

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
2.4.4	Combinaison de <i>features</i>	5
2.5	Les <i>features</i> de type <code>string</code>	5
2.6	Pipeline	6
3	Modèles et résultats	7
3.1	Partie 1 - Regression	7
3.1.1	Cross-validation	7
3.1.2	<i>Grid search</i>	7
3.1.3	<i>Randomized Search</i>	8
3.1.4	Intervale de confiance	8
	Appendices	9
A	Position des stops	9
B	Première visualisation des données	10
C	Recherche de corrélation	13
D	Regression	15
E	Bonus	16

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 fournies 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 détection d'*outliers*, la sélection 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 une phase de sélection 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 aux défis de la mobilité. Plusieurs *opendata* nous étaient 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()` indique 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 valeur. On peut également voir que des retards on été enregistré 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 du au fonctionnement de l'API de la STIB. On peut également

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](#)

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 les quelles 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 sont pas été jugée utile. 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 stratifier 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 visualisée et nous devons encore en apprendre plus avant de choisir quel algorithmes 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 heure 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 heure jusqu'à dix-neuf heure.

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

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

2.4.3 Recherche de corrélation

Étant donné que le 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) 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.

2.4.4 Combinaison de *features*

La *feature* `hour_minute` est une nouvelle *feature* que nous avons créée en combinant les *features* `hour` et `minute`. Ce regroupement des *features* est fait dans un but pédagogique car de manière générale on a plutôt tendance à séparer les *features*. La séparation des *features* est une bonne façon de les rendre plus utilisables par les algorithmes de machine learning car la plus part du temps les datasets possèdent des colonnes du type `string` qui violent le principe de *tidy data*. La division des *features* permet d'améliorer les performances du modèle en découvrant des informations potentielles. Ce qui a déjà été fait en divisant la *feature* `date` en les *features* `year`, `month`, `day`, `hour` et `minute`.

Visualisation de la nouvelle *feature*

Le graphique en annexe C.3 nous montre la répartition des délais en fonction de l'heure de la journée. On remarque que le graphique garde évidemment sa forme de fourche caractéristique.

En réaffichant les coefficients standard de corrélation en prenant en compte cette nouvelle *feature*, on remarque qu'elle est encore moins linéairement corrélée avec la *feature* `delay` que ne le sont les *features* `hour` et `minute`.

2.5 Les *features* de type `string`

La plus part des algorithmes de machine learning préfèrent travailler avec des nombres plutôt qu'avec du texte, c'est pourquoi nous convertissons la *feature* `stop` en différentes catégories repré-

sentée par un nombre.

2.6 Pipeline

Un pipeline a été créé pour faciliter l'exécution des étapes de transformation des données. Le pipeline effectue les opérations suivantes, premièrement amputer les lignes ayant au moins une valeur nulle dans l'une des colonnes. En suite vient l'étape d'ajout de la *feature* `hour_and_minute` qui peut être désactivée à la volée à l'aide d'un paramètre. La dernière étape est la création des catégories des stops comme vu dans la section précédente.

Chapitre 3

Modèles et résultats

Ce chapitre est divisé en deux parties la première concerne la partie de l'équipe "regression" (constituée de Jean A. et de Cyprien L.) et la seconde concerne la partie de l'équipe "classification" (constituée d'Antoine P. et de Jessica D.). Chacune des équipes est partie sur une vision différente du problème, comme leur nom l'indique la première équipe a vu le problème comme un problème de régression tandis que la seconde comme un problème de classification.

3.1 Partie 1 - Regression¹

Différents modèles ont été entraînés afin de sélectionner le meilleur d'entre eux. Le premier à avoir été testé est le modèle de régression linéaire, les prédictions obtenues sont disponibles dans le tableau en annexe D.1. Le modèle fonctionne mais n'est pas précis du tout. Pour nous rendre compte à quel point le modèle se trompe nous calculons la *root-mean-square error* (RMSE), cette dernière est égale à 5.2243. C'est mieux que rien mais étant donné que les délais varient principalement entre 20min et -25min, une erreur moyenne de 5min n'est pas acceptable. Cela signifie que le modèle *underfit* les données d'entraînement ce qui indique que nos données ne sont pas assez pertinentes ou que le modèle n'est pas assez puissant. Lors que l'on utilise donc un modèle d'apprentissage par arbre de décision et que l'on calcule la RMSE on obtient une valeur nulle. Ce qui signifie que le modèle *overfit* les données d'entraînement.

