



**ANTICIPEZ LES BESOINS  
EN CONSOMMATION DE  
BÂTIMENTS**

**Alexandre Delaguillaumie**



# PARTIE I

---

## Introduction

# Introduction

summary

---

**Seattle veut être une ville neutre**

En émission de carbone en 2050

**En réduisant la consommation et l'émission de gaz**

des bâtiments non destinés à l'habitation

**Les relevés des bâtiments sont coûteux à obtenir**

On cherche donc à s'en passer en les prédisant

# Introduction

problem

---

Réaliser un **algorithme** de **prédiction** des **émissions** de gaz à effet de serre et de la **consommation** d'énergie des **bâtiments non destinés à l'habitation** qui n'ont pas encore été mesurés

# Introduction

few things to know

---

**Attention à la fuite de données**

car les variables energetiques ne sont pas connues.

**Un 1er relevé de référence est effectué**

Pour tous les bâtiments.

**On cherche à évaluer l'intérêt de l'EnergyStarScore**

Car il est fastidieux à calculer actuellement.



# PARTIE II

---

## Exploration



# Exploration

## overview

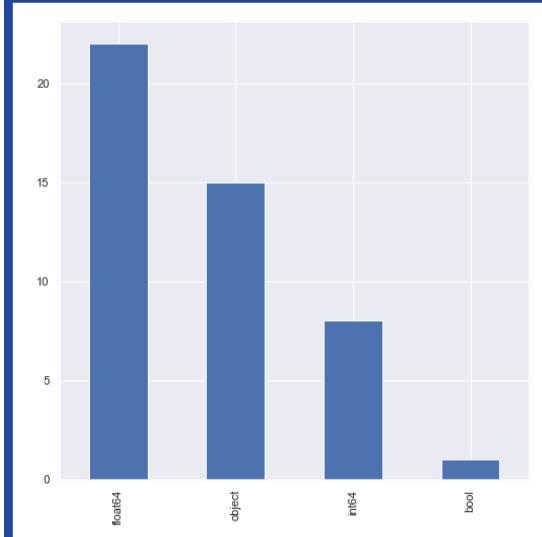
### FILL RATE du dataset brut



**3376 BÂTIMENTS**

Dont 1583 non résidentiels

### RÉPARTITION DES VARIABLES par type



**45 variables**

Dont 20 d'intérêt

# Exploration

targets & features

## Features

**Nombre d'étages**  
Number of Floors

**Nombre de bâtiments**  
Number of Buildings

**Superficie brute**  
PropertyGFABuilding  
SecondLPUT GFA  
ThirdLPUT GFA

**Emplacement**  
Latitude / Longitude

**Ancienneté**  
YearBuilt

**Quartier**  
Neighborhood

**Type de bâtiment**  
LargestPropertyUseType  
SecondLargestPropertyUT  
ThirdLargestPropertyUT

