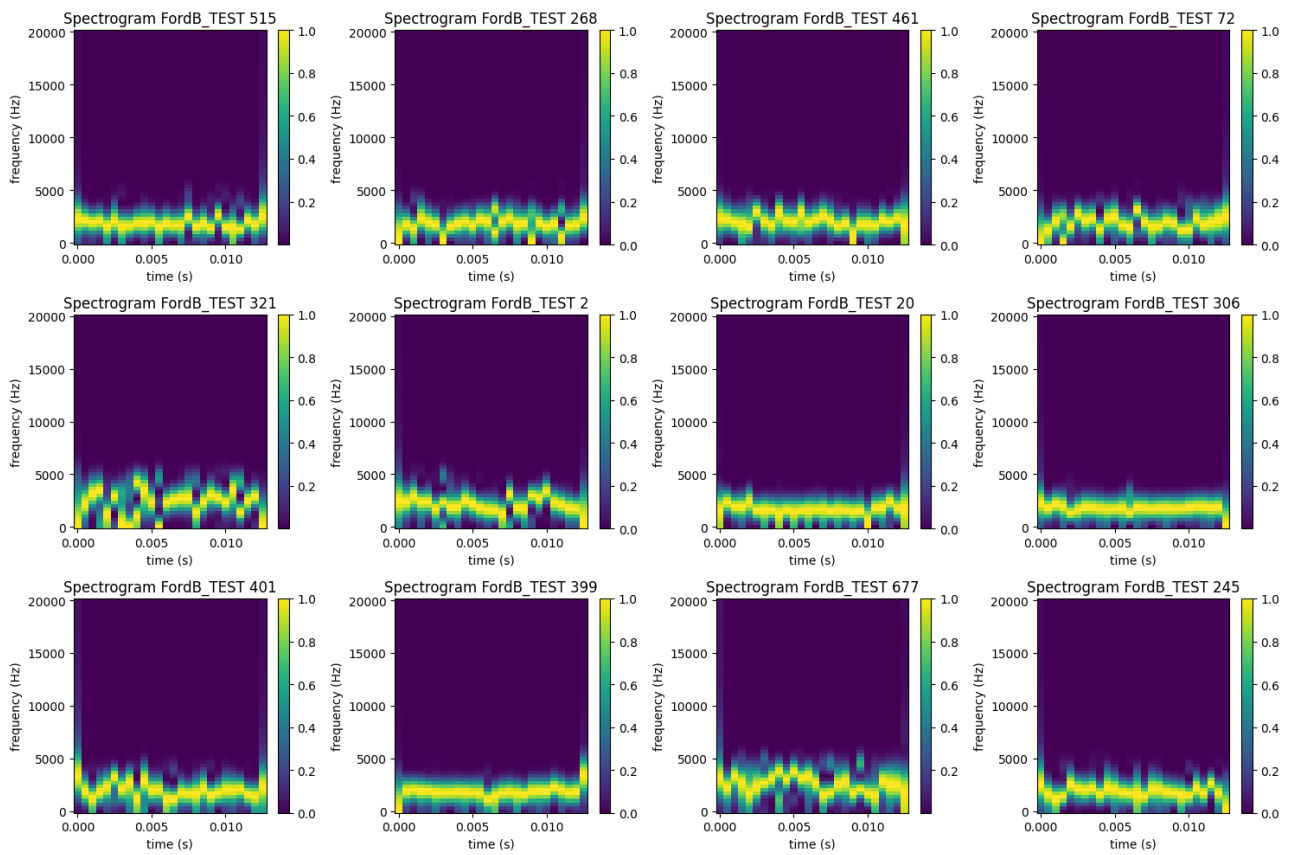


FIGURE 6 – Echantillons de spectrogrammes normalisés du jeu de test

## Segmentation des signaux temporels

Les signaux temporels observés en 3 et 4 ainsi que les spectrogrammes en 5 et 6 semblent indiquer des changements de régimes de bruits dans certains signaux. De plus, d'après les hypothèses faites sur la fréquence d'échantillonnage en partie, les signaux sont mesurés sur une période suffisamment courte pour supposer le régime moteur constant. De fait, des ruptures intempestives dans les enregistrements de bruit de moteur peuvent être indicatrices de la présence de sources de bruit momentanées s'ajoutant au bruit général du moteur. Ainsi on peut raisonnablement supposer que les signaux présentant de nombreuses ruptures ont de bonnes chances d'appartenir à la classe des moteurs défectueux.

On décide donc de faire un traitement de segmentation sur chaque signal afin d'estimer de manière automatique la quantité de changement de régime de bruits présent. Pour cela on utilise l'algorithme PELT [3] en distinguant deux segmentations :

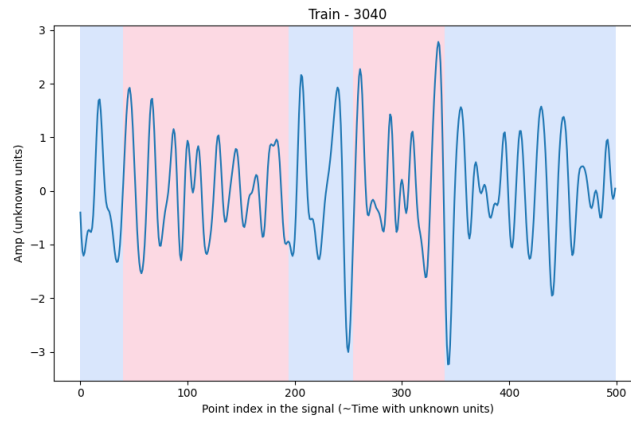
- Une faite sur les changements temporels de moyenne et d'écart-type. La fonction coût utilisée pour la segmentation dans ce cadre est :

$$c_{\Sigma}(t_k; t_{k+1}) = (t_{k+1} - t_k) \log \left( \sigma_{x[t_k; t_{k+1}]}^2 \right) + \frac{1}{\sigma_{x[t_k; t_{k+1}]}^2} \sum_{i \in [t_k; t_{k+1}]} \|x_i - \mu_{x[t_k; t_{k+1}]} \|_2^2$$

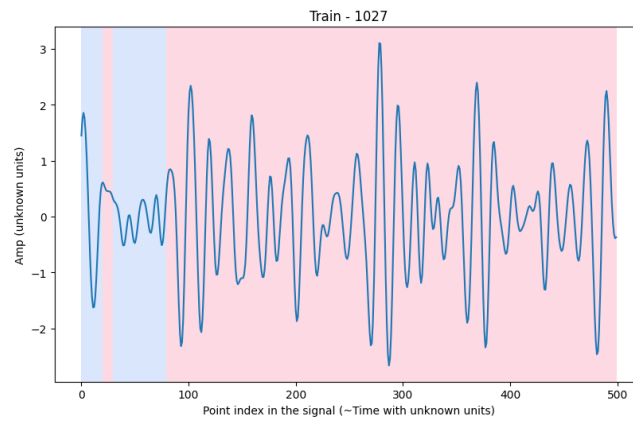
- Une faite sur la cohérence des portions de signal en terme d'information. Pour cela on se base sur la minimisation de l'entropie de Shannon de la densité spectrale de puissance normalisée sur chaque segment. L'intuition est que la superposition d'un régime "sain" et d'un régime "défectueux" doit être porteuse d'informations caractéristiques de ces deux régimes de fonctionnement à la fois. De fait une portion de signal présentant les deux régimes "sain" et "défectueux" aura plus d'information, donc une plus grande entropie, qu'un segment associé uniquement à un régime "sain" ou à un régime défectueux. L'entropie d'un segment sera prise comme l'entropie de Shannon de sa densité spectrale de puissance. La fonction coût utilisée dans ce cadre est :

$$c_H(t_k; t_{k+1}) = - \sum_i \frac{\Gamma_{x[t_k; t_{k+1}]}(f_i)}{\sum_j \Gamma_{x[t_k; t_{k+1}]}(f_j)} \log \left( \frac{\Gamma_{x[t_k; t_{k+1}]}(f_i)}{\sum_j \Gamma_{x[t_k; t_{k+1}]}(f_j)} \right)$$

Avec  $x[t_k; t_{k+1}]$  la portion de signal prise entre les points  $t_k$  et  $t_{k+1}$ .



