

Soutenance du Projet 4

Anticipez les besoins en consommation électrique de bâtiments

Par
Elisée TCHANA

Mentor
Cyril MONTI

Avril 2021

Plan de Soutenance

La présentation de l'application

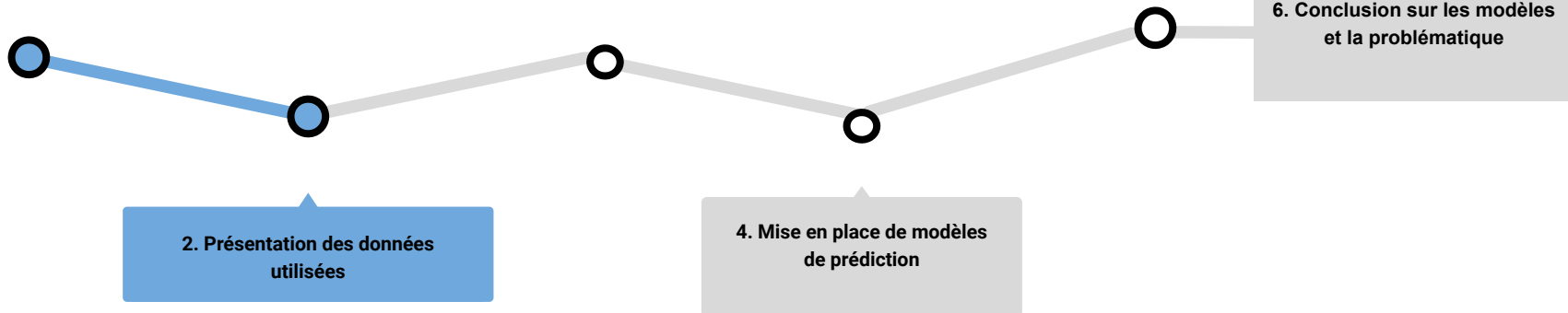
2. Présentation des données
utilisées

3. Nettoyage de la base de
données

4. Mise en place de modèles
de prédiction

5. Comparaison des différents
modèles

6. Conclusion sur les modèles
et la problématique



Problématique du Projet

1. Prédire les performances énergétiques de bâtiments situés à Seattle
2. Se passer des relevés faits sur place (opérations très coûteuses et fastidieuses)



Seattle

Objectifs

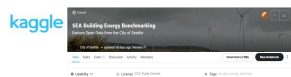
Prédire la consommation en énergie des bâtiments

Prédire les émissions en CO2 des bâtiments

Evaluer l'intérêt de « l' ENERGYSTARScore »



1. Indicateur à l'échelle nationale permettant de refléter les performances énergétiques d'un bâtiment
2. Score allant de 1 (mauvaise performance) à 100 (excellente performance)
3. Un score de 50 représente la médiane national

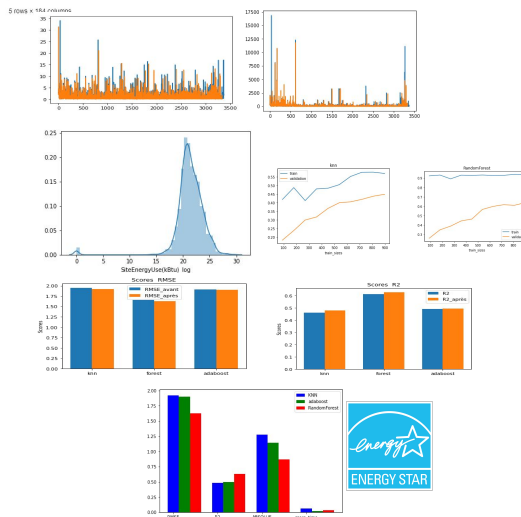


**SEA Building Energy
Benchmarking
Explore Open Data
from the City of
Seattle**

```
#Import packages
import numpy as np
import pandas as pd
import os
import matplotlib.pyplot as plt
import seaborn as sns
import math
import pickle

from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix
from sklearn.preprocessing import StandardScaler
from sklearn import svm, datasets

#Data set
pd.read_csv('openfoldtest.csv', sep='|')
data.head()
```