**Superficie**  
LargestPropertyUseGFA

**Energies utilisées**  
NaturalGas,  
SteamUse,  
Electricity

## Target 1

**Consommation d'énergie du site**  
SiteEnergyUse(kBtu)

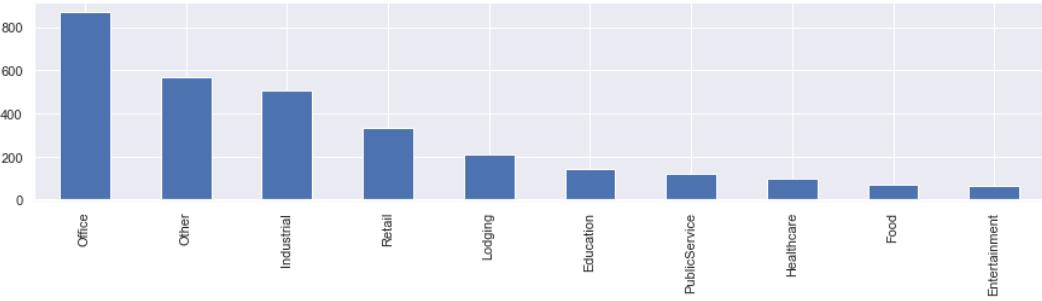
## Target 2

**Emissions de gaz à effet de serre**  
TotalGHGEmissions

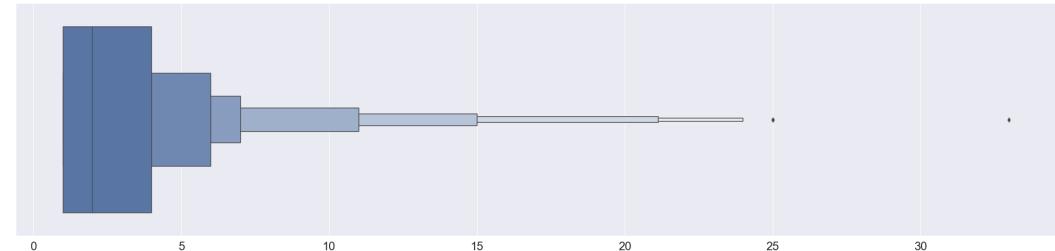
# Exploration

## building characteristics

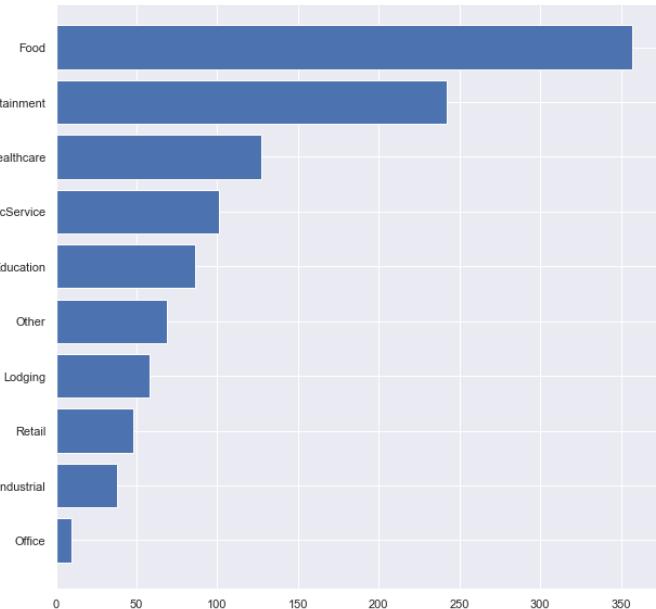
**UTILISATION TOUT CONFONDU**



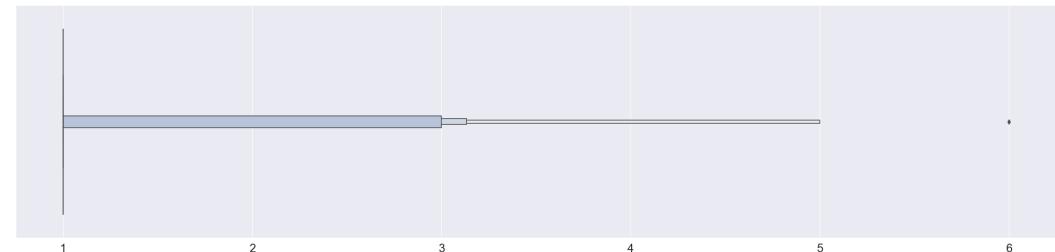
**NOMBRE D'ÉTAGES**



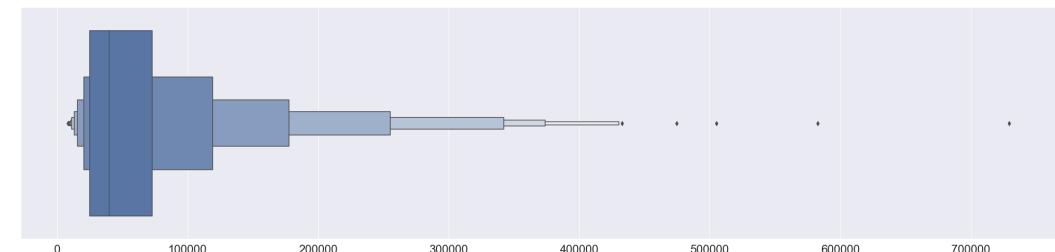
**UTILISATION PRINCIPALE**



**NOMBRE DE BÂTIMENTS**



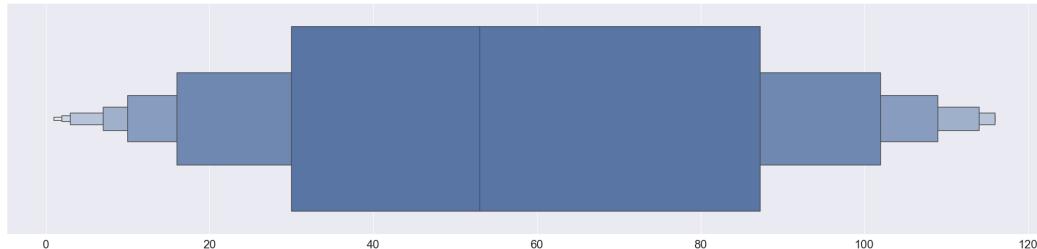
**SUPERFICIE DE L'ACTIVITÉ PRINCIPALE**



# Exploration

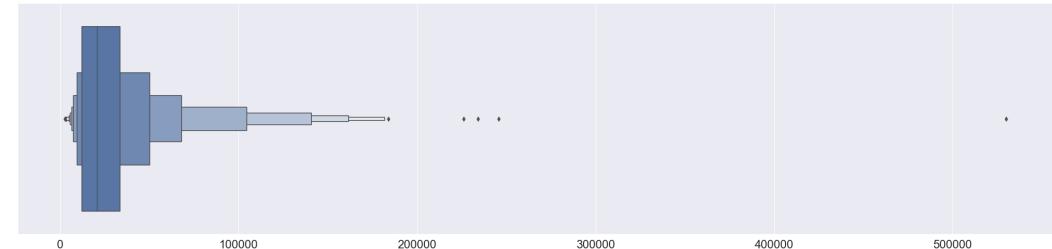
## feature engineering

### ÂGE DU BÂTIMENT



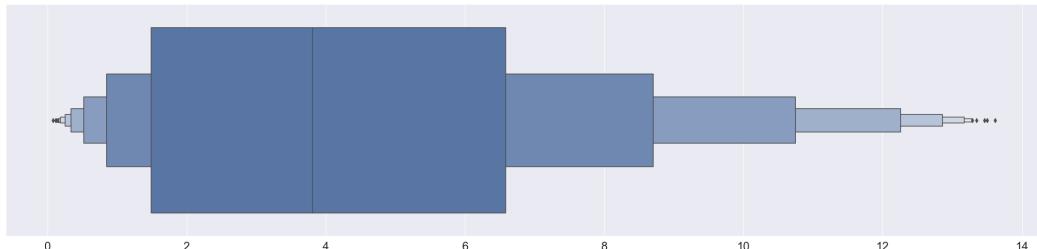
```
# Âge de la construction  
data['Age'] = 2016 - data['YearBuilt']
```

### SUPERFICIE MOYENNE



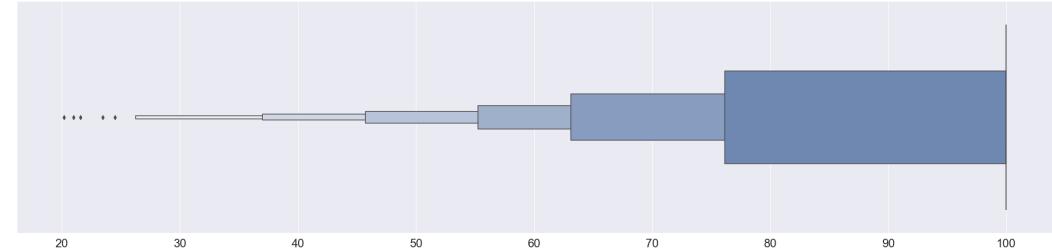
```
# Surface par étage par bâtiment  
data['mean_GFA'] = data['PropertyGFATotal']/(data['NumberofFloors']*data['NumberofBuildings']).round(2)
```

### DISTANCE DU CENTRE-VILLE (EN KM)



```
# Distance par rapport au centre  
downtown = (47.608056, -122.336111)  
data['CBD_distance_km'] = data.apply(lambda x: haversine((x['Latitude'], x['Longitude']), downtown), axis=1)
```

### PROPORTION DE BATIMENTS

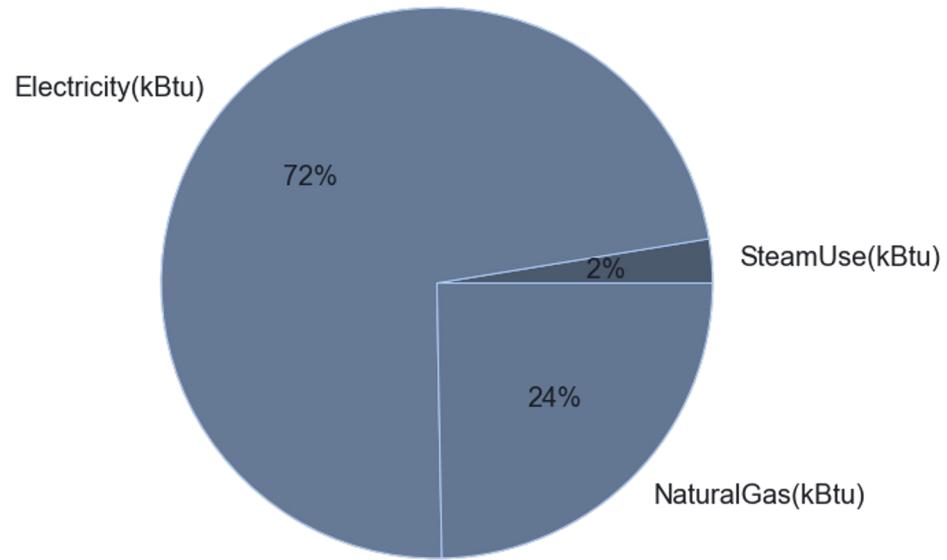


```
# Proportion de surface non-parking  
data['GFABuildingRate'] = (round((data['PropertyGFABuilding(s)'].fillna(0)  
/data['PropertyGFATotal'].fillna(0)),5))*100
```

# Exploration

energy

Répartition des différentes sources d'énergie



À savoir :