(a) Exemple de segmentation sur les changement de moyenne et de variance



(b) Exemple de segmentation par minimisation de l'entropie de chaque segment

FIGURE 7 – Distribution des `rec_err` dans les différents jeux de données

## Conversion des observation issue de la segmentation en indice scalaire

Le scalaire qui sera interprété comme feature est le nombre de ruptures détectées. Intuitivement, on s'attend à ce que plus un grand nombre de ruptures est détecté, plus il y a de changement de régimes dans le signal et donc plus celui-ci est susceptible de traduire la présence d'un défaut dans le moteur.

## Largeur de la bande d'émission spectrale

### Visualisation des densités spectrales de puissance

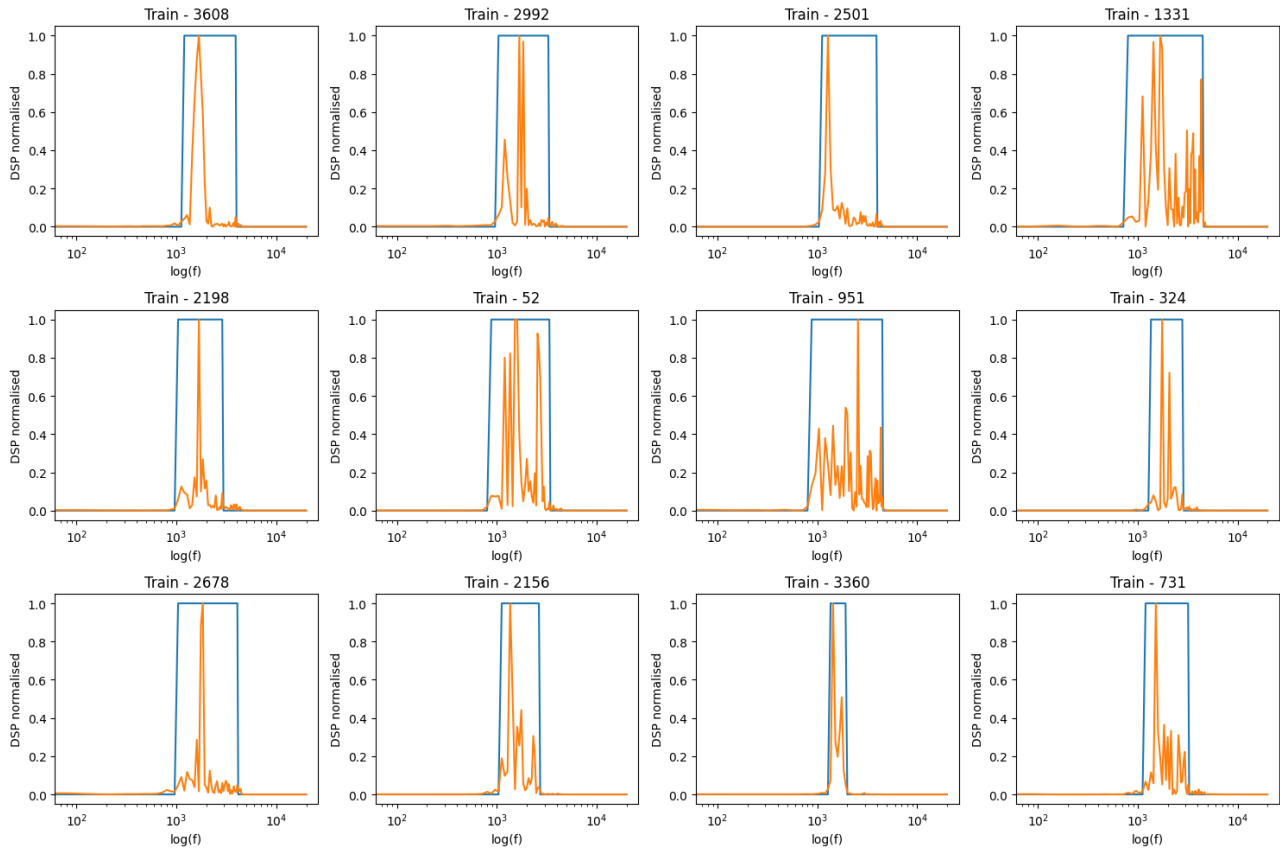


FIGURE 8 – Échantillons de DSP normalisée du jeu d'entraînement avec estimation de la bande d'émission spectrale

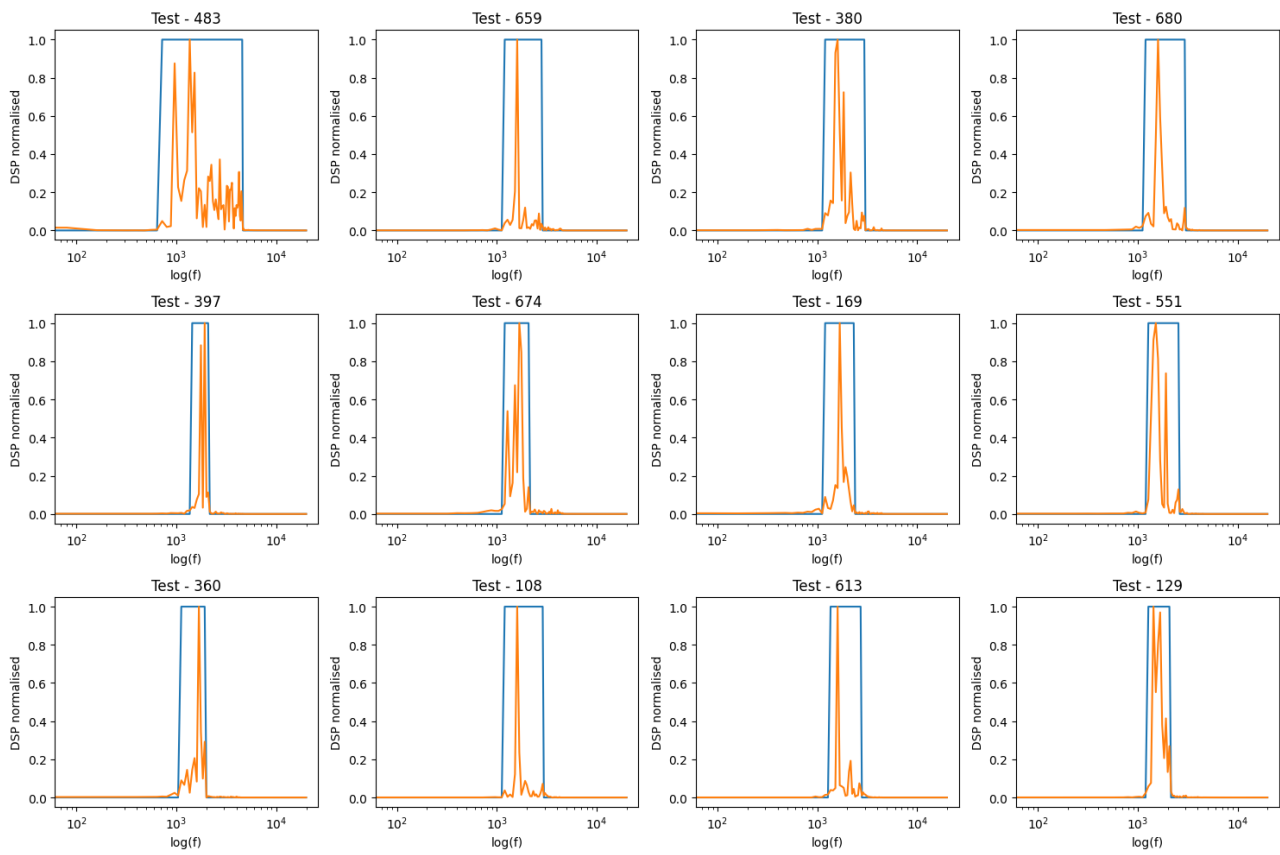


FIGURE 9 – Échantillons de DSP normalisée du jeu de test avec estimation de la bande d'émission spectrale

## Largeur de la bande d'émission spectrale

Les DSP en 8 en 9 mettent en évidence le fait que la bande d'émission spectrale est plus large pour certains signaux que pour d'autre. Plus précisément, certains signaux présentent plus d'émission à hautes fréquences que d'autre. La présence d'émission à hautes fréquences peut être interprété comme de la présence de bruits hors de la bande d'émission d'un moteur "normal". Ainsi, on peut raisonnablement supposer qu'une bande d'émission plus large que la normale indique la présence de défaut.

