# Projet 3

Anticiper les besoins Energétiques de Seattle

## Problématique

Prédire consommation et émission







2050





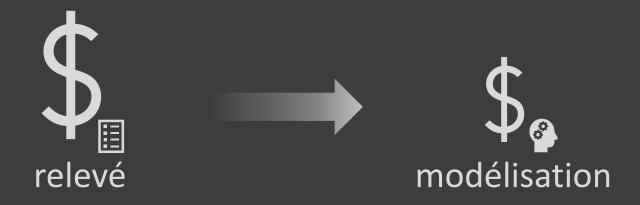






# Problématique

Réduire les coûts



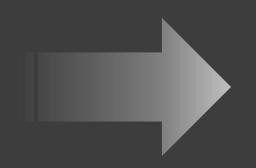
### Problématique

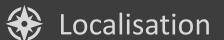
Les données

#### **Dataframe**

47 Variables

Immeubles 3340 + 3376 = **6716** 

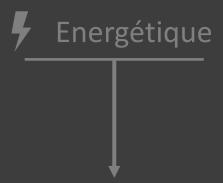


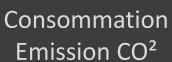


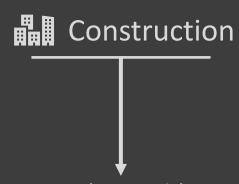




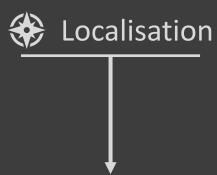
Sélection des Features







- Caractéristique bâtiment
- Type d'usage
- Commentaires
- Conformité données



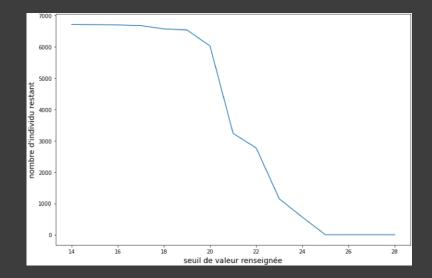
- District
- Quartier
- Lon/Lat
- Zip Code
- Adresse

Labels

Features

Filtrage des individus

filtrage des individus moins renseigné



Retire 172 individus

Suppression des doublons

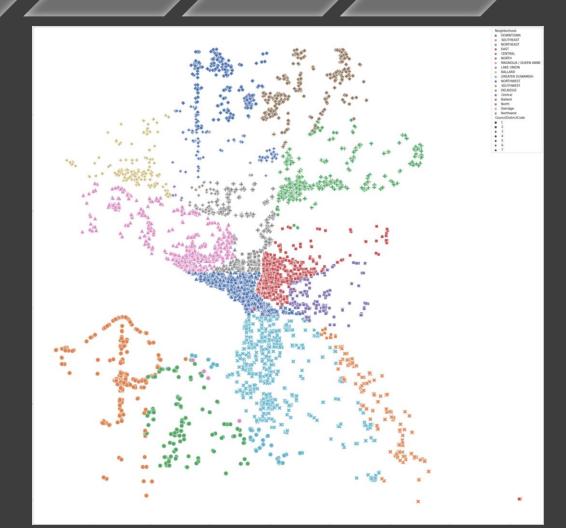
Retire 127 individus

Variable de localisation

Ses districts et quartiers

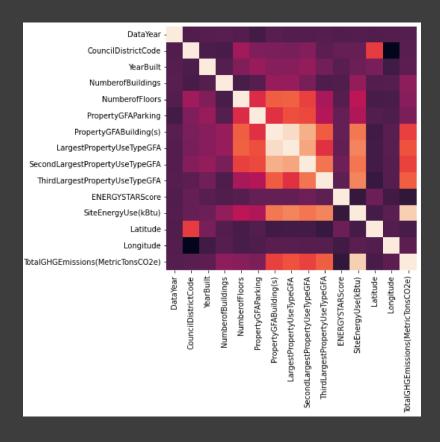
ZipCode: 83 différents

Adresse: non exploitable



Corrélations

#### Corrélations



Valeurs manquantes

Peu de valeurs manquantes distribution concentrée distribution asymétrique

Combler par la médiane

Nouvelles variables

### Avant

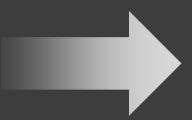
LargestProperty UseType	LargestProperty UseTypeGFA
Hotel	88434
Office	51029
Hotel	348329
Other	385274
Office	103501