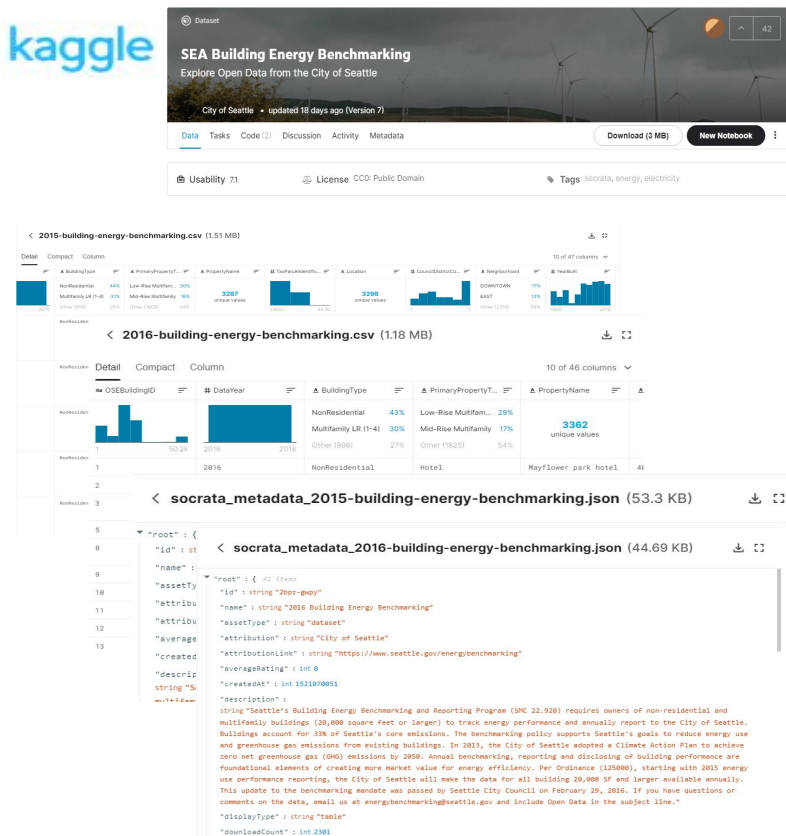
[illegible]

Data Cleaning & Analysis



Data Sharing

Pr sentation du jeu de donn es



Donn es de 2015 :

- 3340 lignes
- 47 colonnes

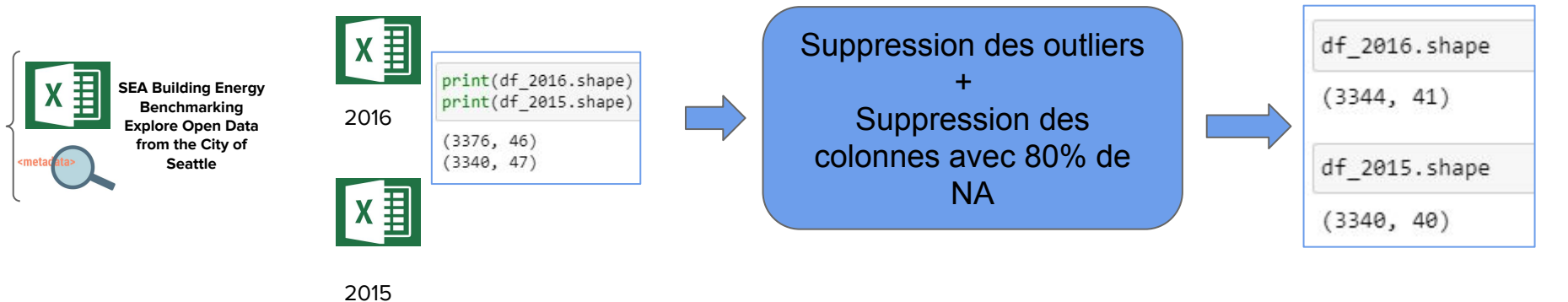
```
print(df_2016.shape)
print(df_2015.shape)
```