## Conversion en indice scalaire

La largeur de la bande d'émission spectrale est déjà un scalaire donné par :

$$\Delta f = f_{max} - f_{min}$$

Avec

$$\begin{cases} f_{min} = \sup\{f \geq 0 & / \quad \forall v \in [0, f] \quad ; \quad \Gamma_x(v) = 0\} \\ f_{max} = \inf\{f \geq 0 & / \quad \forall v \geq f \quad ; \quad \Gamma_x(v) = 0\} \end{cases}$$

En pratique les points où la DSP est réellement nulle sont rares du fait du formalisme numérique. De fait on se ramène à une largeur de bande d'émission estimée par un paramètre  $\epsilon$ . De plus, pour pouvoir utiliser la même valeur de  $\epsilon$  pour chaque signal, on considère désormais la DSP normalisée.

$$\hat{\Delta}_\epsilon f = \hat{f}_{max}^\epsilon - \hat{f}_{min}^\epsilon$$

Avec

$$\begin{aligned} \hat{f}_{min}^\epsilon &= \sup\{f \geq 0 & / \quad \forall v \in [0, f] \quad ; \quad \overline{\Gamma}_x(v) \leq \epsilon\} \\ \hat{f}_{max}^\epsilon &= \inf\{f \geq 0 & / \quad \forall v \geq f \quad ; \quad \overline{\Gamma}_x(v) \leq \epsilon\} \\ \overline{\Gamma}_x(v) &= \frac{\Gamma_x(v)}{\max_{f \geq 0}(\Gamma_x(f))} \end{aligned}$$

## Profil d'autocorrélation

### Visualisation des profils d'autocorrélation

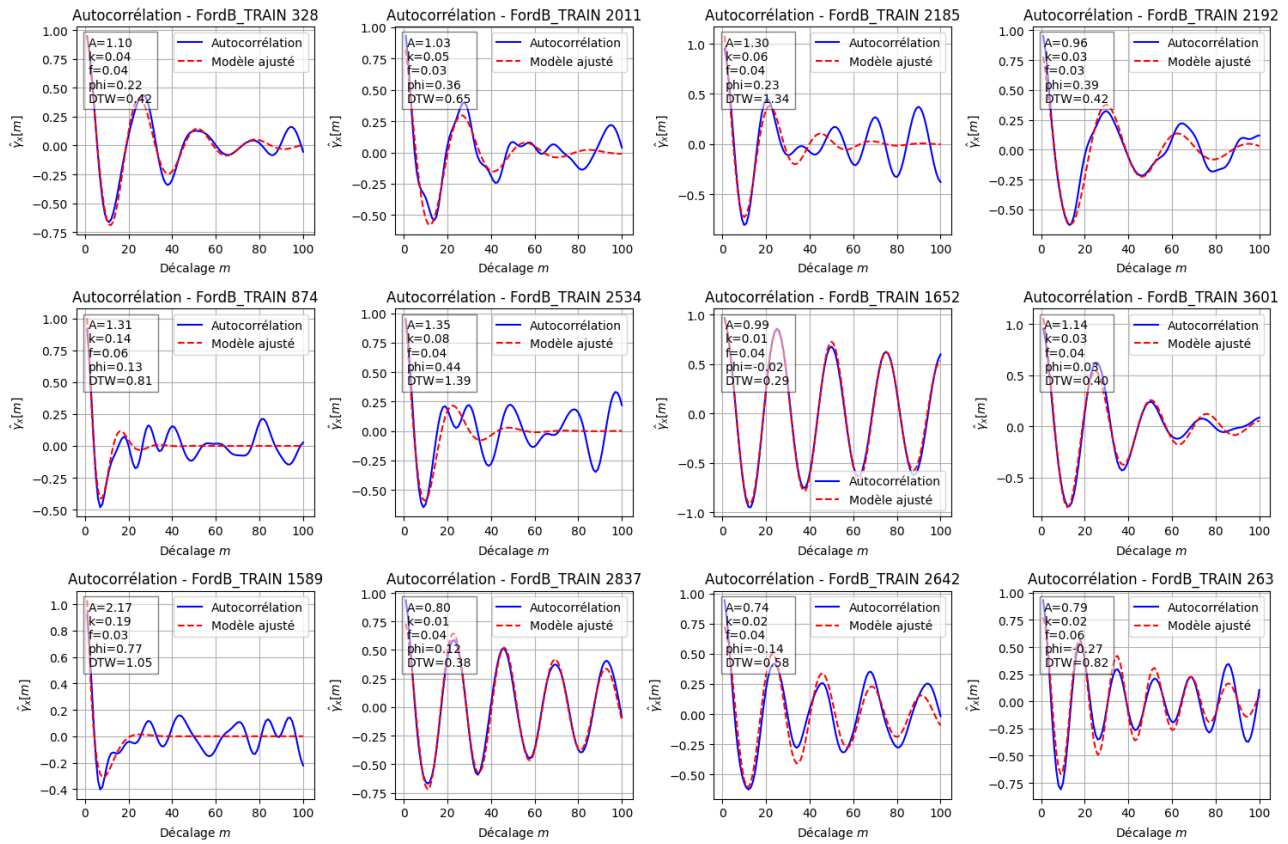


FIGURE 10 – Échantillons de profils d'autocorrélation avec modèle ajusté issus du jeu d'entraînement

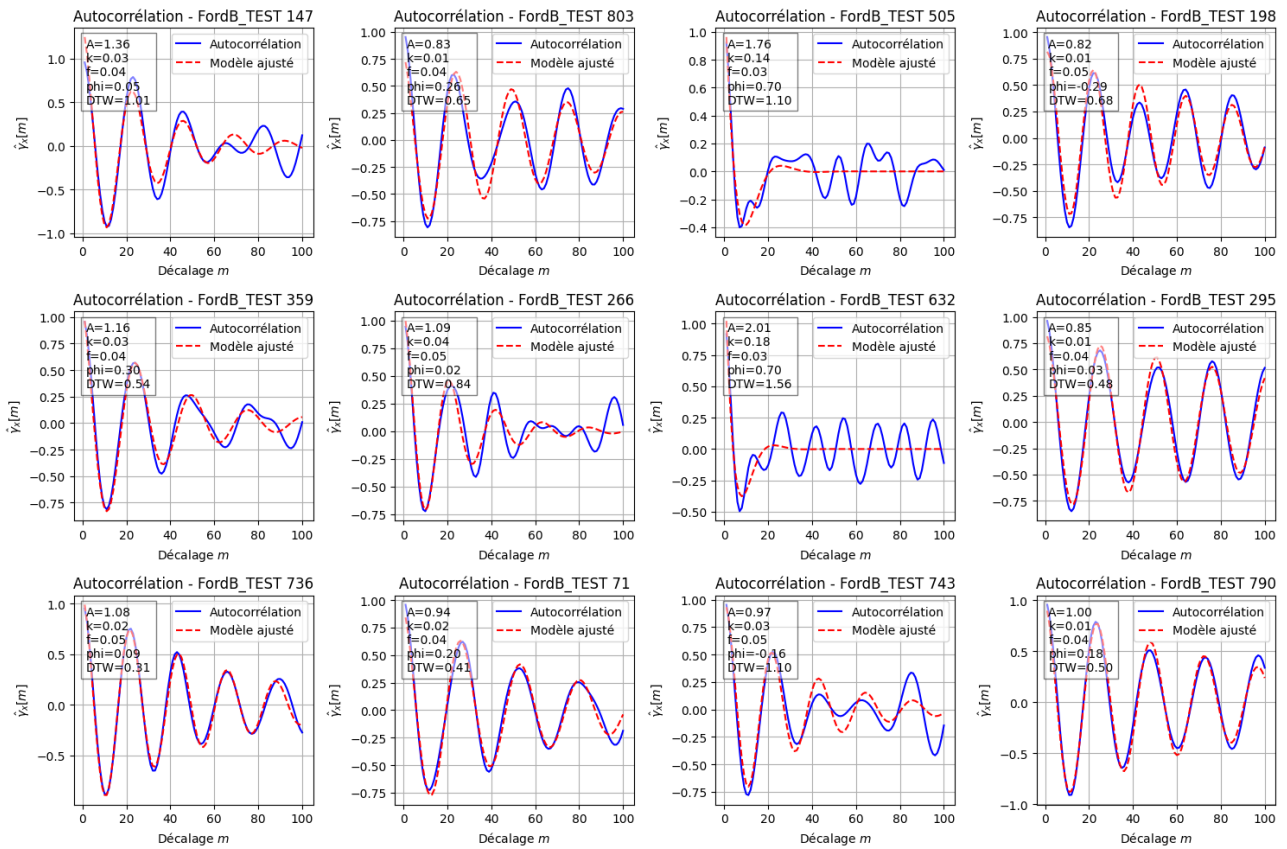


FIGURE 11 – Échantillons de profils d'autocorrélation avec modèle ajusté issus du jeu de test