### Après

Hotel	Office	Other
88434	0	0
0	51029	0
348329	0	0
0	0	385274
0	103501	0

Nouvelles variables

- Liste de tout les types d'usage
- Type de bâtiment
- Quartier
- District

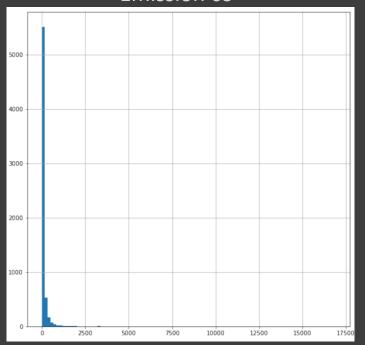


One Hot Encoding

Nouvelles variables

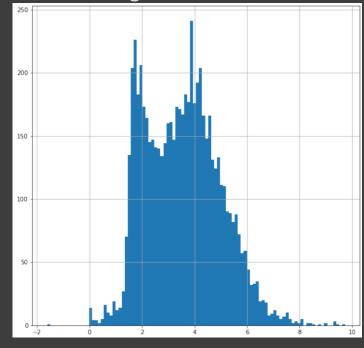
#### Avant

Consommation Emission co<sup>2</sup>



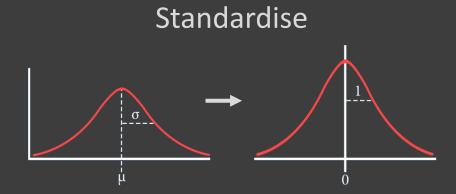
### Après

Log Consommation Log Emission co<sup>2</sup>



Préparation des données





Modèles Linéaire

#### Baseline

Régression linéaire

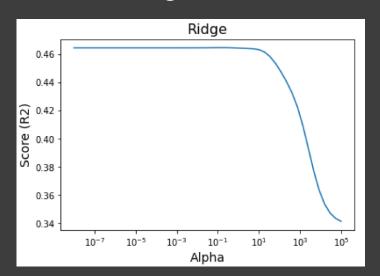
Score (R<sup>2</sup>): 0.481

Learning time: xxx

### Régression Ridge

Score (R<sup>2</sup>): 0.481

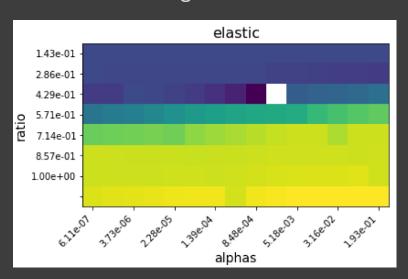
Learning time: 0.02



#### Régression Elastic Net

Score (R<sup>2</sup>): 0.481

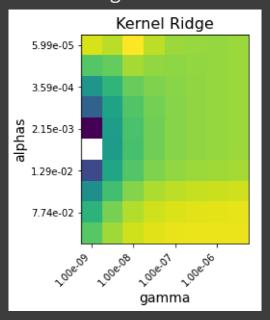
Learning time: 0.54



Modèles à Noyau

#### Kernel Ridge

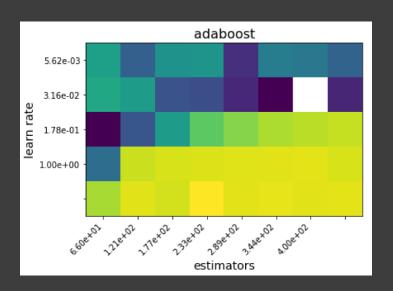
Score (R<sup>2</sup>): 0.447 Learning time: 0.89