Donn es de 2016 :

- 3376 lignes
- 46 colonnes

```
(3376, 46)
(3340, 47)
```

Nettoyage des données - 1



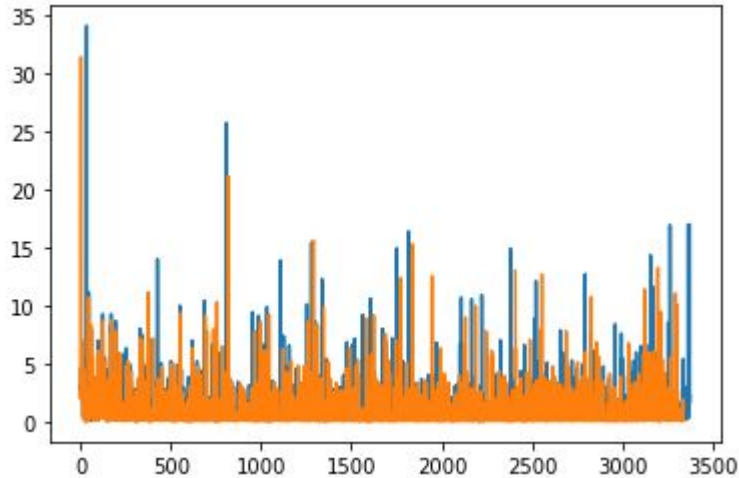
colonnes présentes sur données 2015 et absentes sur 2016 :

```
{'GHGEmissionsIntensity(kgCO2e/ft2)', 'Seattle Police Department Micro Community Policing Plan Areas', 'SPD Beats', 'GHGEmissions(MetricTonsCO2e)', 'OtherFuelUse(kBtu)', 'Zip Codes'}
```

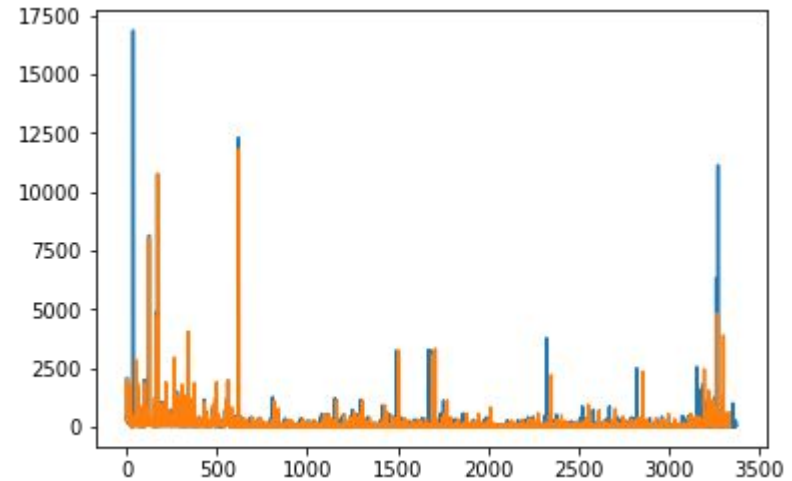
colonnes présentes sur données 2016 et absentes sur 2015 :

```
{'GHGEmissionsIntensity', 'TotalGHGEmissions'}
```

Nettoyage des données - 2



```
2016['GHGEmissionsIntensity'])  
2015['GHGEmissionsIntensity(kgCO2e/ft2)']
```



```
2016['TotalGHGEmissions'])  
2015['GHGEmissions(MetricTonsCO2e)']
```