3.1.1 Cross-validation

Afin de mieux évaluer les modèles testés, nous avons utilisé une *10-fold cross-validation*. En utilisant la cross validation sur le modèle linéaire nous obtenons toujours une erreur moyenne de 5.22 pour les dix *folds*. Le modèle d'apprentissage par arbre de décision obtient lui une valeur de 6.13, pire que le modèle de régression linéaire donc. Ce qui indique que l'arbre de décision *overfit* tellement qu'il fonctionne moins bien que le modèle de régression linéaire. Après avoir entraîné un modèle de type *Random Forest* et du type *Support Vector Regression* sur 10 *folds* nous obtenons respectivement une erreur moyenne de 4.89 et de 5.27. Pour la suite, nous avons choisi le modèle *Random Forest*.

3.1.2 Grid search

Nous avons donné une liste de valeurs pour chaque métaparamètre (`n_estimators`, `max_features` et `bootstrap`) à tester. *Grid search* nous indique que la meilleure combinaison de valeurs est {'

1. Le code concernant cette section est disponible [ici](#)

`max_features': 8, 'n_estimators': 30}` avec une erreur moyenne de 4.758. Ce résultat peut encore être amélioré car la valeur des métaparamètres `n_estimators` et `max_features` sélectionnées ont tous deux la valeur maximale que nous avons donnée. Nous changeons donc la liste des valeurs de `n_estimators` et `max_features` pour des valeurs plus grandes que 30 et 8 respectivement. La meilleur combinaison était `{'max_features': 8, 'n_estimators': 70}` avec une erreur moyenne de 4.726.

3.1.3 *Randomized Search*

Nous avons testé une *Randomized Search* avec des valeurs entre 1 et 200 pour `max_features` et entre 1 et 8 `n_estimators`. La meilleur combinaison de métaparamètres est `{'max_features': 7, 'n_estimators': 180}` avec une erreur moyenne de 4.719.

3.1.4 Intervale de confiance

L'intervale de confiance (95%) pour la RMSE est entre 3,84 et 4,63min

Annexe A

Position des stops

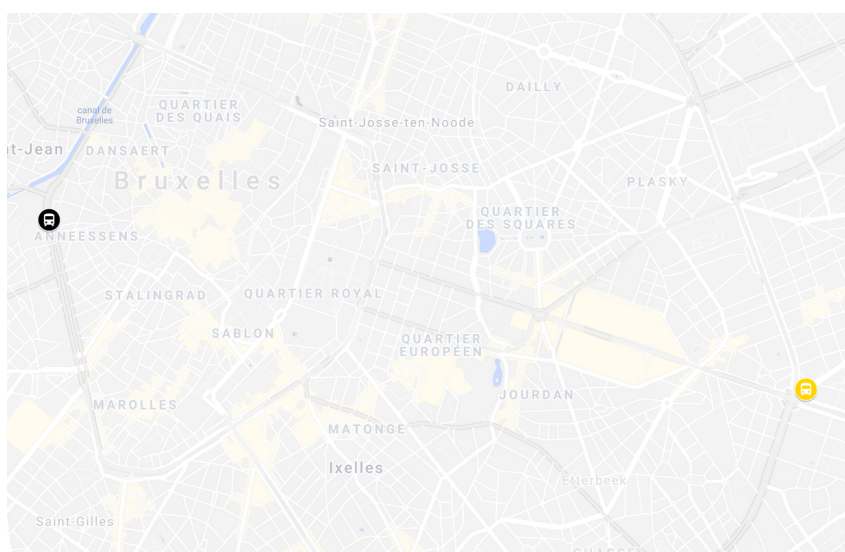


FIGURE A.1 – Position des stops numéro 0089 et 6608G en jaune et noir respectivement.