Modèles Ensembliste

#### Ada Boost

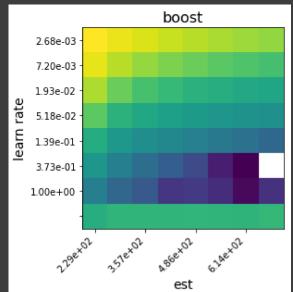
Score (R<sup>2</sup>): 0.489 Learning time: 7.93



# Gradiant Boosting

Score (R<sup>2</sup>): 0.754

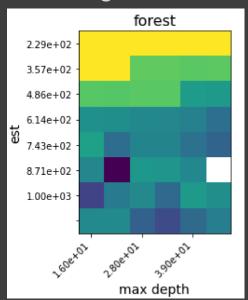
Learning time: 16.2



#### Random Forest

Score (R<sup>2</sup>): 0.801

Learning time: 26.8

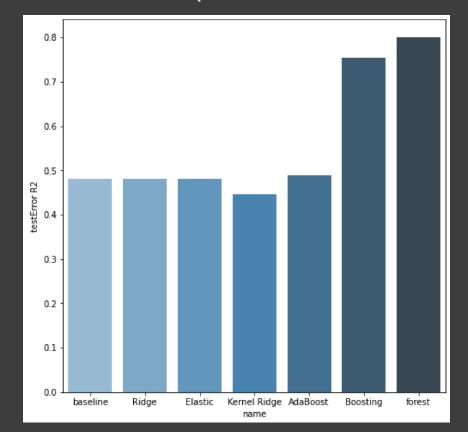


Résumé

#### Résumé des modèles

name	trainError	testError R2	testError RMSE	learningTime	predictTime
baseline	0.524629	0.481074	1.026192	NaN	0.00509963
Ridge	0.524629	0.481074	1.026192	0.017230	0.00682005
Elastic	0.524629	0.481087	1.026180	0.542988	0.00499864
Kernel Ridge	0.623419	0.446675	1.059660	2.043311	0.890072
AdaBoost	0.491015	0.489334	1.017992	7.931864	0.517571
Boosting	0.902375	0.753696	0.706987	16.209746	0.0691218
forest	0.974181	0.800998	0.635484	26.763308	0.633284