Nettoyage des données - 3

- Dans la colonne **buildingtype** : garder seulement les « non residential », « sps-district k12 », « Campus ».
- Dans certaines colonnes comme « Neighborhood » rassembler les mots avec ou sans majuscule comme : north & NORTH
- Remplacer les dernières valeurs manquantes par un KNNimputer

```
df['Neighborhood'].replace('North', 'NORTH', inplace=True)
df['Neighborhood'].replace('CENTRAL', 'Central', inplace=True)
df['Neighborhood'].replace('Northwest', 'NORTHWEST', inplace=True)
df['Neighborhood'].replace('Ballard', 'BALLARD', inplace=True)
df['Neighborhood'].replace('DELRIDGE NEIGHBORHOODS', 'DELRIDGE', inplace=True)
df['Neighborhood'].replace('Delridge', 'DELRIDGE', inplace=True)
df['City'].replace('SEATTLE', 'Seattle', inplace=True)
df['ComplianceStatus'].replace('Error - Correct Default Data', np.nan, inplace=True)
df['ComplianceStatus'].replace('Not Compliant', 'Non-Compliant', inplace=True)
df['ComplianceStatus'].replace('Missing Data', np.nan, inplace=True)
```

TotalGHGEmissions

SiteEnergyUse(kBtu)

Analyse et exploitation

1. Quelle énergie est utilisée : steam, natural et stream
2. Quelle énergie est la plus utilisée : steam, natural ou electric
3. L'âge des bâtiments = datayear – yearbuild

```
df['AgeBat'] = df['DataYear']-df['YearBuilt']
df = df.drop(['DataYear','YearBuilt'],axis=1)

#créer 3 colonne avec des 0
df['most_use_stream'] = 0
df['most_use_natural'] = 0
df['most_use_electric'] = 0

### L'énergie qui est la plus utilisé (electricité/natural/streamuse), rajouter 1
df.loc [(df['SteamUse(kBtu)'] > df['Electricity(kBtu)']) & (df['SteamUse(kBtu)']
df.loc [(df['NaturalGas(kBtu)'] > df['Electricity(kBtu)']) & (df['NaturalGas(kBtu)']
df.loc [(df['Electricity(kBtu)'] > df['NaturalGas(kBtu)']) & (df['Electricity(kBtu)']

# quels sont les énergies utilisés
df.loc [df ['SteamUse(kBtu)'] == 0, 'SteamUse(kBtu)'] = 0
df.loc [df ['SteamUse(kBtu)'] != 0, 'SteamUse(kBtu)'] = 1
df.loc [df ['Electricity(kBtu)'] == 0, 'Electricity(kBtu)'] = 0
df.loc [df ['Electricity(kBtu)'] != 0, 'Electricity(kBtu)'] = 1
df.loc [df ['NaturalGas(kBtu)'] == 0, 'NaturalGas(kBtu)'] = 0
df.loc [df ['NaturalGas(kBtu)'] != 0, 'NaturalGas(kBtu)'] = 1

### renommer pour ne pas croire que j'utilise des données fuites
df = df.rename(columns = {'SteamUse(kBtu)': 'use_stream', 'Electricity(kBtu)': '
                        'NaturalGas(kBtu)': 'use_natural'})
#df=df.drop(['SteamUse(kBtu)', 'Electricity(kBtu)', 'NaturalGas(kBtu)'],axis=1)
```

Suppression des données qui sont pas utilisables

```
df = df.drop(['OSEBuildingID', 'Address', 'City', 'State', 'TaxParcelIdentificationN
            'PropertyName', 'SourceEUI(kBtu/sf)', 'PropertyGFATotal', 'TotalGHGEmi
```

Analyse et exploitation

Faire un pipeline pour faire un preprocessor et transformer les variables qualitatives en 0 et 1 grâce à `BinaryEncoder` et pour les variables quantitatives standardiser les données.