### Profil d'autocorrélation

Les profils d'autocorrélations des signaux semblent se partager en deux classes. L'une avec des profils sinusoïdaux plutôt régulier et l'autre avec des profils sinusoïdaux très amorti ou très irréguliers. D'après notre hypothèse sur la fréquence d'échantillonnage faite en les signaux sont issus de mesures sur une période très courte. De fait le régime moteur est censé être cyclique et régulier sur toute la durée du signal. Par conséquent, un profil d'autocorrélation irrégulier ou amorti est potentiellement indicateur d'un fonctionnement défectueux du moteur auquel il est associé.

### Conversion en indice scalaire

Les indices scalaires servant à caractériser le profil d'autocorrélation sont les paramètres  $A, k, \phi$  et  $f$  d'un modèle de sinusoïde amortie :

$$\hat{\gamma}_x(m) = Ae^{-km} \cos(2\pi fm + \phi)$$

De plus, pour tenir compte des signaux au profil irréguliers ne pouvant pas être suffisamment bien modélisé par un sinus amorti, nous prendrons également l'erreur DTW entre le modèle et le profil réel comme indice scalaire.

## Solution choisie

### Sélection des features intéressantes pour la classification

#### Nécessité de la sélection

Les consignes du mini-projet imposent de classer les signaux suivant trois features uniquement. Si les trois premières composantes principales sont les plus efficaces pour maximiser la variance entre les données [2] elles présentent un inconvénient de taille. En effet, elles sont issues de combinaisons linéaires de chaque features, les coefficients de cette combinaison linéaire sont obtenu uniquement dans le but de maximiser la variance entre les données. De fait les composantes principales sont difficilement interprétable or, dans le cadre de la détection de défauts industriel, il est crucial de pouvoir interpréter les critères de détection.

De fait, on cherche à extraire les trois features qui semblent être les plus pertinentes parmi celles proposées en .

#### Observation des profils de distribution

L'observation des distributions individuelles de chaque feature donne déjà un bon aperçu de la capacité de chaque feature à séparer efficacement les données. De plus, cette observation permet de faire un premier tri, ce qui permettra par la suite d'avoir une sélection plus fine vis-à-vis du critère par PCA. Plus particulièrement :

- Les features dont la distribution présente un unique pic de forme gaussienne sont éliminées.
- Les features dont la distribution est trop faiblement étalée sont éliminées.
- Les features dont la distribution est trop uniforme sont éliminées.

Finalement, l'observation des figures en et sous le critère de sélection énoncé précédemment permet de ne garder que :

- L'entropie de Shannon du signal complet "H".
- Le facteur d'amortissement du modèle ajusté sur le profil d'autocorrélation du signal complet "k".
- L'erreur de reconstruction du modèle ajusté sur le profil d'autocorrélation du signal complet "rec\_err".
- La largeur de bande d'émission spectrale "width".

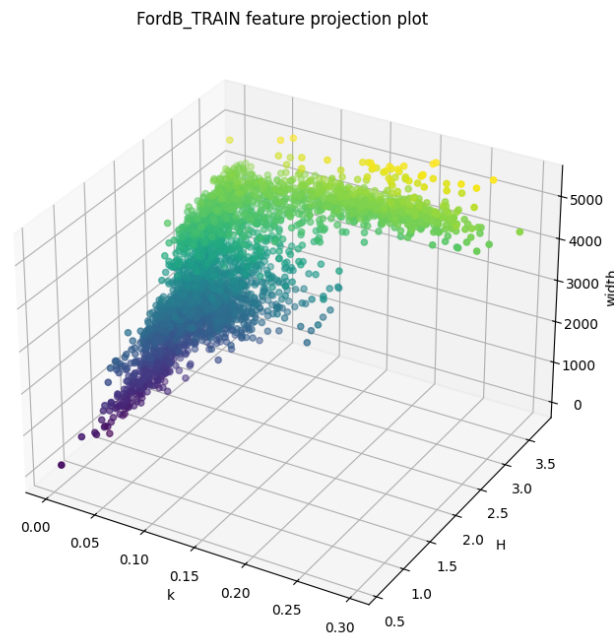
#### Analyse par Composante Principale

L'analyse en composante principales des features retenues sur le jeu de données d'entraînement et de test permet de conclure que la feature "rec\_err" n'est pas fiable. En effet, si elle était fiable alors sa pertinence ne devrait pas changer avec le type de jeu de données, ce qui serait mit en évidence par un ressemblance entre les deux affichage des features dans le plan en composante principales. On en déduit donc par élimination que les features les plus pertinentes sont :

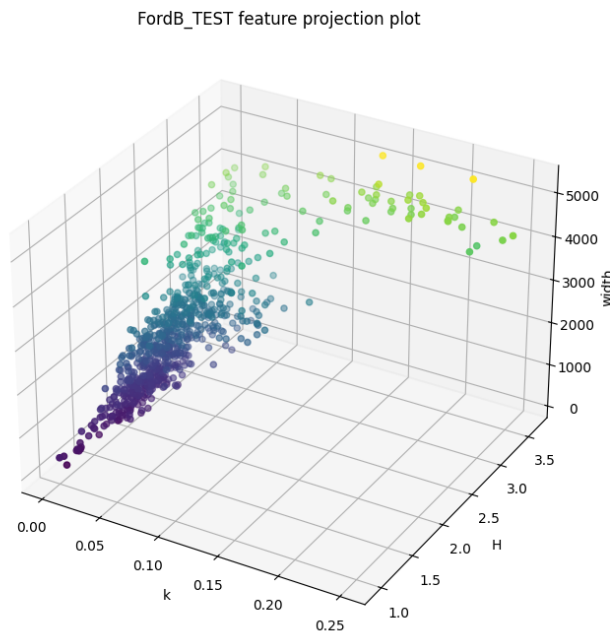
- L'entropie de Shannon du signal complet "H".
- Le facteur d'amortissement du modèle ajusté sur le profil d'autocorrélation du signal complet "k".
- La largeur de bande d'émission spectrale "width".

## Résultats

### Visualisation des données dans l'espace des features sélectionnées



(a)



(b)

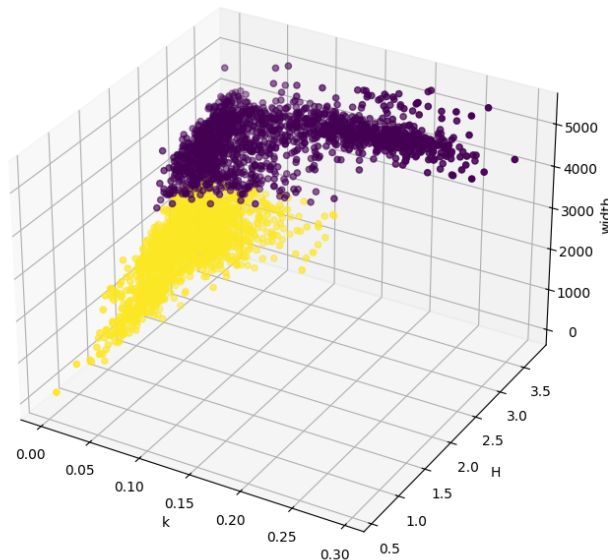
FIGURE 12 – Données projetées dans le plan des deux premières composantes principales

La visualisation des jeux de données projetées dans l'espace engendré par les trois features sélectionnées met en évidence la présence de deux cluster. L'un "en haut" sur la figure, associé à un de grande valeurs de width, k et H. L'autre "en bas", associé à un de faible valeurs de width, k et H.