**kBtu**  
(kilo-British thermal unit)

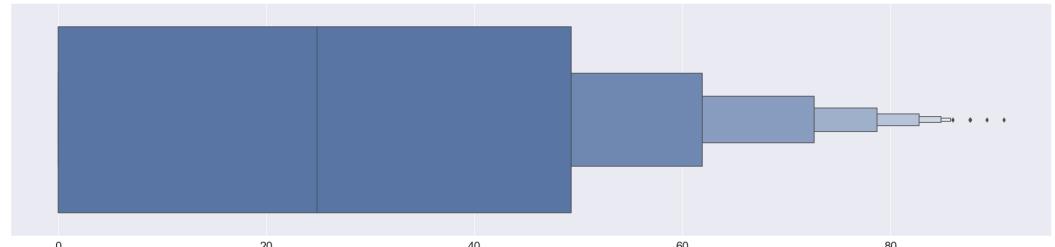
3.142

=

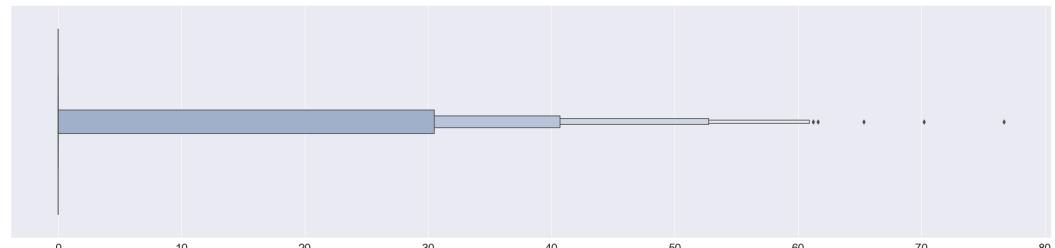
**kWh**  
(kilo-Watt-hour)