Tester plusieurs modèles différents en comparaison avec le modèle de base `DummyRegressor`.



Couper le dataset en deux parties : train set (80%) et un test set (20%)

Évaluer les modèles grâce à différentes métriques telles que : RMSE, R2, valeur absolue, le temps de traitement... avec une validation croisée (3) et 40n_iter.

RMSE : l'erreur quadratique moyenne, permet de voir la différence entre la valeur prédite et la valeur réelle.

R2 : c'est le coefficient de détermination (entre 0 et 1). Permet de juger la qualité de la régression linéaire.

Absolue : mesure des erreurs entre les observations prédites et réelles.

Fit time : le temps d'ajustement de l'estimateur sur l'ensemble de trains pour chaque fractionnement de CV.

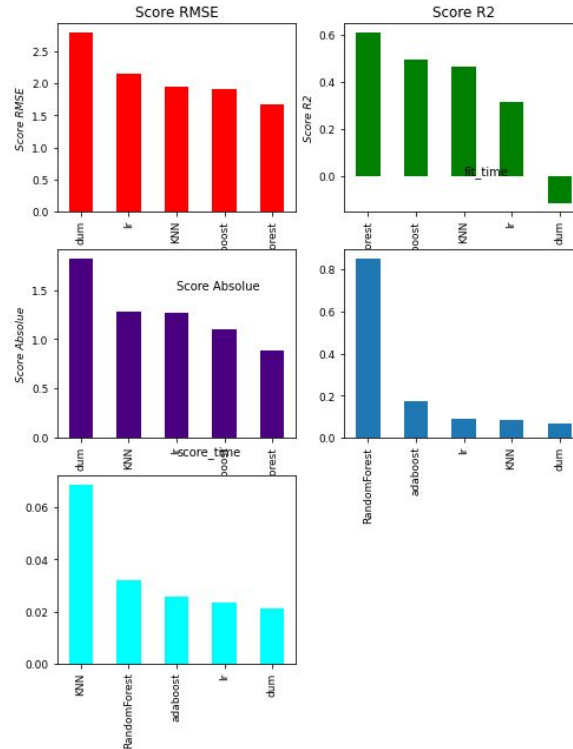
Score time : le temps de notation de l'estimateur sur l'ensemble de tests pour chaque fractionnement de CV.

TotalGHGEmissions



- Random Forest
- adaboost
- KNN

Résultat métrique des différents modèles



le modèle dummy de base à un R2 qui se rapproche de 0 c'est-à-dire que peu importe l'entrée il va sortir la même réponse

Le modèle avec le meilleur R2 et RMSE est le random forest mais c'est aussi celui qui met le plus de temps à s'entraîner

SiteEnergyUse(kBtu)

le modèle dummy de base à un R^2 qui se rapproche de 0 c'est-à-dire que peu importe l'entrée il va sortir la même réponse

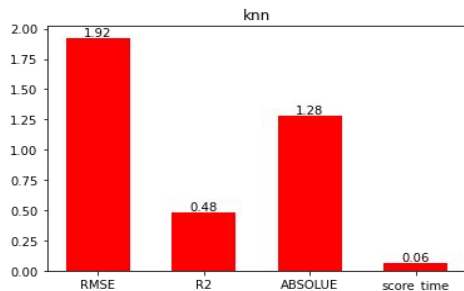
Le modèle avec le meilleur R^2 et RMSE est le random forest mais c'est aussi celui qui met le plus de temps à s'entraîner

Les 3 modèles retenus sont random forest/ adaboost et KNN.