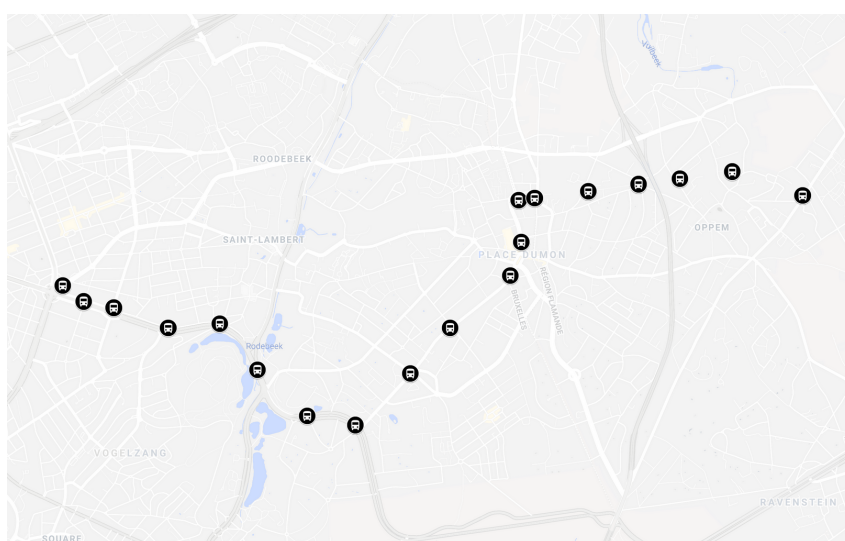


FIGURE A.2 – Position des stops présent sur la ligne 39

Annexe B

Première visualisation des données

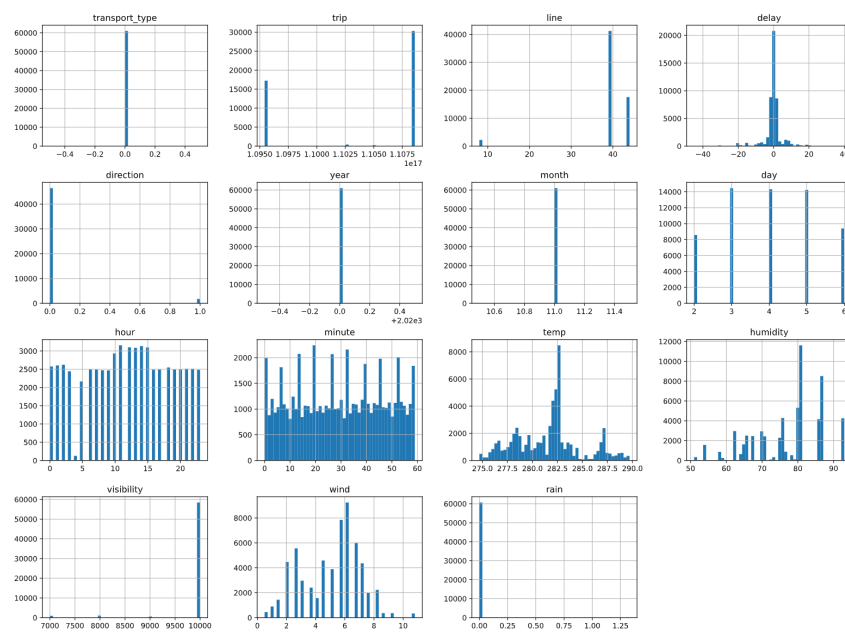


FIGURE B.1 – Plot des features par nombre d'occurrences

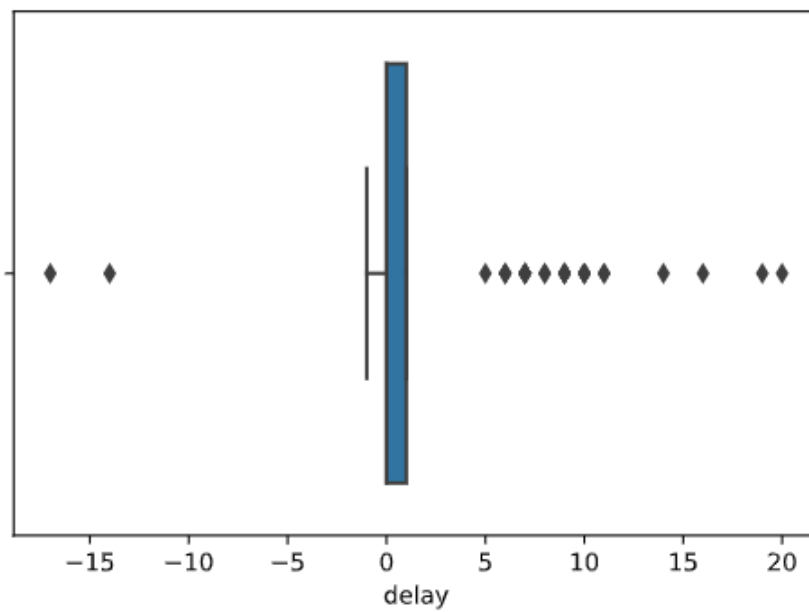


FIGURE B.2 – Boxplot des délais du 19 septembre sur la ligne 39 au stop 0089

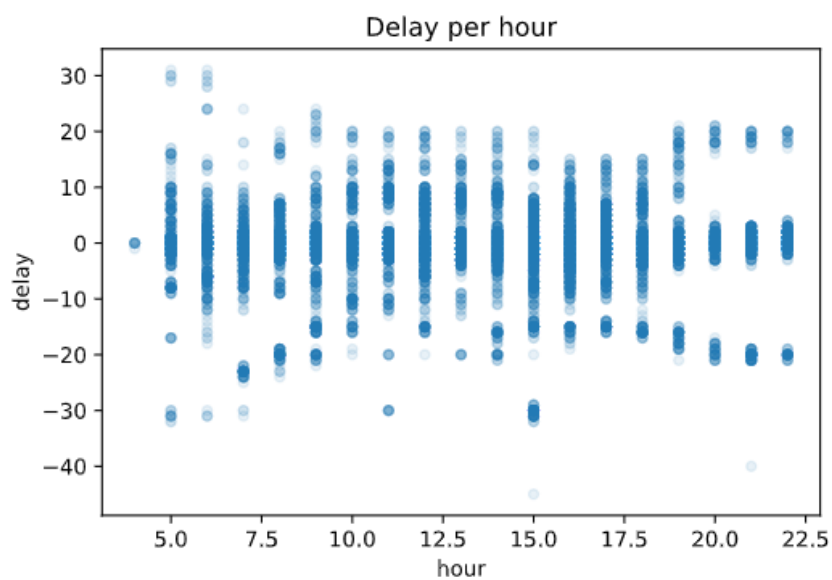


FIGURE B.3 – Délai par heure pour tout le dataset