## Clustering

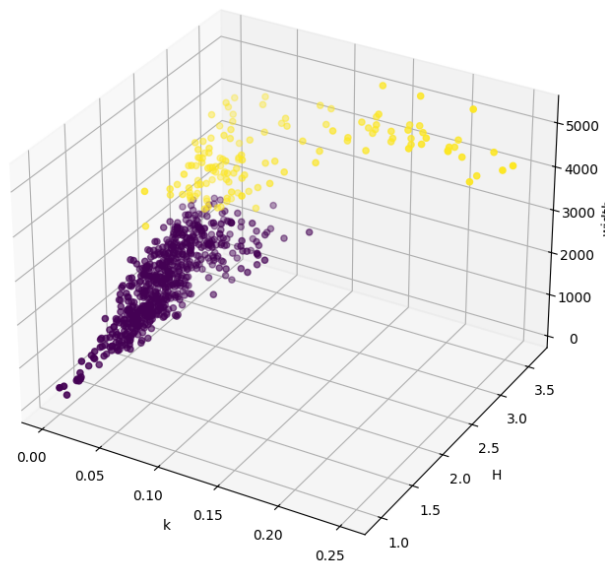
La visualisation des données dans l'espace engendré par les trois features sélectionnées a permis de mettre en évidence la présence de deux clusters. On choisit de les mettre encore plus en évidence via un clustering par l'algorithme k-mean.

FordB\_TRAIN feature projection plot with k\_mean cluster coloring



(a)

FordB\_TEST feature projection plot with k\_mean cluster coloring



(b)

FIGURE 13 – Données projetées dans l'espace engendré par les feautres chosies avec clustering

Comme les deux cluster sont suffisamment distinct pour qu'on puisse les distinguer à l'œil nu on en déduit qu'il est raisonnable d'utiliser l'algorithme k-means pour effectuer le clustering. On obtient alors les résultats ci-dessous :

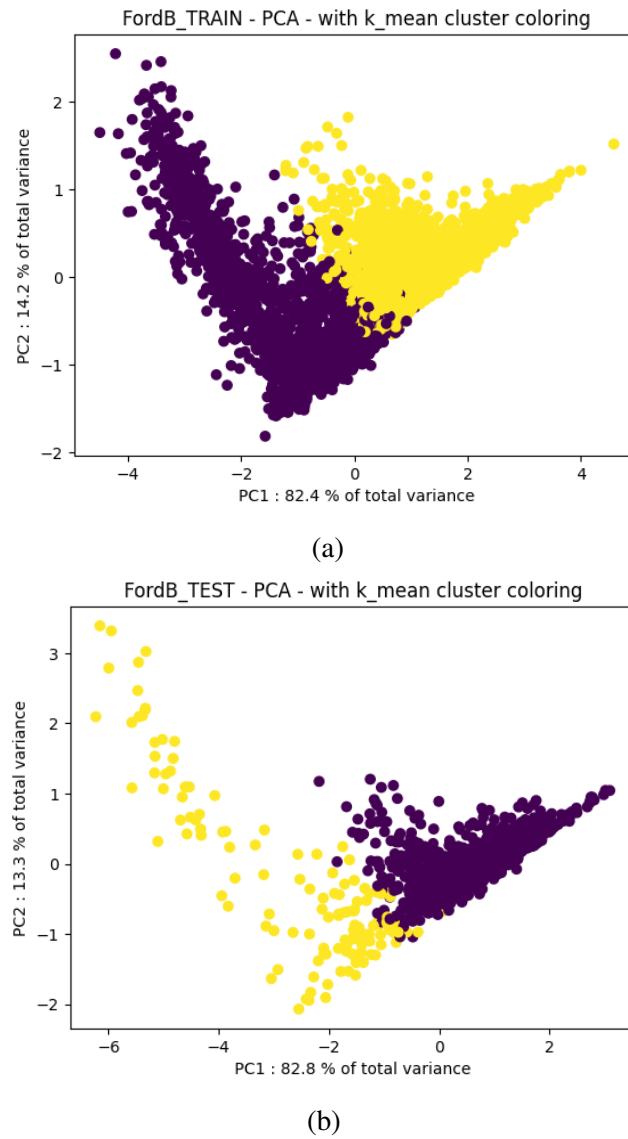


FIGURE 14 – Données projetées dans le plan des deux premières composantes principales avec clustering

Si la visualisation des données projetée dans l'espace engendré par les features choisies semble indiquer que la séparation par k-mean est pertinente, on peut cependant noter que la frontière entre les deux cluster est floue. Cela est notamment mit en évidence par la visualisation de la projection des données sur les deux premières composantes principales.

## Conclusion

Le mini-projet qui a été effectué avait pour objectif de séparer en deux classes un groupe de signaux sans utiliser de méthodes supervisées (sans avoir connaissance de la classe des signaux de la base d'entraînement). Les signaux que nous avons choisis de traiter étaient des bruits de moteurs et les classes de ces signaux étaient la présence d'un défaut ou non.

Une première analyse qualitative sur la nature audible des signaux a permis d'estimer la fréquence d'échantillonnage utilisée pour les enregistrer.

Ensuite, un raisonnement qualitatif sur la visualisation de caractéristiques générale des signaux telle que leur propriété temporelle, spectrale et leur autocorrélation ont permis de trouver des pistes sur les caractéristiques prometteuses dans l'objectif de trouver des features mettant en évidence la présence de défauts.

Puis, des analyses plus détaillées sur chaque caractéristiques prometteuses ont permis d'extraire des features potentiellement pertinentes pour séparer les signaux en deux classes. Les trois features les plus pertinentes parmi celles proposées ont été choisies. D'abord, un premier tri grossier a été effectué en éliminant toutes les features dont les distributions ne semblaient pas porter d'informations pertinentes quant-à la séparation des données en deux classes. Ensuite une analyse plus fine basée sur de l'analyse en composantes principales a permis de faire la sélection des features finales.

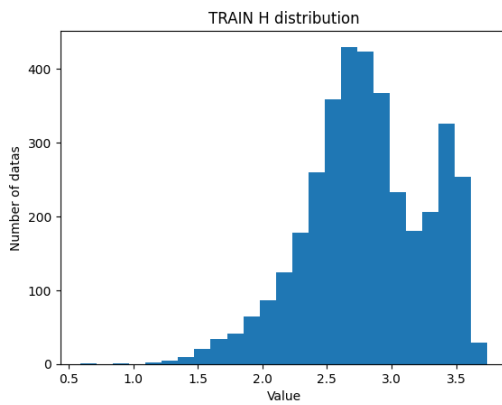
Enfin un clustering a été effectué et les données, avec leur appartenance à leur cluster respectifs, ont été visualisées dans l'espace engendré par les 3 features choisies pour évaluer la capacité des features sélectionnées à séparer les données de manière cohérente.

En somme, nous pensons avoir trouvé des features permettant de faire une séparation pertinentes des données. La pertinence des features choisies peut être interprétée intuitivement en observant la figure ??, notamment par le fait qu'un cluster soit associé à de fortes valeurs de H, k et width. En effet, la présence d'un défaut devrait ajouter du bruit supplémentaire au bruit naturel du moteur, ce qui se traduit par un H plus élevé pour les moteurs défectueux. De plus le bruit ajouté par la présence du défaut ne devrait pas, à priori être totalement compris dans les émissions spectrales naturelles du moteur, ce qui se traduit par des valeurs de width plus grandes pour les moteurs défectueux. Enfin, la présence de défauts induit des régimes de bruit irréguliers, ce qui se traduit par un profil d'autocorrélation diminuant avec le décalage, ce qui est traduit par de forte valeur de k pour les moteurs défectueux. Cependant, la redondance des informations contenues dans les features sélectionnées, notamment entre H et width, semble brouiller la frontière entre les deux classes. De fait, nous pensons que ces features ne sont pas parfaitement adaptées pour effectuer la séparation, bien qu'elles soient tout de même pertinentes.