Même conclusion que la target précédente mais avec des scores moins performants

Analyse et exploitation

KNN



Les meilleurs paramètres :

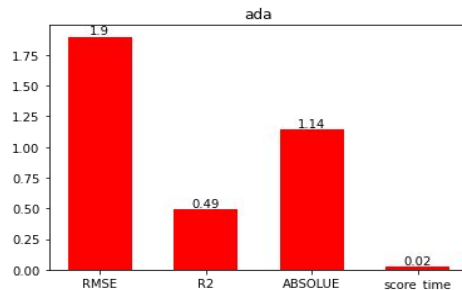
Algorithm = brute

Leaf size = 39

Neighbors = 10

Modèle basé sur les voisins les plus proches. Le score est prédit par interpolation local des cibles associées des voisins les plus proches dans l'ensemble d'apprentissage.

Adaboost



Les meilleurs paramètres :

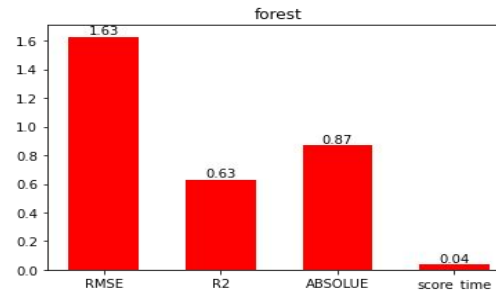
N_estimator = 95

Learning rate = 4

loss = exponential

Un régresseur adaboost, s'ajuste sur l'ensemble des données d'origine puis ajuste les copies supplémentaires du régresseur sur le même ensemble de données mais le poids varie en fonction des erreurs

Random Forest



Les meilleurs paramètres :

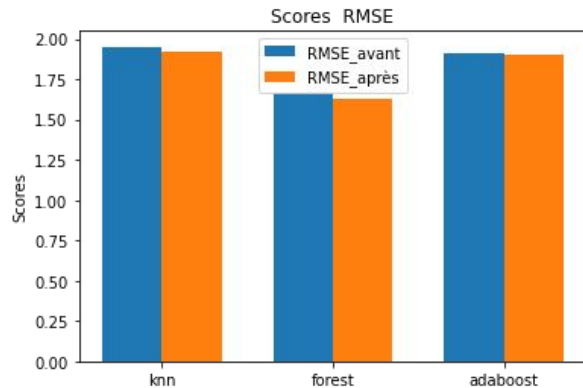
N_estimator = 115

Bootstrap = true

Max sample leaf = 1

C'est un méta-estimateur qui ajuste un certain nombre d'arbres de décision sous divers sous-échantillons et utilise la moyenne pour améliorer la précision prédictive et contrôler le sur-ajustement.

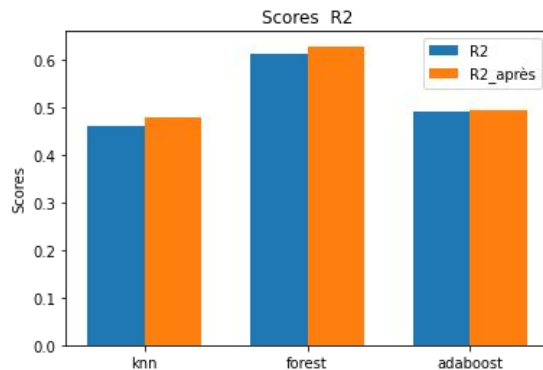
Analyse et exploitation



Plus le score est proche de 0, plus le modèle est performant

Les modèles sont plus performants en cherchant les meilleurs paramètres

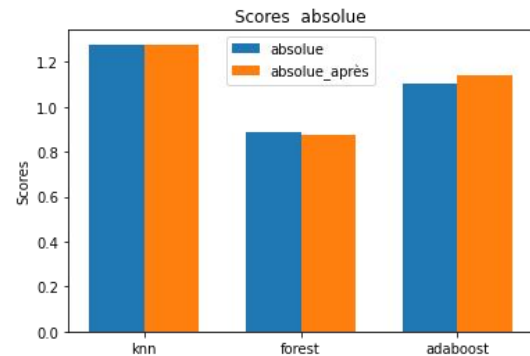
Le meilleur modèle par rapport au RMSE est le random forest