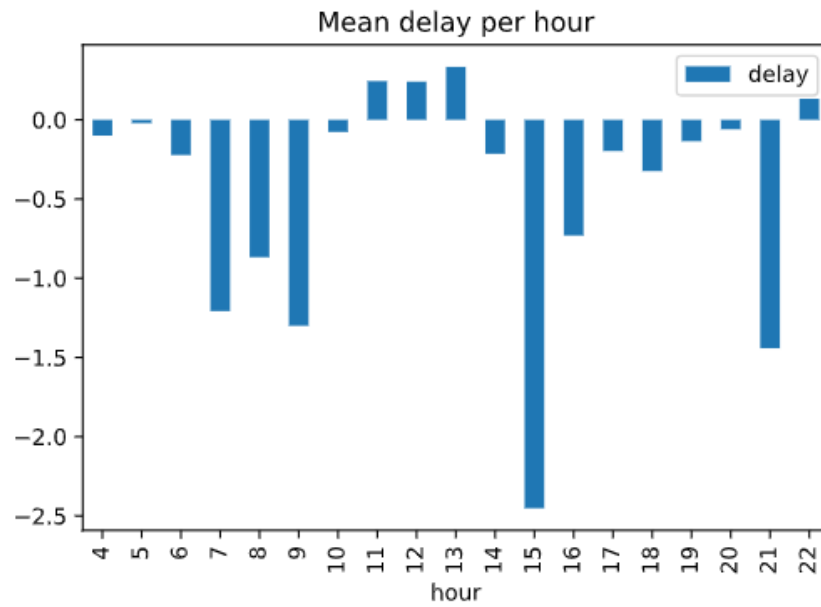


FIGURE B.4 – Délai moyen par heure pour tout le dataset

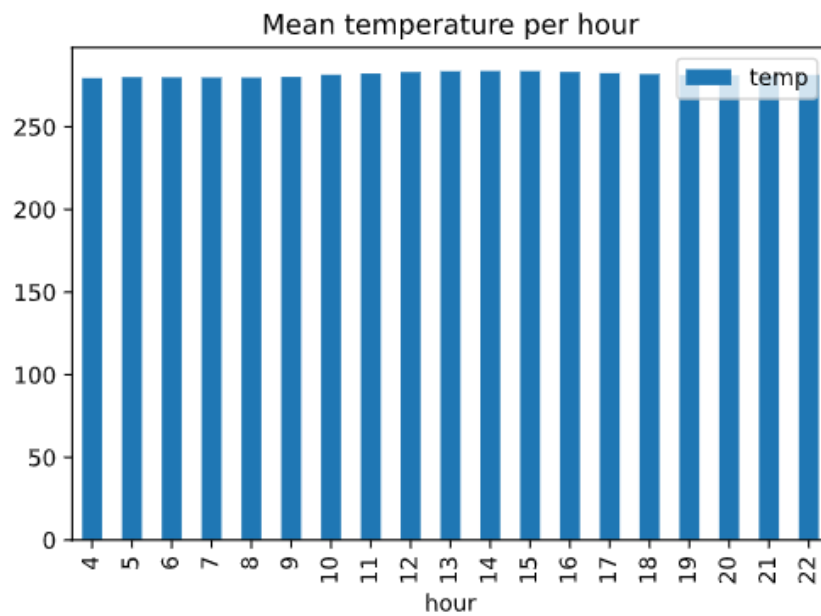


FIGURE B.5 – Température moyenne par heure pour tout le dataset

Annexe C

Recherche de corrélation

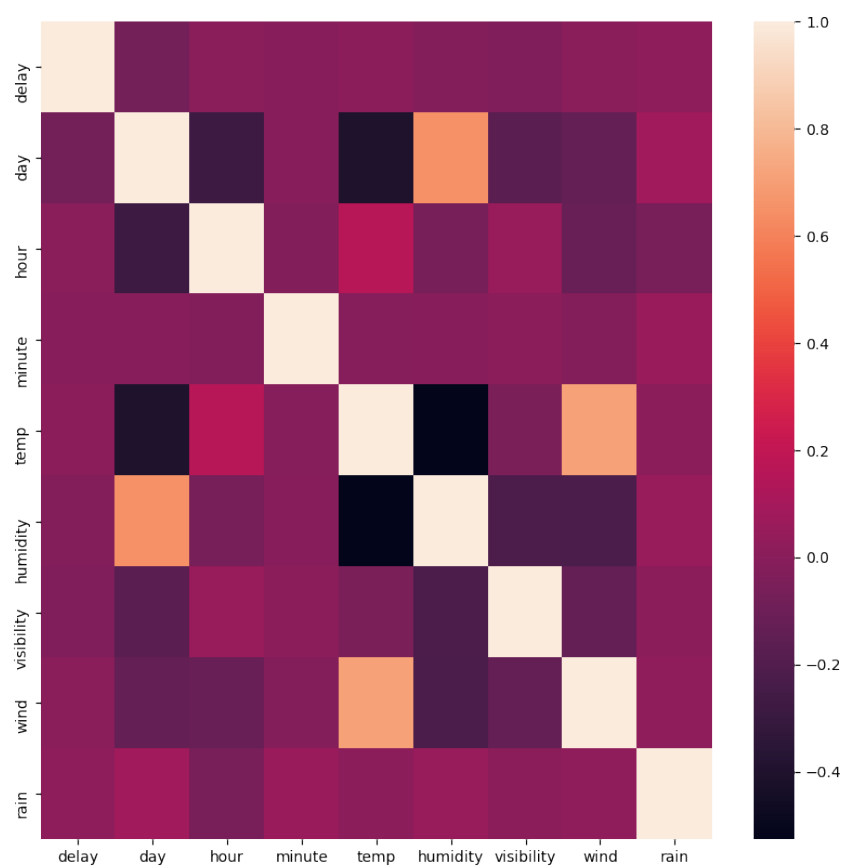


FIGURE C.1 – Matrice de corrélation linéaire entre les différentes *features*

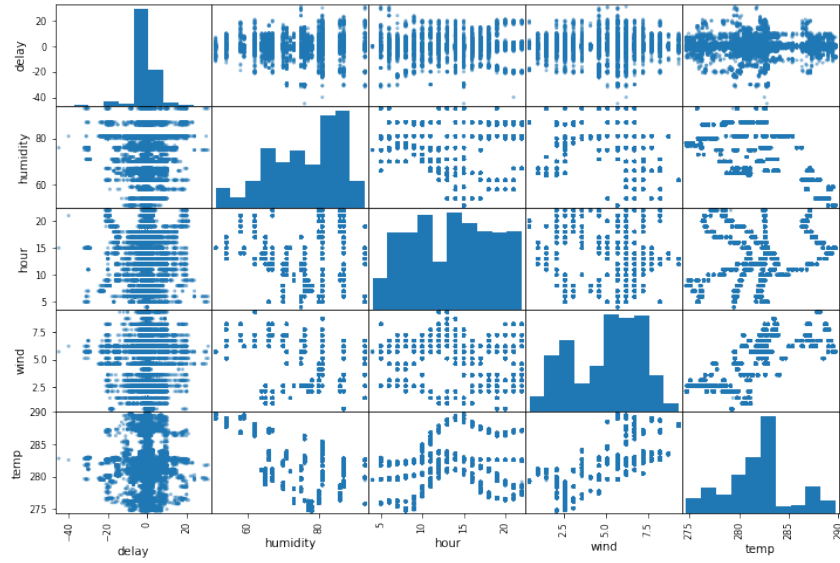