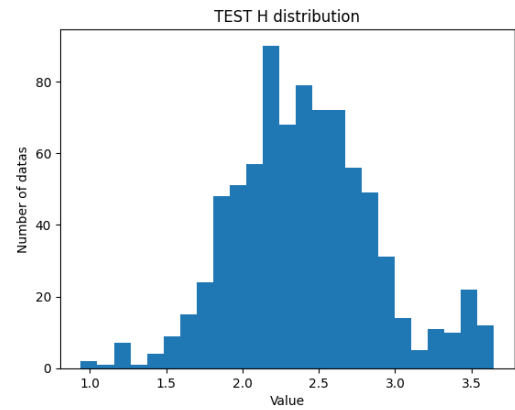
# Bibliographie

- [1] Laurent OUDRE. « Intelligence Artificielle & Machine Learning pour la modélisation de séries temporelles et de signaux - Séance 5 : Pré-traitements des séries temporelles ». fr. In : ().
- [2] Laurent OUDRE. « Intelligence Artificielle & Machine Learning pour la modélisation de séries temporelles et de signaux - Séance 6 : Extraction et sélection de caractéristiques ». fr. In : ().
- [3] Laurent OUDRE. « Intelligence Artificielle & Machine Learning pour la modélisation de séries temporelles et de signaux - Séance 8 : Détection de ruptures et d'anomalies ». fr. In : ().

## Annexe : Distributions des features retenues

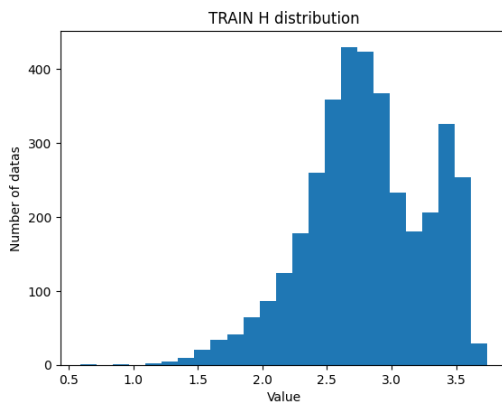


(a)

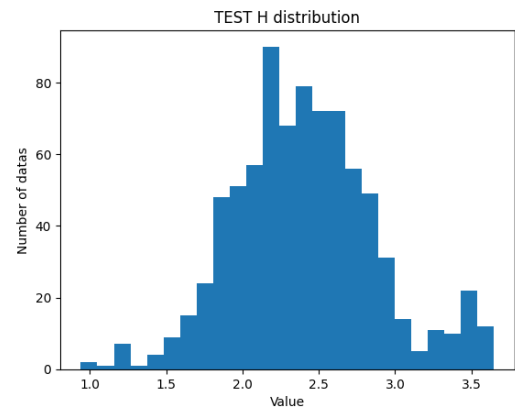


(b)

FIGURE 15 – Distribution des k dans les différents jeux de données

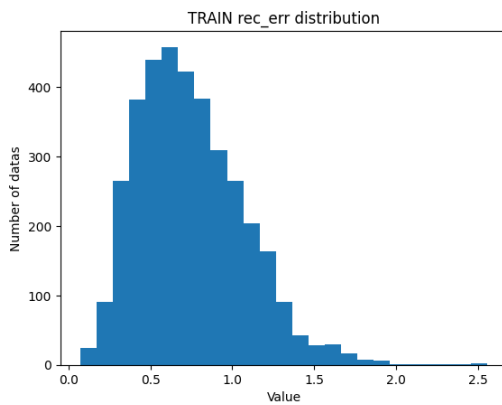


(a)

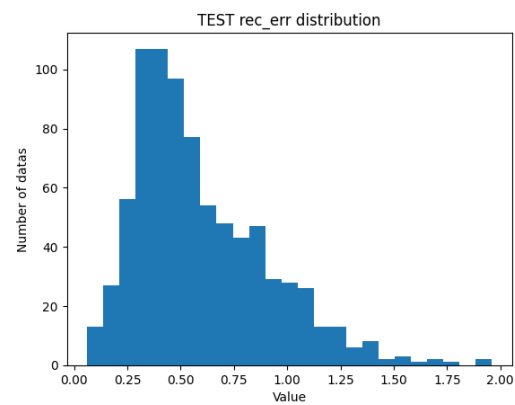


(b)

FIGURE 16 – Distribution des k dans les différents jeux de données



(a)



(b)

FIGURE 17 – Distribution des rec\_err dans les différents jeux de données



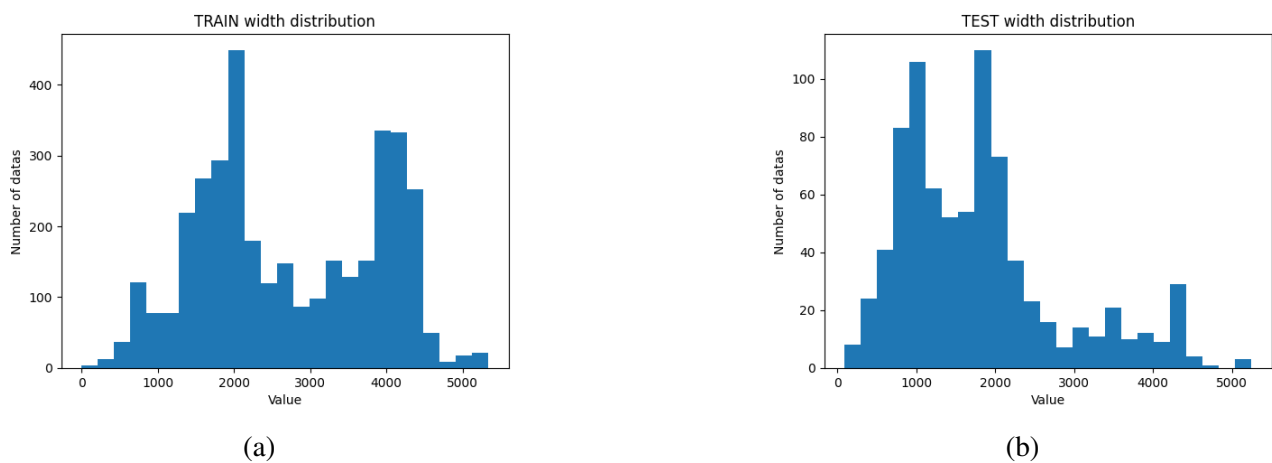
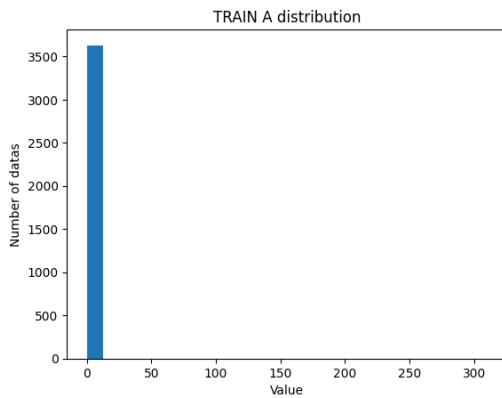
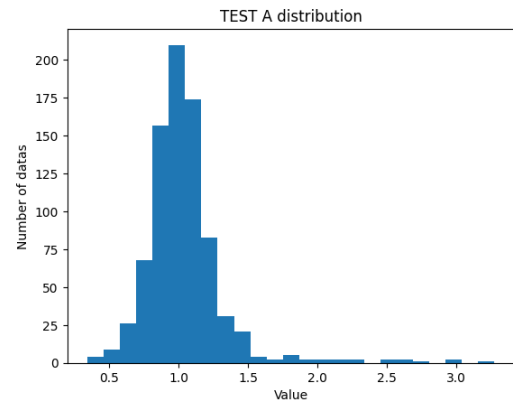


FIGURE 18 – Distribution des width dans les différents jeux de données

## Annexe : Distributions des features éliminées

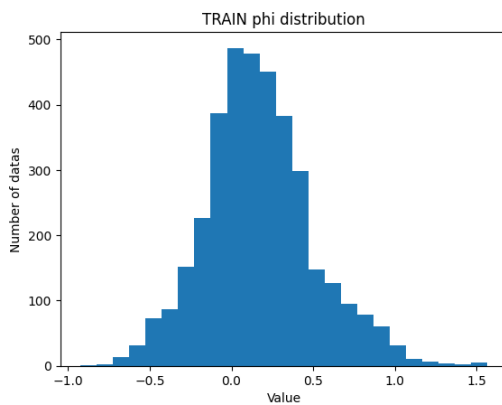


(a)

