

Can You Catch It ?

IDASM104 : Projet interdisciplinaire



Jean ALBRECQ, Antoine P.,
Cyprien L., Jessica D.

Professeurs : S. Faulkner B. Frénay V. Salnikov.
2020-2021

Table des matières

1	Introduction	2
1.1	Présentation du projet	2
2	Préparation et visualisation des données	3
2.1	Récolte des données	3
2.2	Premier coup d'œil à la structure des données	3
2.3	Préparation des données	4
2.3.1	Stratification des données	4
2.4	Visualisation des données	4
2.4.1	Création du test-set	4
2.4.2	Premières visualisations	4
2.4.3	Recherche de corrélation	5
2.5	Combinaison de <i>features</i>	6
	Appendices	7
A	Position des stops	7
B	Première visualisation des données	8
C	Recherche de corrélation	11

Chapitre 1

Introduction

Le but de ce rapport est d'expliquer la démarche et la méthodologie qui a guidé l'élaboration des modèles de machine learning et de l'analyse des données fournie par les *opendata stib-mivb*. Ce rapport est constitué des différentes parties : l'analyse et la préparation des données, l'entraînement des modèles de régression et de classification et l'analyse de leurs résultats. Nous nous sommes concentré uniquement sur une ligne de bus mais le système pourrait facilement être étendu au reste du réseau.

L'étape d'analyse et la préparation des données met en lumière les notions de normalisation, la détecteur d'*outliers*, la selection de *features*. La visualisation des données est également une partie importante de l'analyse des données. En suite dans l'étape d'entraînement des modèles passe par un phase de selection des méta-paramètres et d'optimisation des prédictions.

1.1 Présentation du projet

Il nous a été demandé de développer un nouveau service ou une analyse pertinente par rapport au défis de la mobilité. Plusieurs *opendata* nous était proposées, nous avons décidé de choisir celle de la STIB. Nous avons choisi de mettre en place un service permettant de savoir si prochain bus qui arrivera à un stop que l'on attend aura du retard ou non.

Chapitre 2

Préparation et visualisation des données

2.1 Récolte des données

La première étape est de toute évidence la récolte des données. Notre projet nous demandais d'avoir accès à un historique de retard mais malheureusement cette historique ne fait pas partie des datasets des opendata STIB. Nous avons donc développé un script python¹ nous permettant de constituer cet historique de retard. Pour obtenir le délais, le script compare le temps d'arrivée théorique (qui nous est fournis par les fichiers GTFS²) et l'heure d'arrivée prévue (qui nous est fournie par l'api "*waiting time*"). Le délais est enregistré dans un fichier csv. En plus du délais, le script enregistre la température, la vitesse du vent, l'humidité et la visibilité grâce à l'api OpenWeather³. Un nouveau fichier csv est généré chaque jour.

Nous avons dans un premier temps récolté les données pour deux stop (les numéros 0089 et 6608G, voir leur emplacement dans l'annexe A.1) du premier novembre au douze novembre. Les fichiers csv sont disponibles sur le *repository*⁴ du projet. Dans un second temps, nous avons récolté les données de tout les stops d'une ligne de bus (la ligne 39). La position de tout les stops de la ligne sont visible sur l'annexe A.2⁵. Les données récoltées durant cette deuxième phase sont disponible sur le *repository* du projet⁶.

2.2 Premier coup d'œil à la structure des données

La commande `data.info()` qu'il y a 12798 lignes sans valeur pour la colonne `delay`. Pour chacune des *features* un graphique du nombre d'occurrences par valeur a été créé (voir l'annexe B.1). On peut voir sur ces graphiques que plusieurs *features* ont toujours la même

1. Le code de ce script est disponible à l'adresse suivante : [lien github du script](#)

2. General Transit Feed Specification

3. Documentation disponible [ici](#)

4. Disponible [ici](#)

5. Une carte GoogleMyMaps est également disponible [ici](#)

6. Disponible [ici](#)

valeur. On peut également voir que des retards ont été enregistrés pour d'autres lignes que la numéro 39, il faudra donc supprimer ces dernières.

Dans le graphique en annexe B.1 on remarque qu'il y a moins d'occurrences de valeur quatre pour la *features* `hour`. Cela est dû au fonctionnement de l'API de la STIB. On peut également voir qu'il y a plus de d'occurrences pour les valeurs entre dix et quinze de la *features* `hour`, cette différence est due à l'heure de démarrage et d'arrêt du script de récolte des données.