FIGURE C.2 – Plot de chaque *features* par rapport à chacune des autres

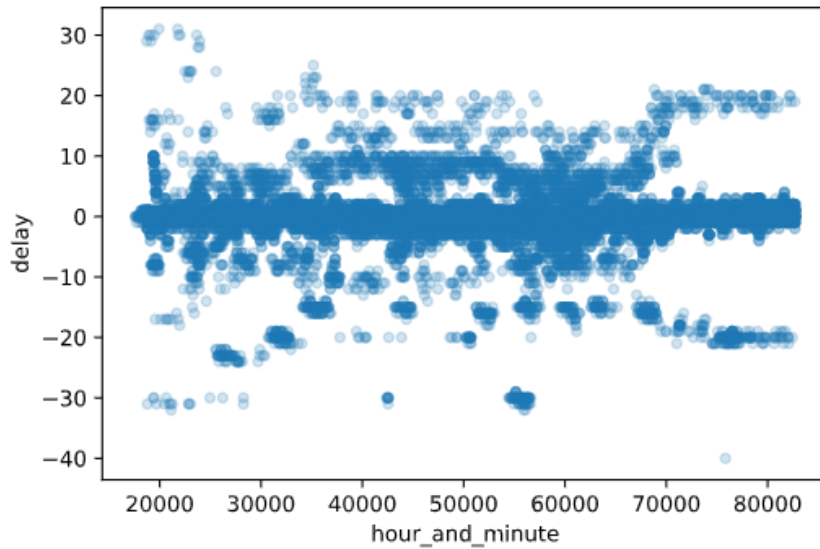


FIGURE C.3 – Délai par heure et minute pour tout le dataset

Annexe D

Regression

Valeur réelle	Prédiction	Différence
-1.0	1.02270508	-2,023
-1.0	-1.04394531	0,044
0.0	-1.015625	-1.016
-1.0	0.03173828	-1,032
-18.0	0.48291016	-18,483

TABLE D.1 – Échantillon de prédiction du modèle de régression linéaire

Annexe E

Bonus



FIGURE E.1 – Bonus

Listings

2.1	Coefficient standard de corrélation pour la <i>feature</i> <code>delay</code>	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] Hadley WICKHAM. « Tidy data ». In : The Journal of Statistical Software 59 (10 2014). URL : <http://www.jstatsoft.org/v59/i10/>.
- [10] Zixuan ZHANG. Understand Data Normalization in Machine Learning. Août 2019. URL : <https://towardsdatascience.com/understand-data-normalization-in-machine-learning-8ff3062101f0>.