Plus le score est proche de 1, plus le modèle est performant

Les modèles sont plus performants en cherchant les meilleurs paramètres

Le meilleur modèle par rapport au R2 est le random forest



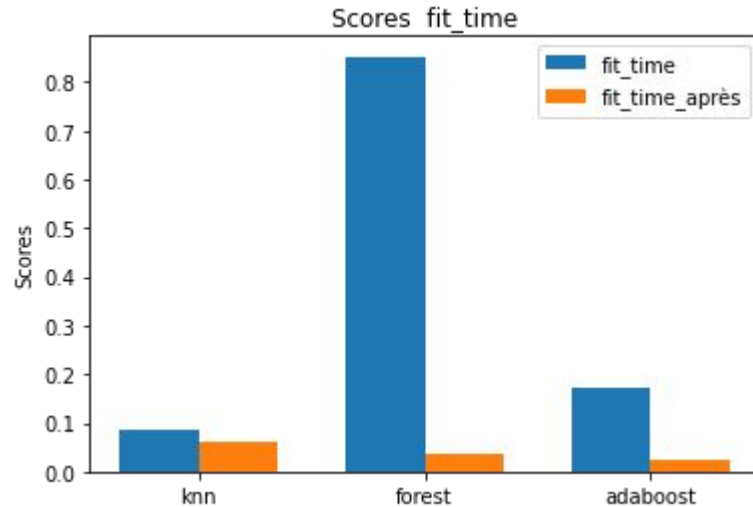
Plus le score est proche de 0, plus le modèle est performant

Les modèles sont plus performants en cherchant les meilleurs paramètres

Le meilleur modèle par rapport aux absolus est le random forest

Comparaison avant après

Analyse et exploitation

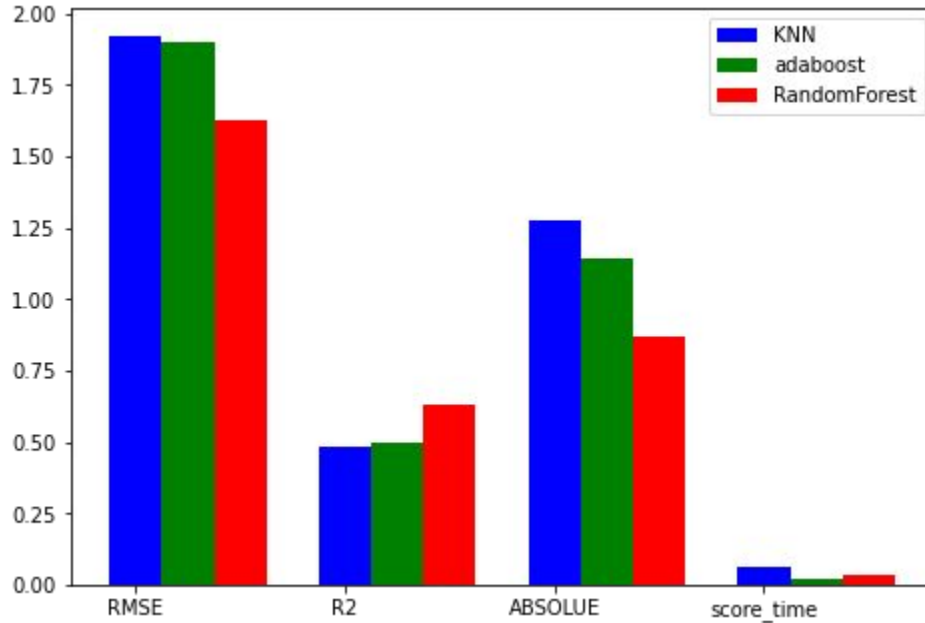


Lorsqu'on cherche les meilleurs paramètres, les modèles sont plus longs.