Le boxplot en annexe B.2 montre la répartition des retards du 19 septembre sur la ligne 39 au stop 0089. On remarque que la majorité des valeurs de délais se situe entre une minute d'avance et une minute de retard.

2.3 Préparation des données

Les lignes du dataset pour lesquelles la colonne `delay` n'avait pas de valeur ont été supprimées. Les lignes dont la valeur de la colonne `line` n'était pas égale à 39 ont également été supprimées. Les *features* ayant toujours la même valeur ont été supprimées. Les colonnes `trip`, `theoretical_time`, `expectedArrivalTime` et `date` ont été supprimées car elles n'ont pas été jugées utiles. Les colonnes `theoretical_time`, `expectedArrivalTime` et `date` ont été supprimées car les valeurs étaient du type `string`. La colonne `date` a été considérée comme redondante, sa valeur étant déjà stockée dans les colonnes `hour`, `minute` et `day`.

2.3.1 Stratification des données

La stratification des données ne nous a pas semblé utile car le dataset n'est assez grand pour rendre cette dernière nécessaire. Cependant dans un but pédagogique nous avons quand même stratifié la colonne `hour`, afin que la répartition des différentes valeurs reste identique dans le dataset ainsi que dans le test-set.

2.4 Visualisation des données

2.4.1 Création du test-set

Cela pourrait paraître étrange de mettre de côté une partie des données à ce moment. Les données n'ont même pas encore été vraiment visualisées et nous devons encore en apprendre plus avant de choisir quel algorithme utiliser. Cependant si le test-set est créé maintenant c'est pour éviter le *snooping bias*. Le test-set est constitué de vingt pourcent des données du dataset.

2.4.2 Premières visualisations

La première visualisation générée (annexe B.3) est la répartition des délais en fonction de l'heure. On remarque qu'après dix-huit heures on a soit une avance ou un retard de vingt minutes ou une variation de cinq minutes par rapport à l'horaire théorique, sans valeur

intermédiaire. On remarque également que la majorité des délais ont une valeur nulle. On peut également voir qu'il y a une augmentation des délais après quinze heures jusqu'à dix-neuf heures.

Le second graphique (annexe B.4) indique le délai moyen par heure. On y remarque une augmentation des délais entre six et neuf heures, à quinze heures et ainsi qu'à vingt-et-une heures. Le pic de retard de vingt-et-une heures vient sans doute du couvre-feu de vingt-deux heures, les autres pics quand à eux sont à priori dus au trafic de Bruxelles.

Sur la dernière infographie (annexe B.5) on remarque une hausse des températures de midi à seize heures.

2.4.3 Recherche de corrélation

Étant donné que notre dataset n'est pas trop grand nous pouvons facilement calculer le *coefficient standard de corrélation* entre chaque paire de *feature*. Comment on peut le voir avec le bout de code suivant (Voir listing 2.1), le retard est très peu linéairement corrélé avec les autres *features*.

```
1 >>> corr_matrix = data.corr()
2 >>> corr_matrix["delay"].sort_values(ascending=False) # warning: this
    check only linear correlation
3 delay      1.000000
4 rain       0.018600
5 temp       0.009013
6 wind       0.002024
7 hour      -0.000627
8 minute     -0.006874
9 humidity   -0.016347
10 visibility -0.034723
11 day       -0.079071
12 Name: delay, dtype: float64
```

Listing 2.1 – Coefficient standard de corrélation pour la *feature* delay.

La *heatmap* (disponible en annexe C.1 et générée avec le code 2.2) nous permet de confirmer que la *feature* delay n'est pas linéairement corrélée avec les autres *features*. Elle nous fournit cependant des informations supplémentaires comme le fait que la *feature* temp est linéairement corrélée avec la *feature* wind, ainsi qu'inversément linéairement corrélée avec la *feature* humidity. Cependant ces corrélations restent faibles, comme on peut le voir sur la figure C.2. On voit bien que globalement quand la valeur de la *feature* humidity chute quand celle de la *feature* temp augmente mais on reste cependant loin d'une belle ligne droite.

```
1 corr_matrix = data.corr()
2 plt.figure(figsize=(10,10), dpi=100)
3 sns.heatmap(corr_matrix)
```

Listing 2.2 – Matrice des coefficients standard de corrélation.

2.5 Combinaison de *features*

La *feature* `hour_minute` est une nouvelle *feature* que nous avons créé en combinant les *features* `hour` et `minute`.