1

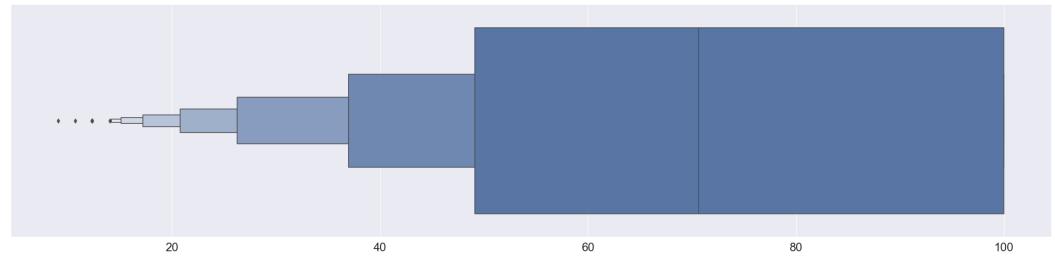
**GAZ NATUREL**



**VAPEUR**



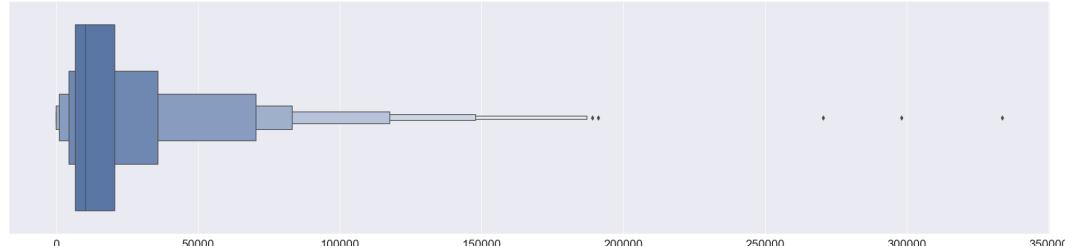
**ELECTRICITÉ**



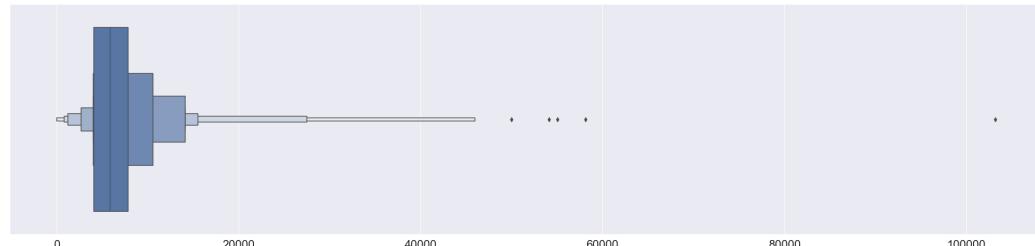
# Exploration

## KNN imputer

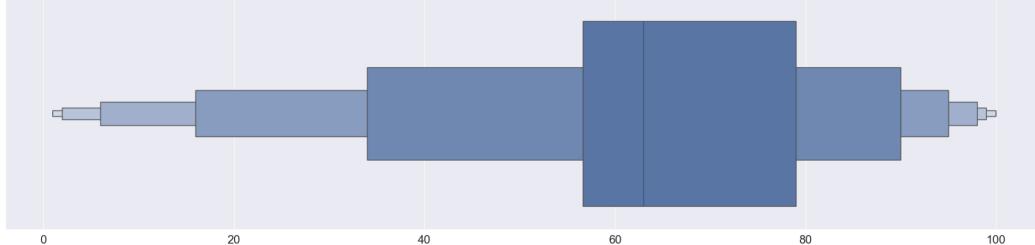
### SECOND LARGEST PROPERTY USE TYPE



### THIRD LARGEST PROPERTY USE TYPE GFA



### ENERGY STAR SCORE



`sklearn.neighbors.KNeighborsRegressor`

`sklearn.impute.KNNImputer`

GridSearchCV

`n_neighbors :`

`SLPUT GFA`

41

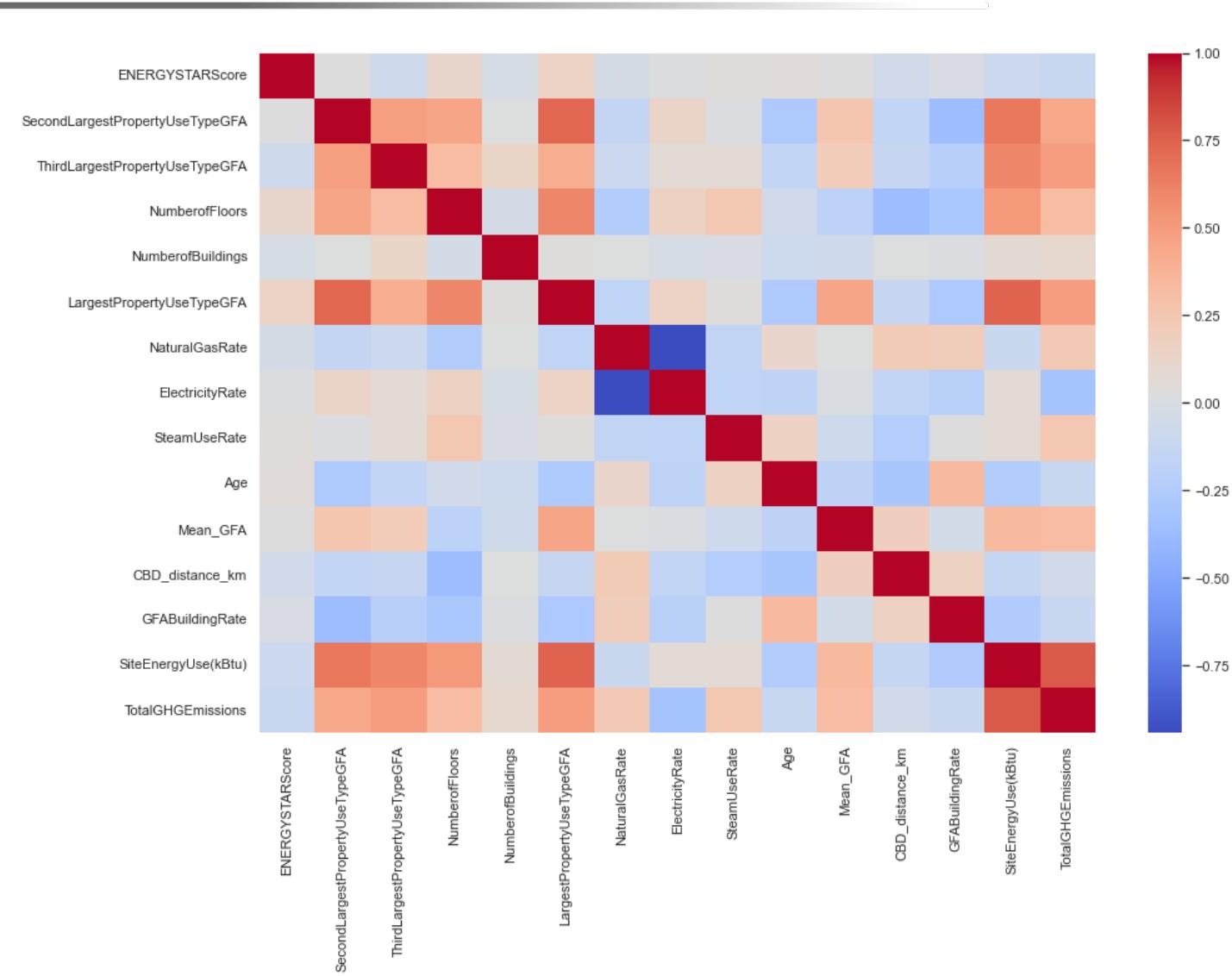
`TLPUT GFA`

54

`ESS`

138

# Exploration correlation





# PARTIE III

---

## Preprocessing

# Preprocessing

GetDummies

- 1 Encodage sur les features catégorielles**  
pour effectuer des opérations de Machine Learning
- 2 Plus simple**  
que OneHotEncoder
- 3 Effectué sur 4 variables catégorielles**  
Neighborhood, SecondLargestPropertyUseType,  
ThirdLargestPropertyUseType et LargestPropertyUseType

# Preprocessing

GetDummies

- 1 **Encodage sur les features catégorielles**  
pour effectuer des opérations de Machine Learning
- 2 **Plus simple**  
que OneHotEncoder
- 3 **Effectué sur 4 variables catégorielles**  
Neighborhood, SecondLargestPropertyUseType,  
ThirdLargestPropertyUseType et LargestPropertyUseType

# Preprocessing

GetDummies

- 1 **Encodage sur les features catégorielles**  
pour effectuer des opérations de Machine Learning
- 2 **Plus simple**  
que OneHotEncoder
- 3 **Effectué sur 4 variables catégorielles**  
Neighborhood, SecondLargestPropertyUseType,  
ThirdLargestPropertyUseType et LargestPropertyUseType

# Preprocessing

MinMaxScaler

## Raison 1

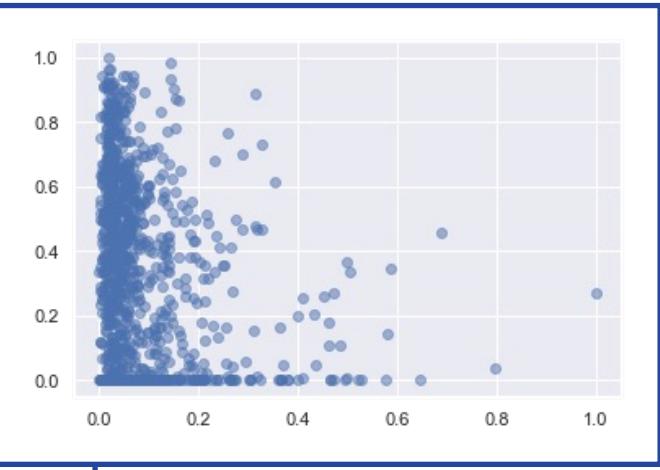
### Différentes échelles entre les variables

Une variable peut prendre le dessus lors de la minimisation de la fonction coût

## Raison 2

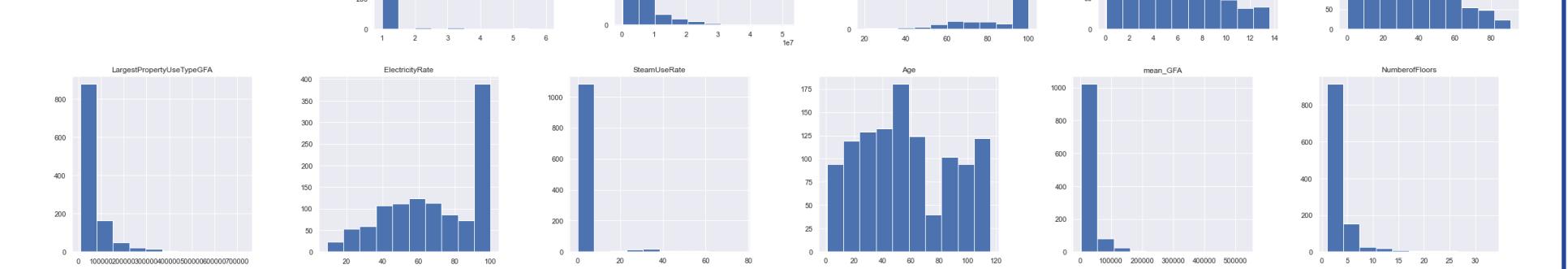
### Pas d'écrasement de données

Variance suffisante lors de la transformation (pas d'outliers)