#### Score par modèles



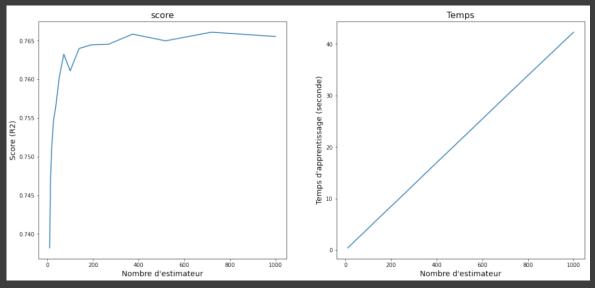
Affinage des hyperparamètres

Sélection Nombre d'estimateur

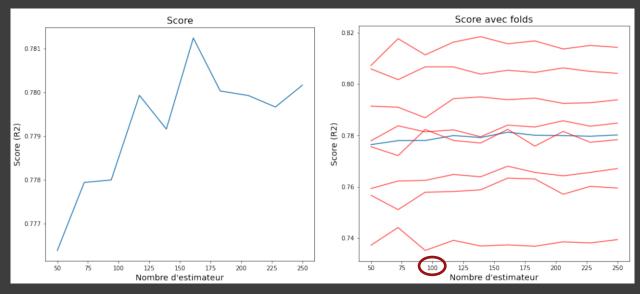
Intervalle large et peu précis

Intervalle étroit et plus précis

#### 5 Folds

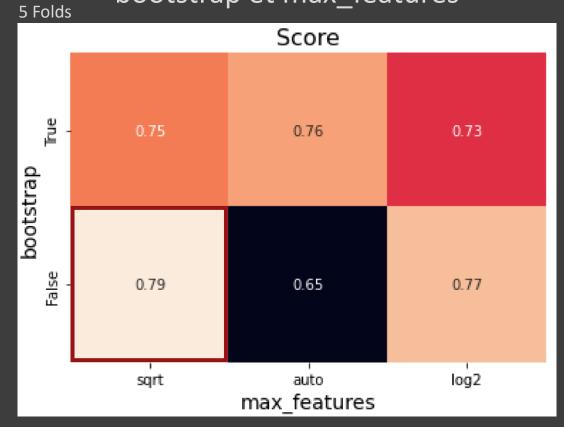


#### 8 Folds



Affinage des hyperparamètres

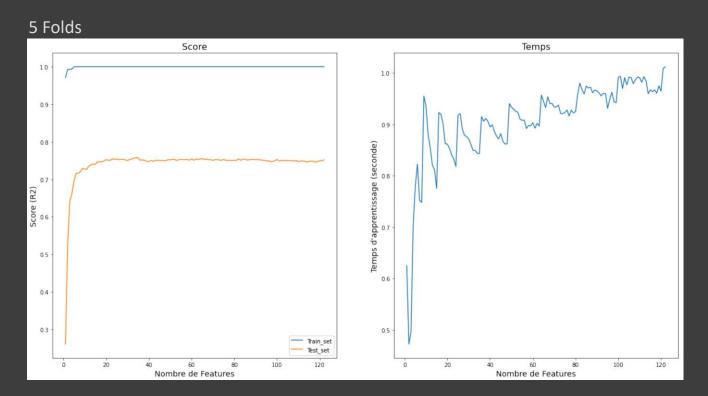
### Sélection bootstrap et max\_features