Annexe A

Position des stops

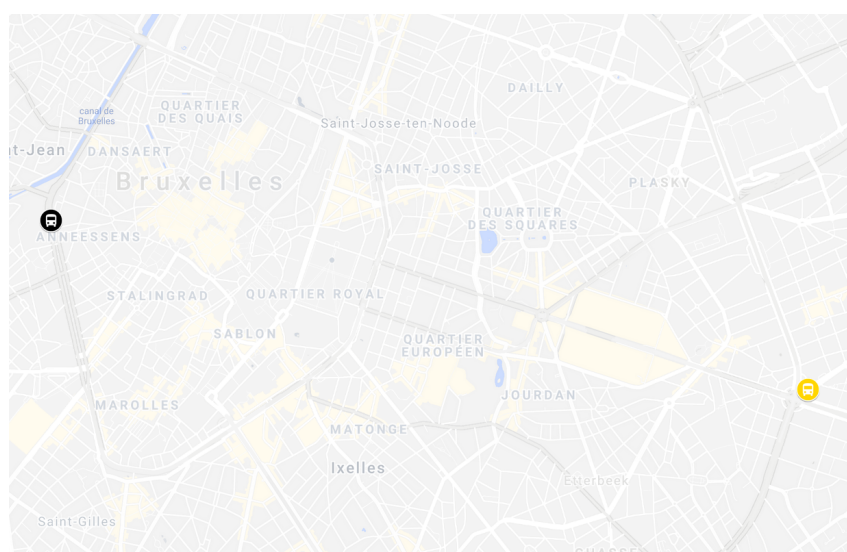


FIGURE A.1 – Position des stops numéro 0089 et 6608G en jaune et noir respectivement.