## Raison 3

### Données non normales

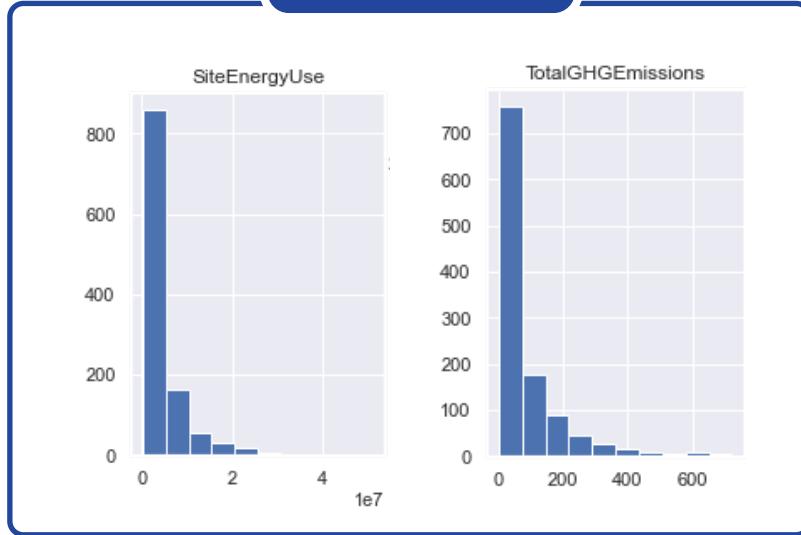


# Preprocessing

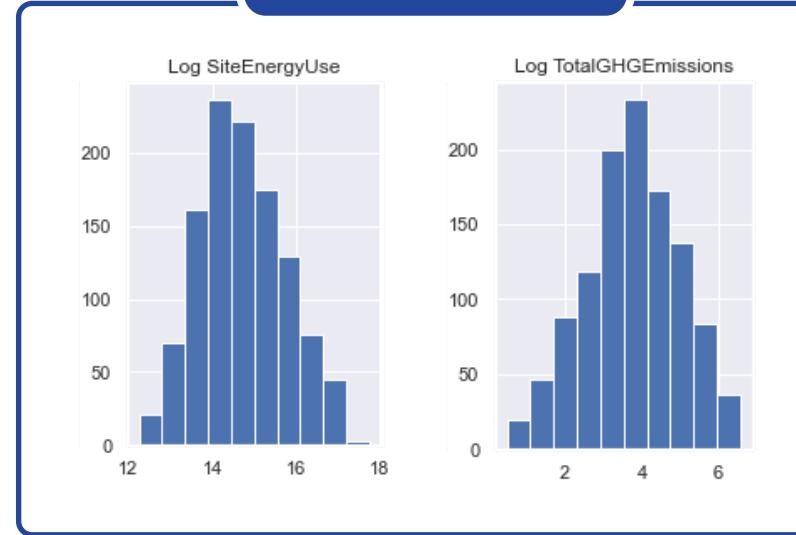
## log trick

```
FunctionTransformer(np.log,inverse_func=np.exp,check_inverse=True)
```

Original



Transformed



Corrélation

SiteEnergyUse(kBtu)	0.51	0.075	0.74	-0.095	-0.11	0.071	0.073	-0.25	0.34	-0.14	-0.25
---------------------	------	-------	------	--------	-------	-------	-------	-------	------	-------	-------

TotalGHGEmissions	0.32	0.1	0.49	-0.12	0.24	-0.32	0.25	-0.12	0.33	-0.062	-0.11
-------------------	------	-----	------	-------	------	-------	------	-------	------	--------	-------

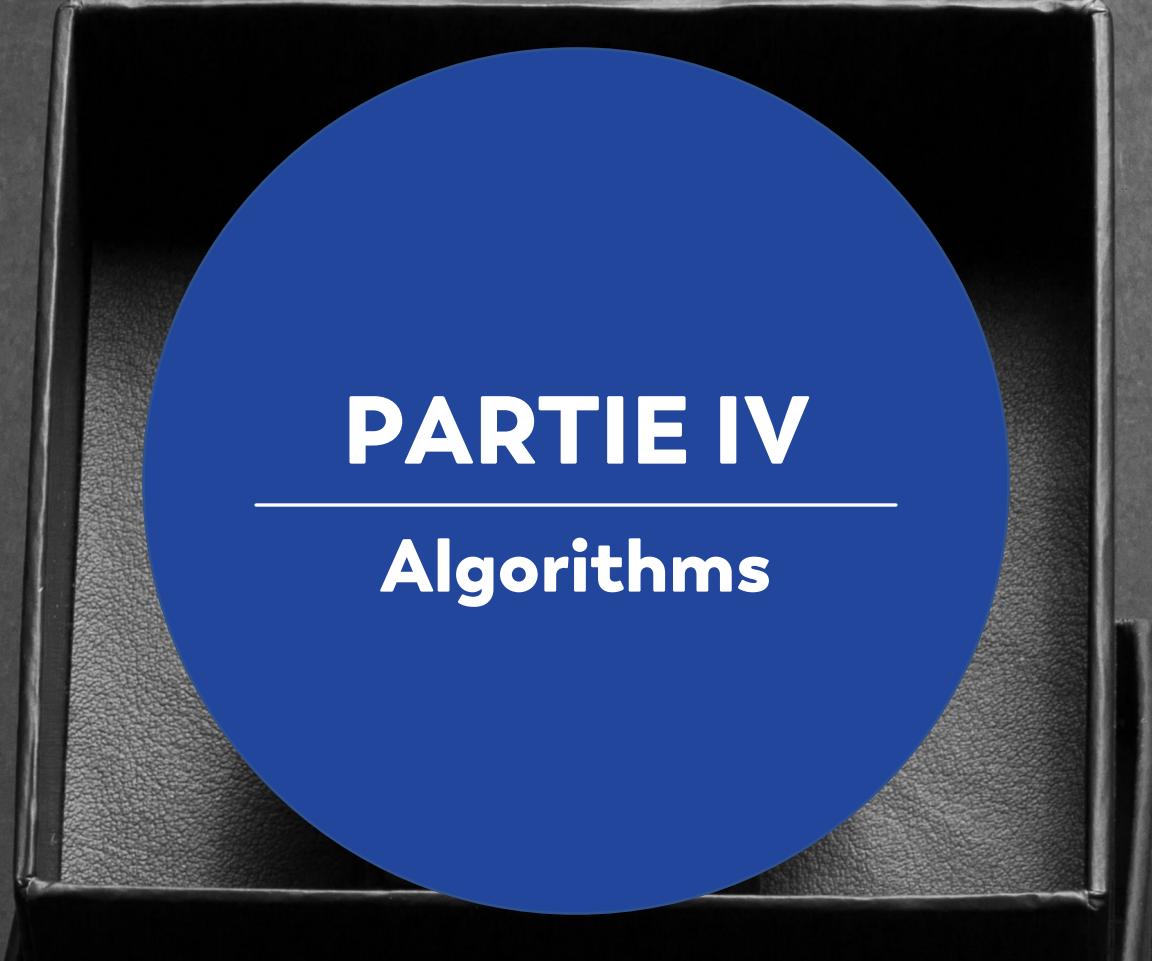
NumberofFloors  
NumberofBuildings  
LargestPropertyUseTypeGFA  
ENERGYSTARScore  
NaturalGasRate  
ElectricityRate  
SteamUseRate  
Age  
mean\_GFA  
CBD\_distance\_km  
GFABuildingRate

Corrélation ajustée

LogSiteEnergyUse	0.46	0.098	0.63	-0.18	-0.056	0.019	0.098	-0.26	0.3	-0.15	-0.3
------------------	------	-------	------	-------	--------	-------	-------	-------	-----	-------	------

LogTotalGHGEmissions	0.31	0.098	0.45	-0.15	0.44	-0.51	0.22	-0.12	0.26	-0.061	-0.13
----------------------	------	-------	------	-------	------	-------	------	-------	------	--------	-------

NumberofFloors  
NumberofBuildings  
LargestPropertyUseTypeGFA  
ENERGYSTARScore  
NaturalGasRate  
ElectricityRate  
SteamUseRate  
Age  
mean\_GFA  
CBD\_distance\_km  
GFABuildingRate