Le random forest est très long car $n_{\text{estimator}}$ est très grand et il y a beaucoup de combinaisons.

Comparaison avant après

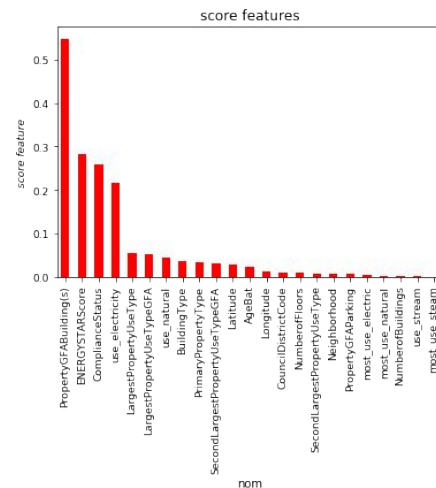
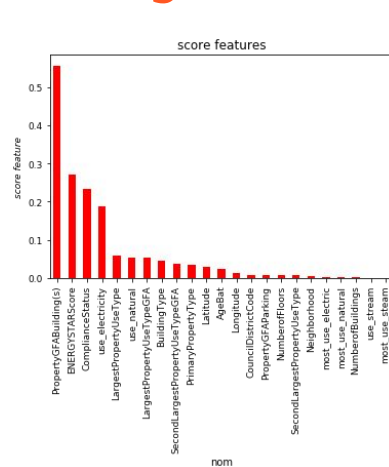
Analyse et exploitation



Le meilleur modèle par rapport aux metrics est le random forest.

Le random forest est le modèle le plus long en terme de temps d'entraînement

Analyse et exploitation



Avec ENERGYSTARScore

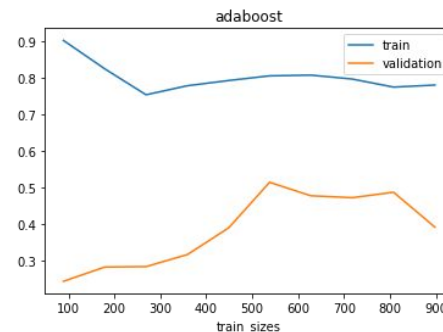
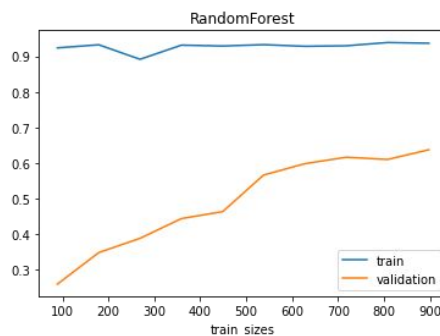
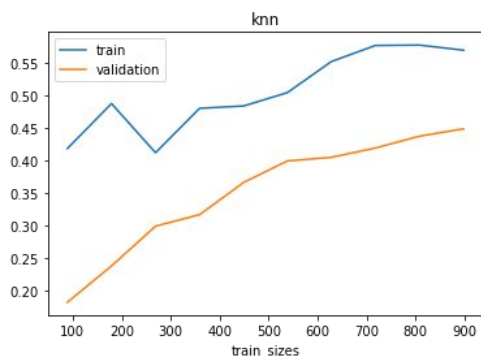
RMSE : 1,88

R2 : 0,66

Sans ENERGYSTARScore

RMSE : 1,89

R2 : 0,65



Les courbes nous montrent que des donn es suppl mentaires nous permettraient d'avoir surement de meilleurs r sultats.

Conclusion

- Le **ENERGYScore** ne joue pas un rôle très important dans l'amélioration des différents modèles et pour les deux targets.
- Le meilleur modèle pour les deux targets est le Random forest mais c'est aussi le plus long à exécuter.
- Malgré les datasets il manque un peu de données pour que les modèles soient plus efficaces et qu'ils soient plus performants.





Merci pour votre attention !!