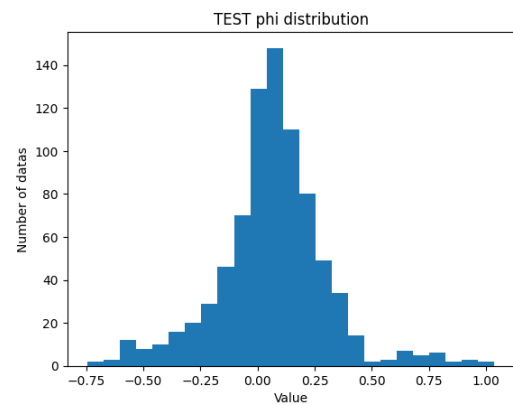


(b)

FIGURE 19 – Distribution des A dans les différents jeux de données

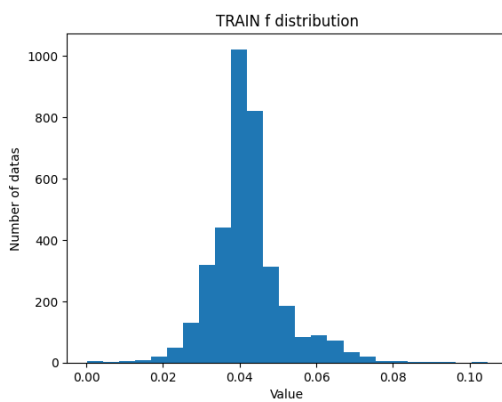


(a)

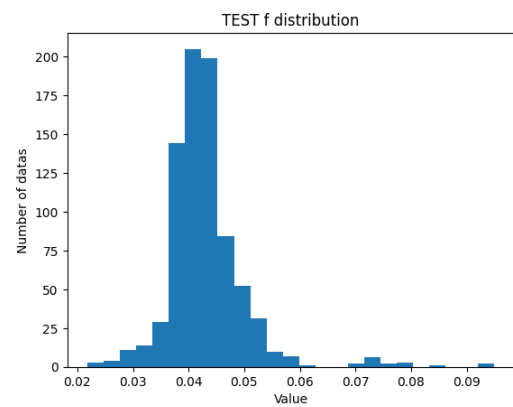


(b)

FIGURE 20 – Distribution des phi dans les différents jeux de données

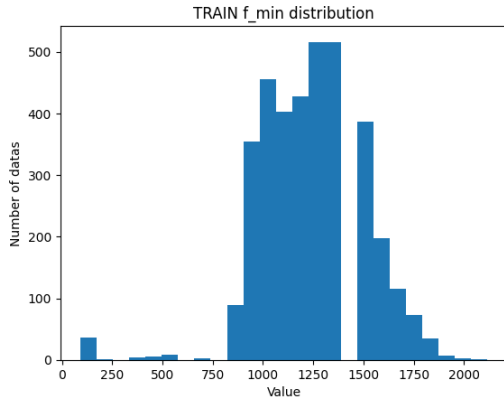


(a)

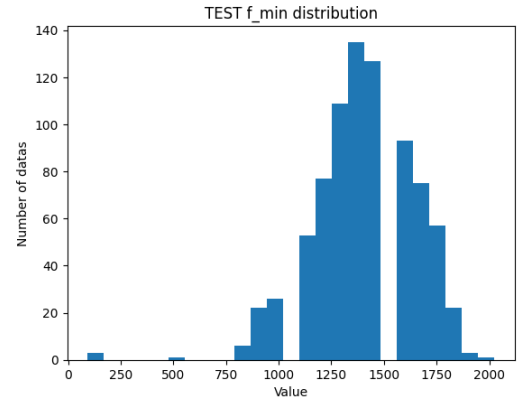


(b)

FIGURE 21 – Distribution des f dans les différents jeux de données

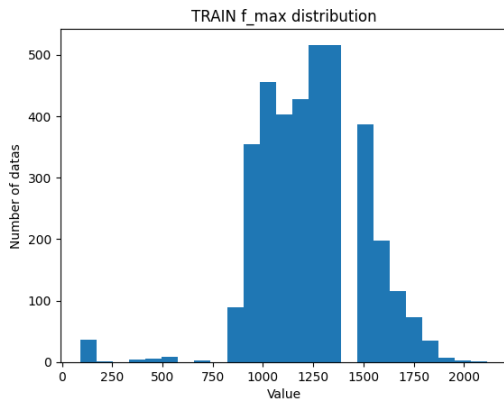


(a)

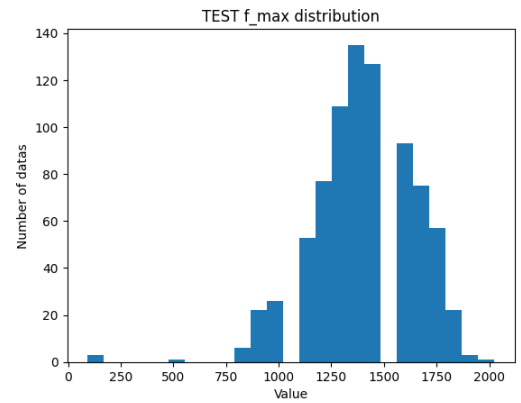


(b)

FIGURE 22 – Distribution des f\_min dans les différents jeux de données

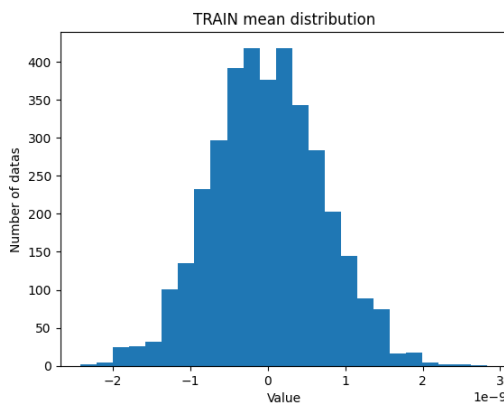


(a)

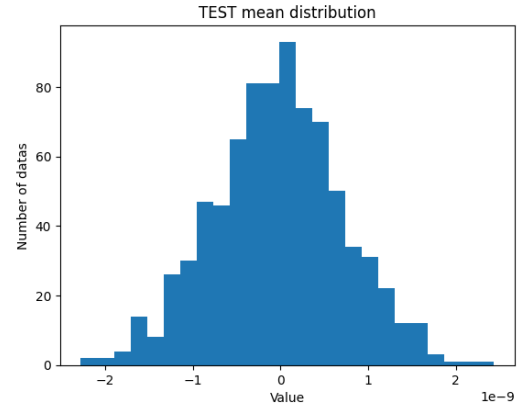


(b)

FIGURE 23 – Distribution des f\_max dans les différents jeux de données

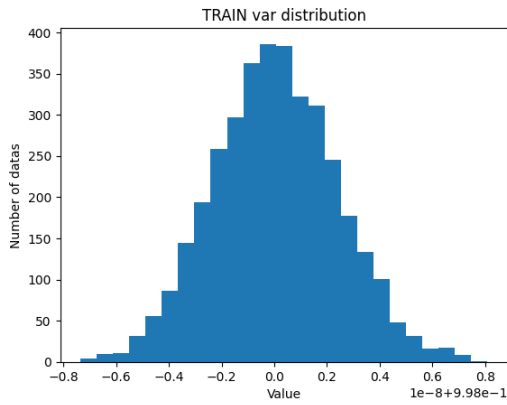


(a)

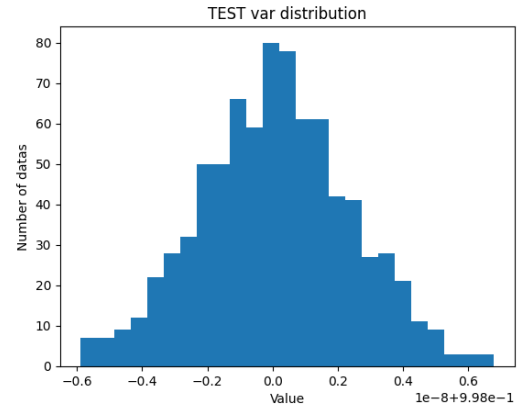


(b)