# PARTIE IV

---

## Algorithms

# Algorithms

naive approaches

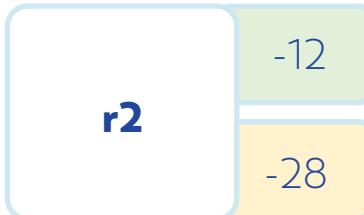
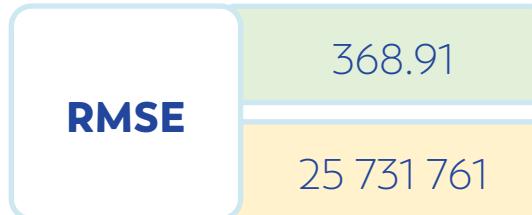
TARGETS

GHGEmmisions

SiteEnergyUse

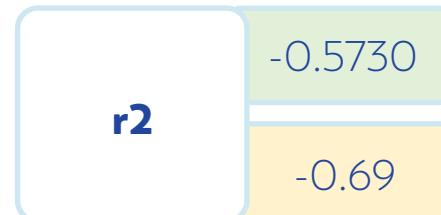
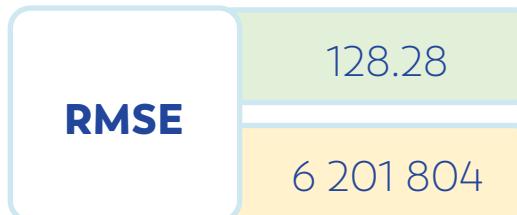
Tenter de prédire avec des valeurs aléatoires :

```
np.random.randint(np.min(y_log), np.max(y_log), y_test_rand.shape)
```



Prédit systématiquement la moyenne du train set:

```
sklearn.dummy.DummyClassifier
```



# Algorithms

hyperparams tuning

TARGETS

GHGEmissions

SiteEnergyUse

`sklearn.linear_model.LinearRegression`

No hyperparams

`sklearn.linear_model.Ridge`

GridSearchCV

Alpha :

0.15

1.0

# Algorithms

hyperparams tuning

TARGETS

GHGEmmissions

SiteEnergyUse

`sklearn.linear_model.Lasso`

GridSearchCV

Alpha :

0.001363

0.00534

`sklearn.linear_model.ElasticNet`

GridSearchCV

Alpha :

0.00129

0.00207

l1 ratio :

1

0.3591

# Algorithms

hyperparams tuning

TARGETS

GHGEmissions

SiteEnergyUse

`sklearn.neighbors.KNeighborsRegressor`

GridSearchCV

**n\_neighbors :**

6      7

**metrics :**

manhattan      manhattan

`sklearn.ensemble.GradientBoostingRegressor`

GridSearchCV

**n\_estimators :**

93      81

# Algorithms

hyperparams tuning

TARGETS

GHGEmissions

SiteEnergyUse

`sklearn.ensemble.RandomForestRegressor`

`RandomizedSearchCV`

`n_estimators` :

99

14

`min_samples_split` :

6

5

`min_samples_leaf` :

4

2

`max_samples` :

0.11

0.48

`max_depth` :