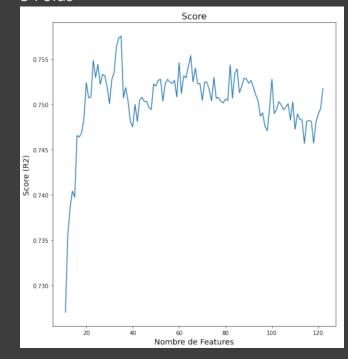
Feature Selection

Intervalle large et peu précis

Intervalle de 10 à plus

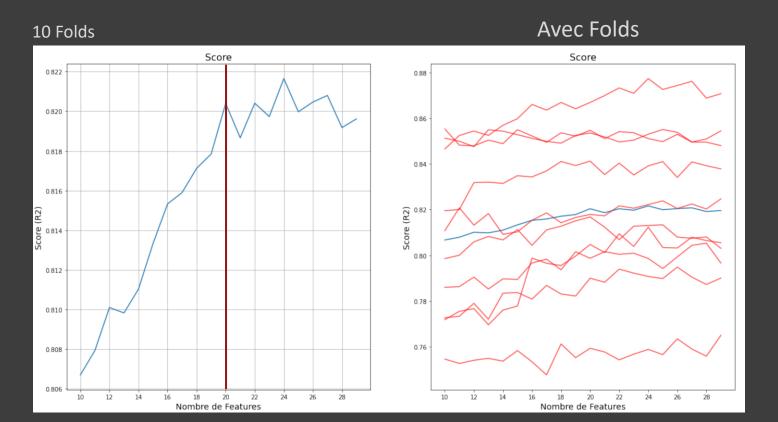








Intervalle étroit et précis



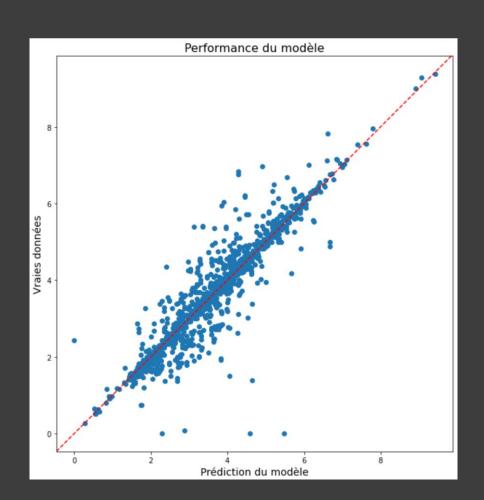
#### Feature Selection

#### Informations retenues

- 1. Surface du bâtiment
- 2. Année de fabrication
- 3. Surface par type d'usage pour :
  - Multifamily Housing
  - Parking
  - Office
  - Non-Refrigerated Warehouse
  - Other
  - Hotel
  - Retail Store
  - Supermarket/Grocery Store
  - Hospital (General Medical & Surgical)

- 4. Nombre d'étage
- 5. Latitude
- 6. Longitude
- 7. Type d'usage pour :
  - Multifamily LR (1-4)
  - Multifamily Housing
  - NonResidential
  - Office
- 8. Surface du parking
- 9. S'il est du centre ville (DOWNTOWN)

EnergyStar Score

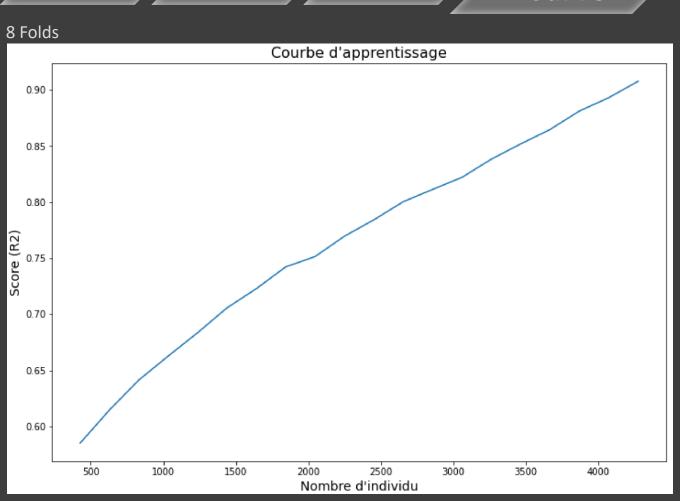


Score R<sup>2</sup> avec

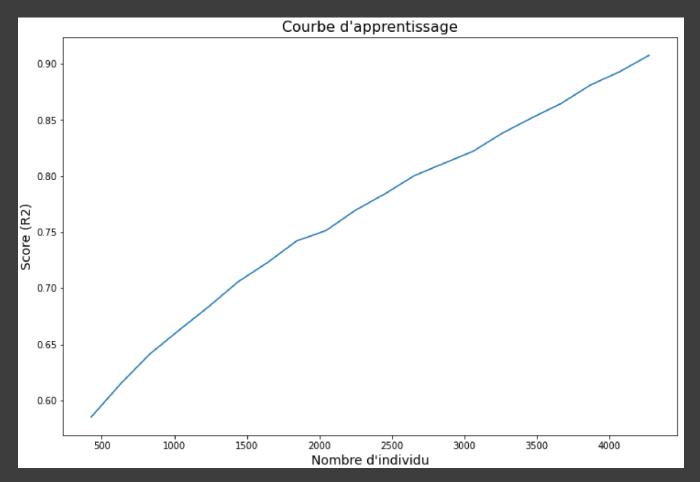
Score R<sup>2</sup> sans

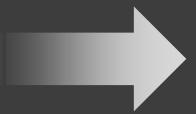
EnergyStarScore: 0.904 EnergyStarScore: 0.903











Plus de données!