FIGURE 24 – Distribution des mean dans les différents jeux de données

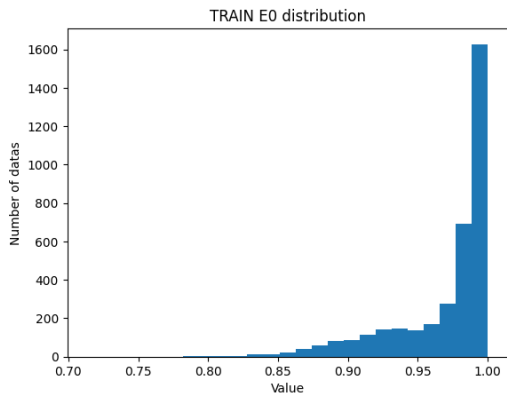


(a)

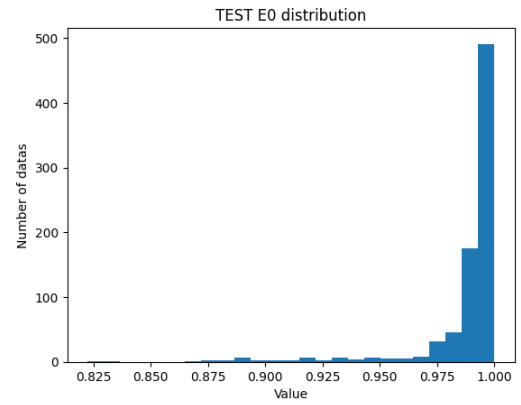


(b)

FIGURE 25 – Distribution des var dans les différents jeux de données

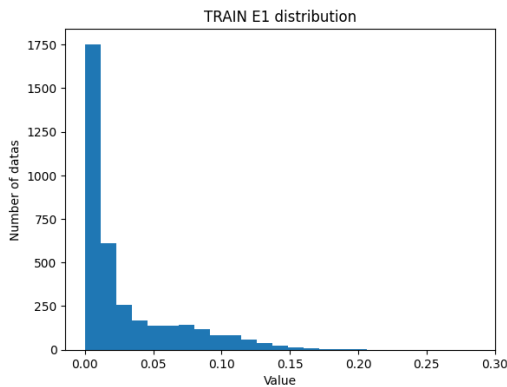


(a)

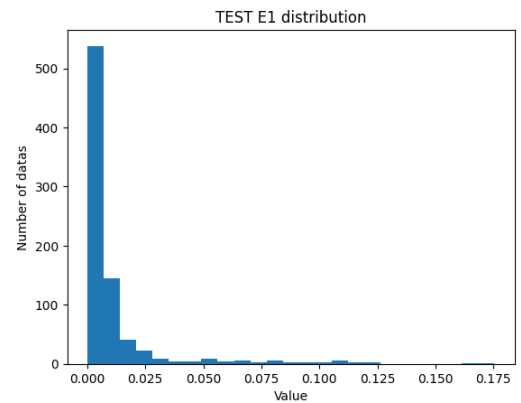


(b)

FIGURE 26 – Distribution des E0 dans les différents jeux de données

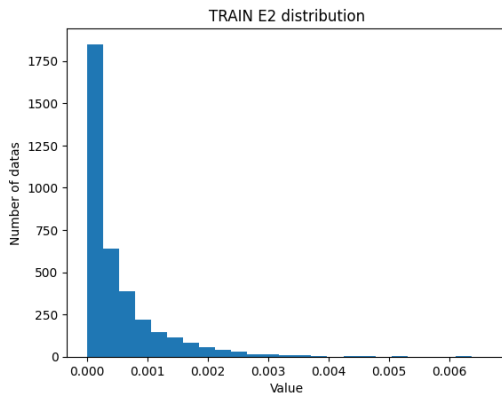


(a)

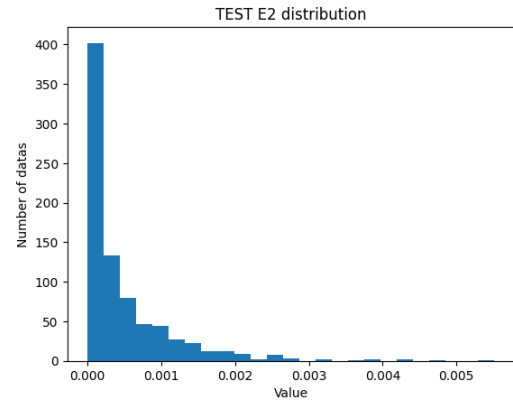


(b)

FIGURE 27 – Distribution des E1 dans les différents jeux de données

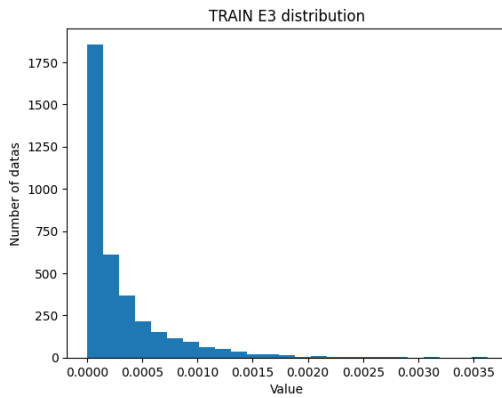


(a)

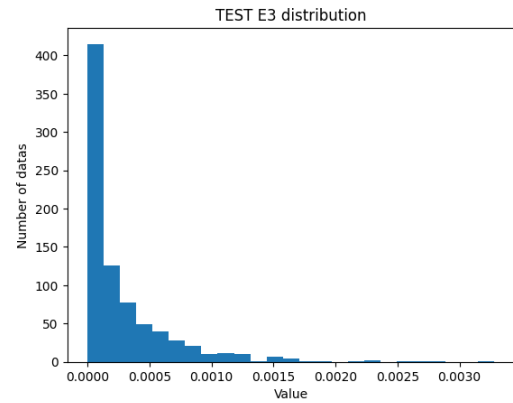


(b)

FIGURE 28 – Distribution des E2 dans les différents jeux de données

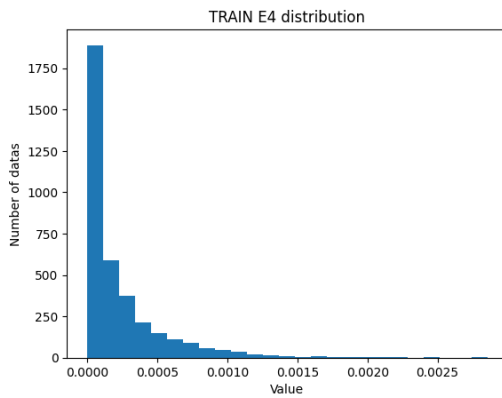


(a)

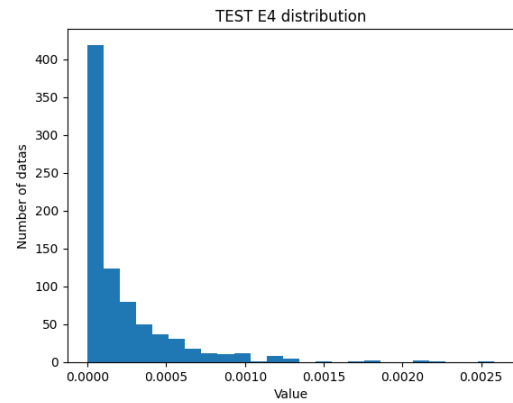


(b)

FIGURE 29 – Distribution des E3 dans les différents jeux de données



(a)



(b)

FIGURE 30 – Distribution des E4 dans les différents jeux de données

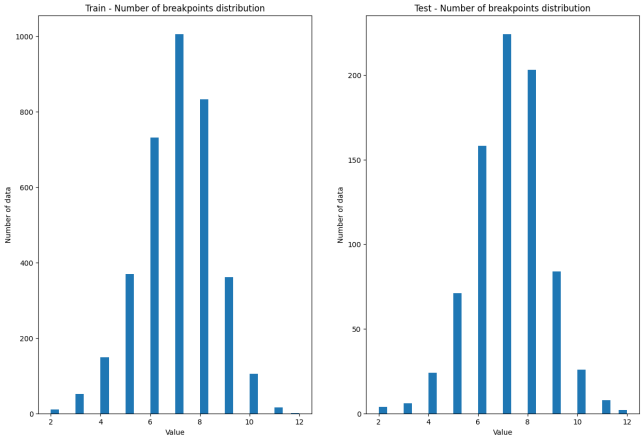


FIGURE 31 – Distribution des breaks segmentées par changement de variance et de moyenne dans les différents jeux de données