28

29

# Algorithms

hyperparams tuning

TARGETS

GHGEmissions

SiteEnergyUse

*dmlc*  
**XGBoost**

RandomizedSearchCV

**base\_score :**

0.3105

0.2789

**learning\_rate :**

0.0578

0.0763

**max\_depth :**

6

6

**min\_child\_weight :**

5

2

**n\_estimators :**

145

91

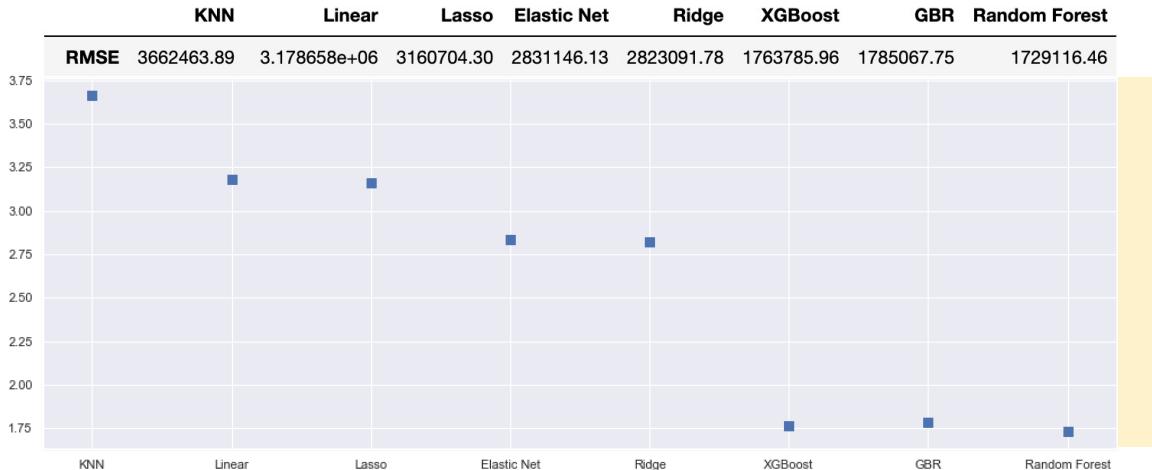
# Algorithms

RMSE score

TARGETS

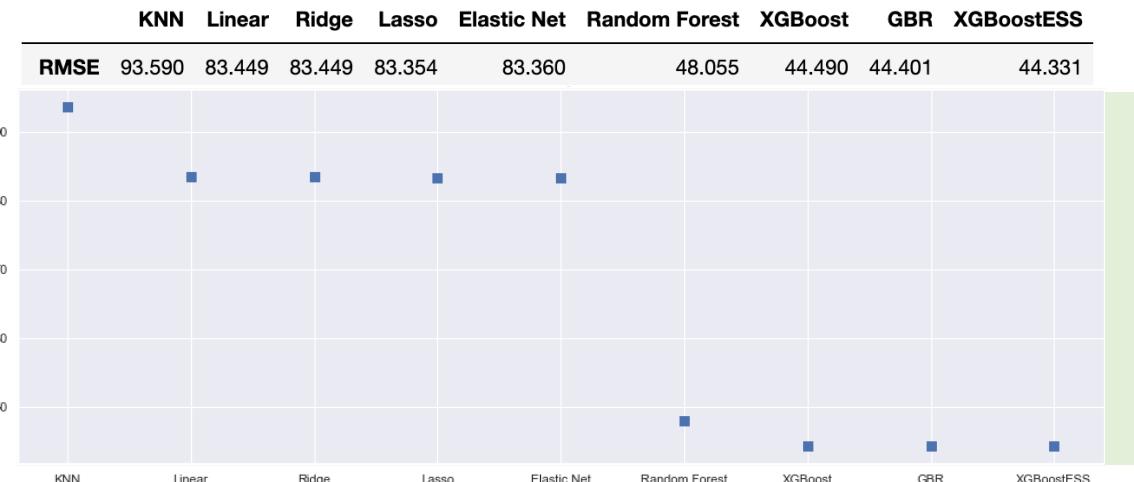
GHGEmmisions

SiteEnergyUse



**Prediction**  
min : 606 895  
moy : 3 957 993  
max : 32 703 654

**Prediction**  
min : -66.15  
moy : 81.18  
max : 541.73



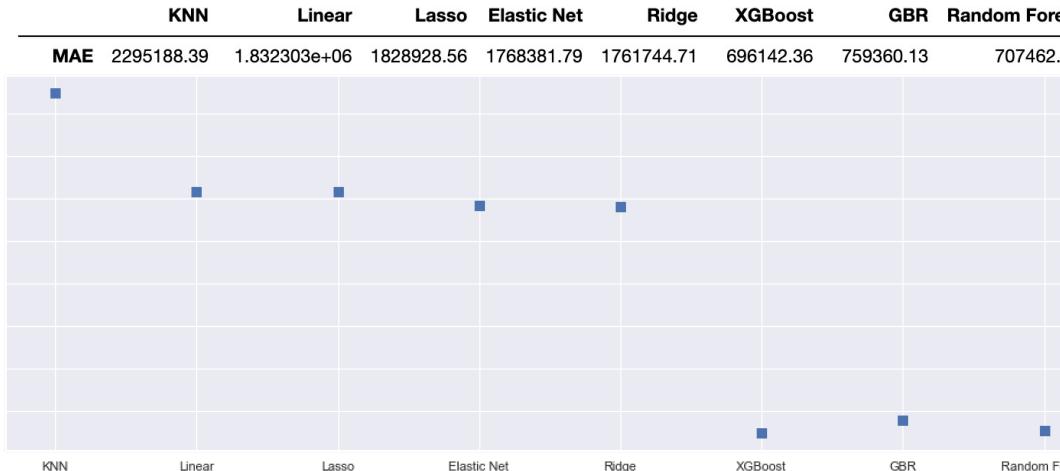
# Algorithms

MAE score

TARGETS

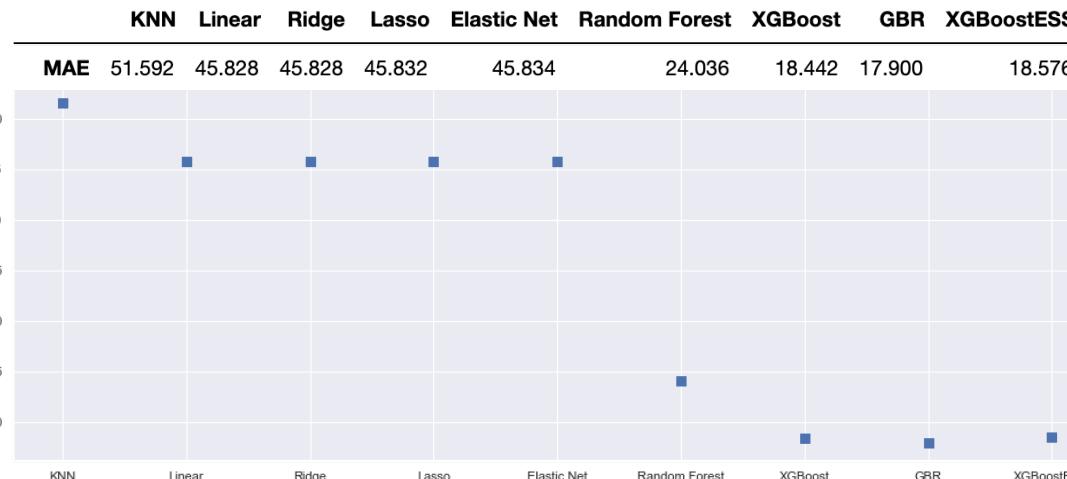
GHGEmmisions

SiteEnergyUse



**Prediction**  
min : 606 895  
moy : 3 957 993  
max : 32 703 654

**Prediction**  
min : -66.15  