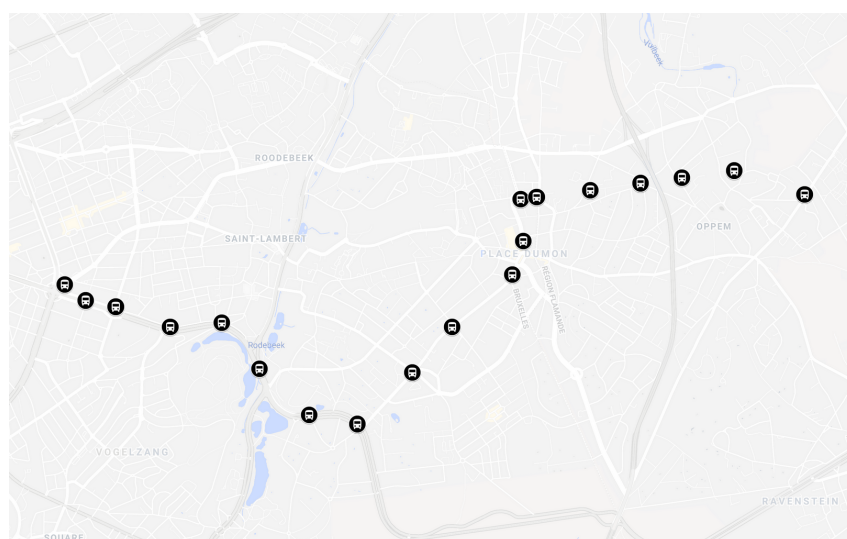


FIGURE A.2 – Position des stops présent sur la ligne 39

Annexe B

Première visualisation des données

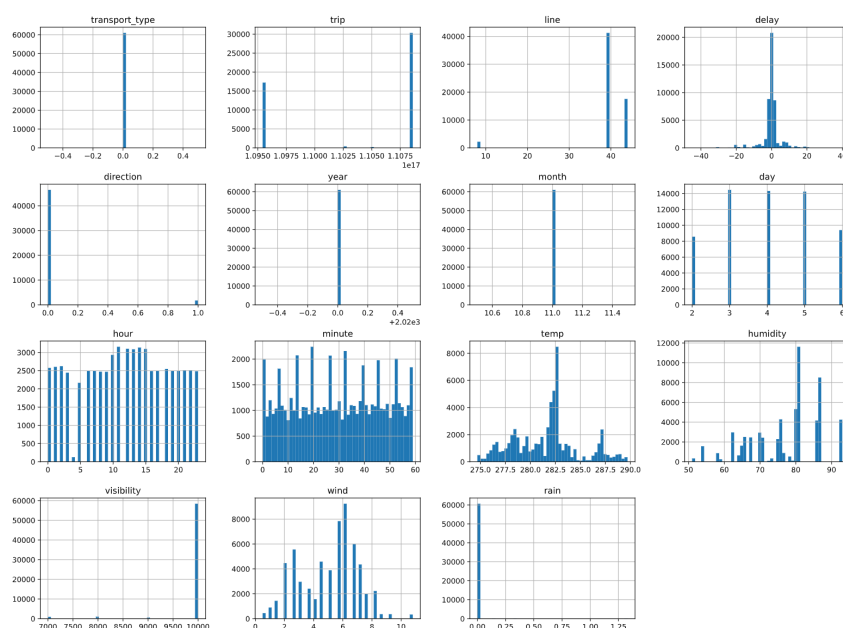


FIGURE B.1 – Plot des features par nombre d'occurrences

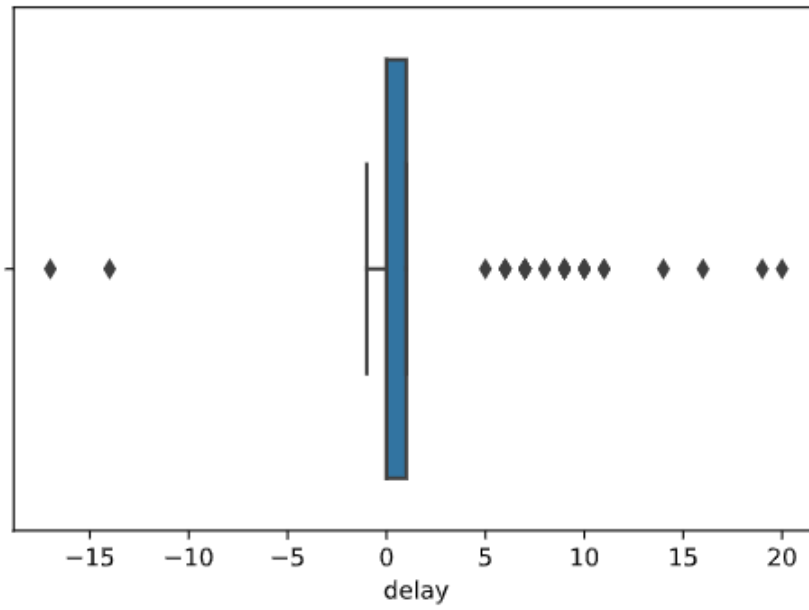


FIGURE B.2 – Boxplot des délais du 19 septembre sur la ligne 39 au stop 0089

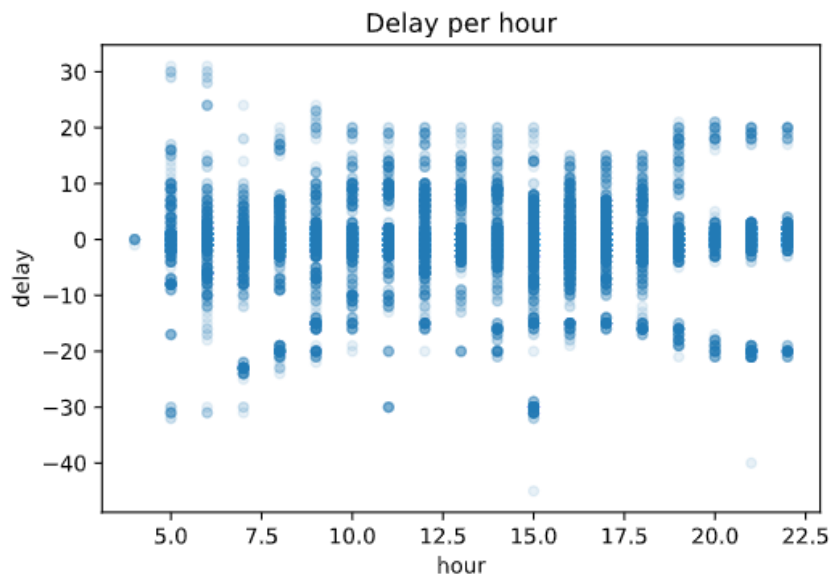


FIGURE B.3 – Délai par heure pour tout le dataset

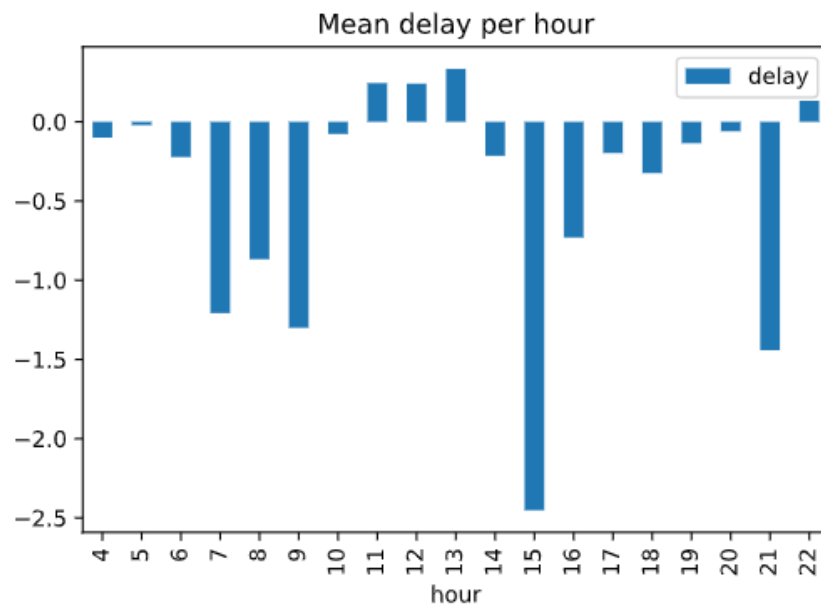


FIGURE B.4 – Délai moyen par heure pour tout le dataset

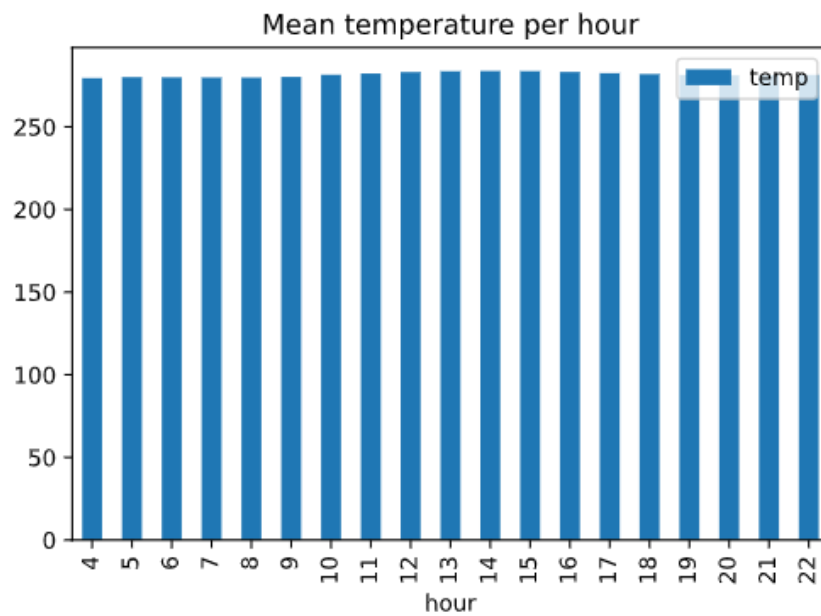


FIGURE B.5 – Température moyenne par heure pour tout le dataset

Annexe C

Recherche de corrélation

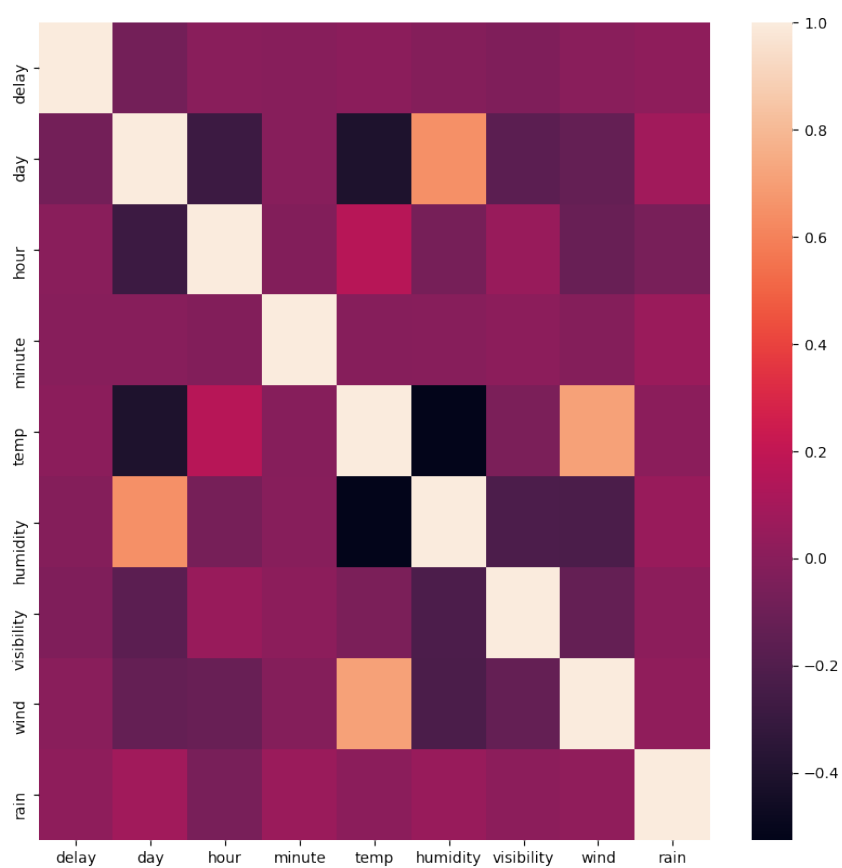


FIGURE C.1 – Matrice de corrélation linéaire entre les différentes *features*