moy : 81.18  
max : 541.73



# Algorithms

MdAE score

TARGETS

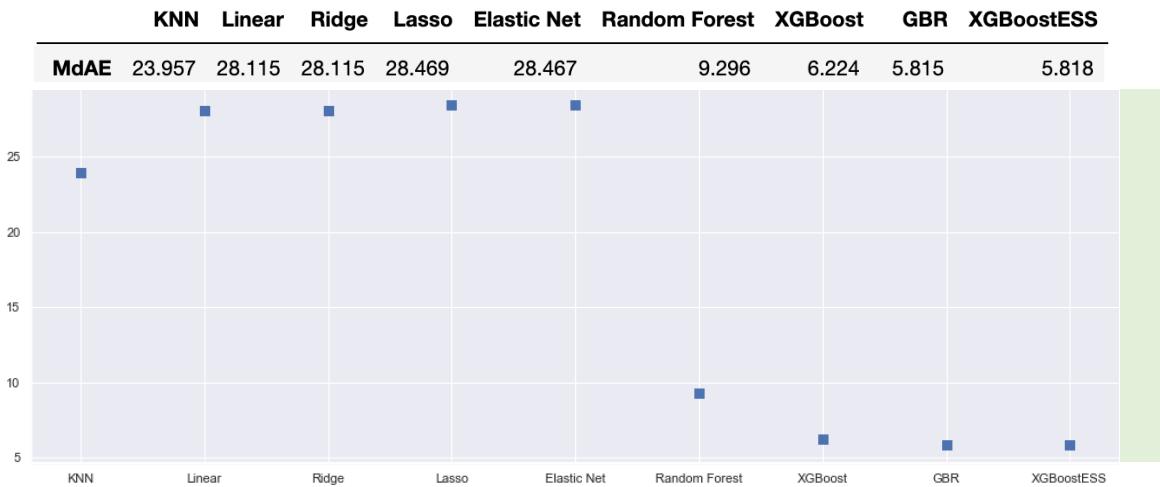
GHGEmmisions

SiteEnergyUse



**Prediction**  
min : 606 895  
moy : 3 957 993  
max : 32 703 654

**Prediction**  
min : -66.15  
moy : 81.18  
max : 541.73



# Algorithms r2 score

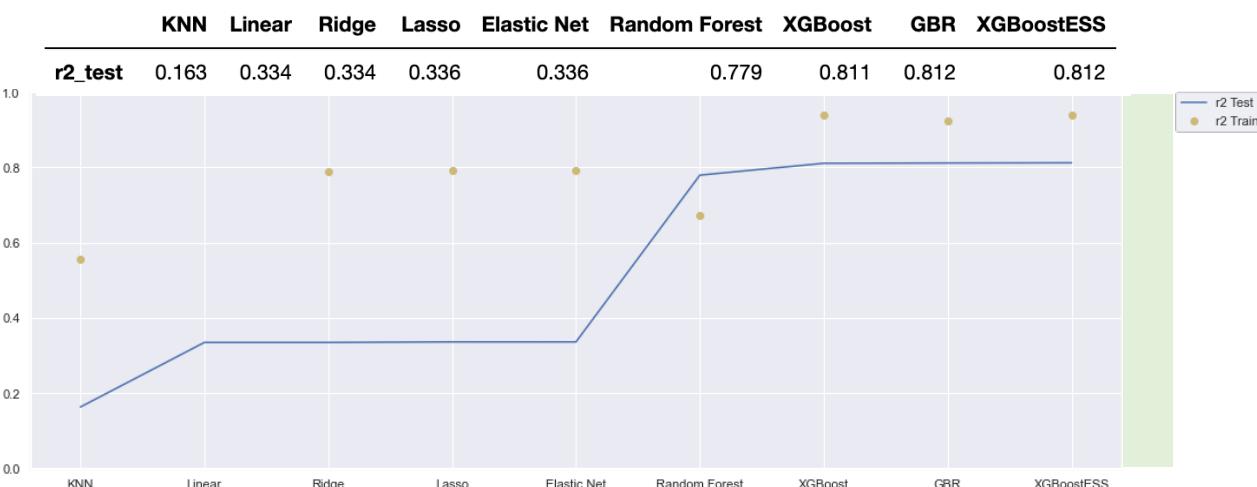
TARGETS

GHGEmissions

SiteEnergyUse



Avec l'ESS  
r2 : 0.812  
r2 ajusté : 0.803





# PARTIE V

---

## Interpretation

# Interpretation SHAP (Force Plot)

## EMISSIONS DE GAZ À EFFET DE SERRE

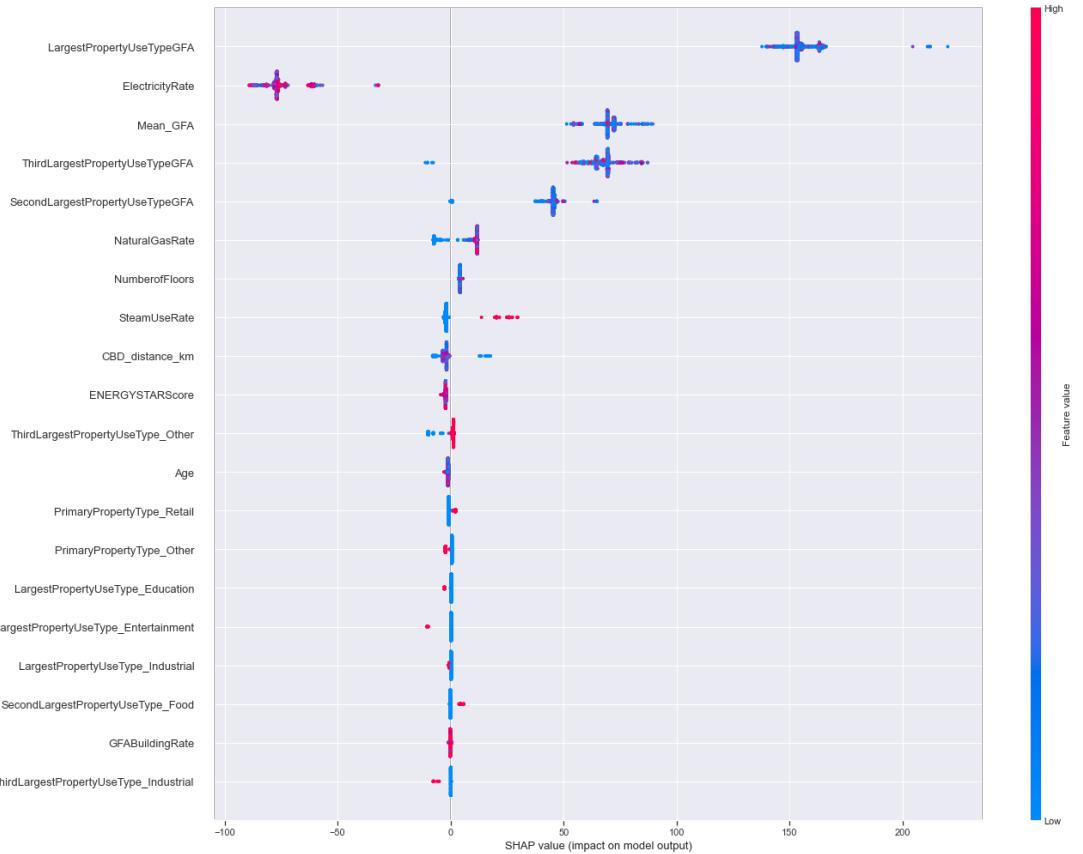


## CONSOMMATION D'ENERGIE

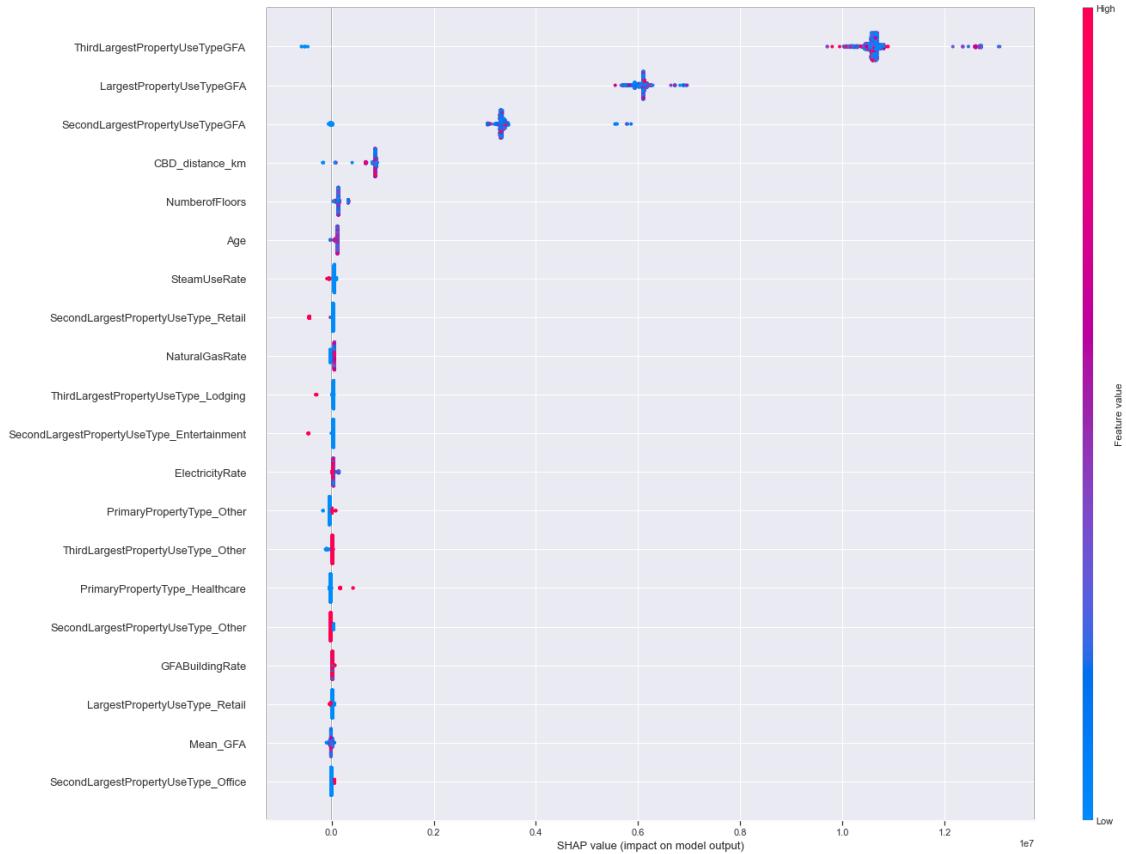


# Interpretation SHAP (Summary Plot)

## EMISSIONS DE GAZ À EFFET DE SERRE



## CONSOMMATION D'ENERGIE



**Merci**

**Alexandre Delaguillaumie**