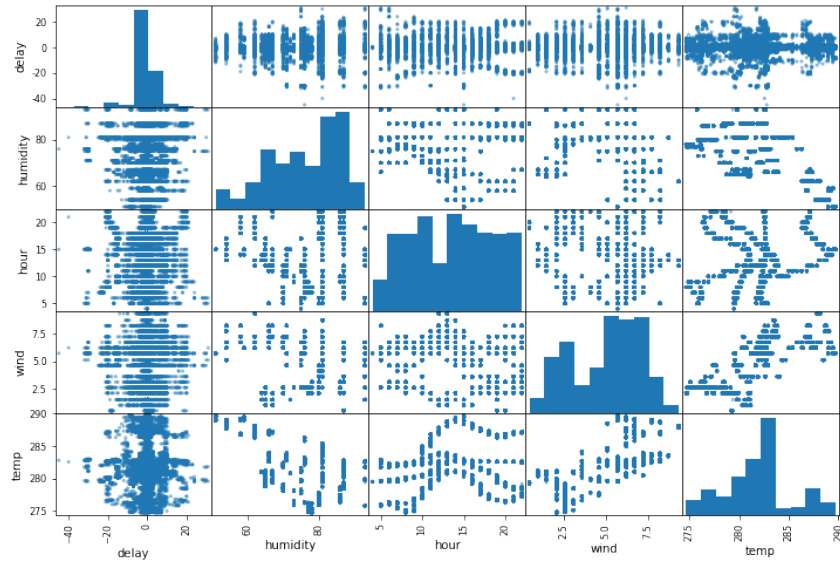


FIGURE C.2 – Plot de chaque *features* par rapport à chacune des autres



FIGURE C.3 – Bonus

Listings

2.1	Coefficient standard de corrélation pour la <i>feature</i> <code>delay</code>	5
2.2	Matrice des coefficients standard de corrélation.	5

Bibliographie

- [1] Robert R. F. DEFILIPPI. Standardize or Normalize Examples in Python. Avr. 2018. URL : <https://medium.com/@rrfd/standardize-or-normalize-examples-in-python-e3f174b65dfc>.
- [2] Aurélien GÉRON. Hands-on machine learning. Avr. 2017. URL : <http://www.amazon.com/exec/obidos/redirect?tag=citeulike07-20%5C&path=ASIN/1491962291>.
- [3] Urvashi JAITLEY. Why Data Normalization is necessary for Machine Learning models. Avr. 2019. URL : <https://medium.com/@urvashilluniya/why-data-normalization-is-necessary-for-machine-learning-models-681b65a05029>.
- [4] machine learning in Python scikit-learn documentation. URL : <https://scikit-learn.org/stable/>.
- [5] Python Data Analysis Library. URL : <https://pandas.pydata.org/>.
- [6] Emre RENÇBEROĞLU. Fundamental Techniques of Feature Engineering for Machine Learning. Avr. 2019. URL : <https://towardsdatascience.com/feature-engineering-for-machine-learning-3a5e293a5114>.
- [7] Natasha SHARMA. Ways to Detect and Remove the Outliers. Mai 2018. URL : <https://towardsdatascience.com/ways-to-detect-and-remove-the-outliers-404d16608dba>.
- [8] statistical data visualization seaborn documentation. URL : <https://seaborn.pydata.org/>.
- [9] Zixuan ZHANG. Understand Data Normalization in Machine Learning. Août 2019. URL : <https://towardsdatascience.com/understand-data-normalization-in-machine-learning-8ff3062101f0>.