Ville de Seattle -Consommation et émissions des bâtiments

Prédiction des émissions de CO2 et de consommation totale d'énergie

Préambule

- Ville de Seattle
 - Objectif neutralité carbone 2050
 - Souhaite comprendre les consommations et émissions de CO2 de ses bâtiments
 - A fait des relevés dans différents bâtiments de la ville
- Objectif
 - Prédire les consommations et émissions de CO2
 - En utilisant les données structurelles
 - Évaluer l'intérêt de l'ENERGYSTARScore dans la modélisation

Sommaire

- Présentation du jeu de données
- Feature engineering
 - Valeurs manquantes
 - Variables quantitatives
 - Variables qualitatives
- Modélisation
 - Méthodologie
 - Résultats
- Feature importance
- ENERGYSTARScore

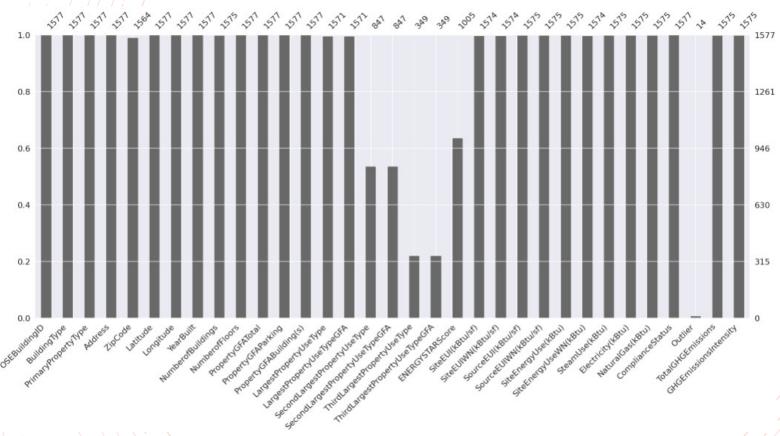
Présentation des données

- Premier passage en revue du jeu de données
 - 3380 bâtiments environ
 - 45 caractéristiques
 - Localisation
 - Nombre de bâtiments et étages
 - Surfaces et utilisations
 - Consommations et émissions de CO2
 - Peu de valeurs manquantes
 - Pas de doublons

Exploration des données

- Type de bâtiments
 - 1665 bâtiments non destinés à l'habitation
- Simplification des features
 - DataYear
 - City
 - State
 - Neighborhood
 - CouncilDistrictCode
 - Comments
 - PropertyName
 - TaxParcelIdentificationNumber
 - Différentes unités
 - YearsENERGYSTARCertified

Valeurs manquantes



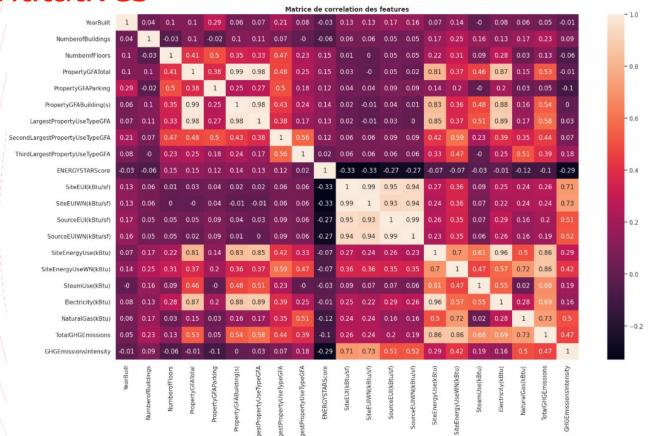
Valeurs manquantes



- Renseignement des valeurs lorsque possible
 - Google Street View
 - Recherche adresse
 - ..
- Données énergétiques
 - Suppression des bâtiments
- Surfaces et utilisations
 - Largest => Suppression
 - Second & Third =>
 - Surface = 0
 - Type = NotUsed

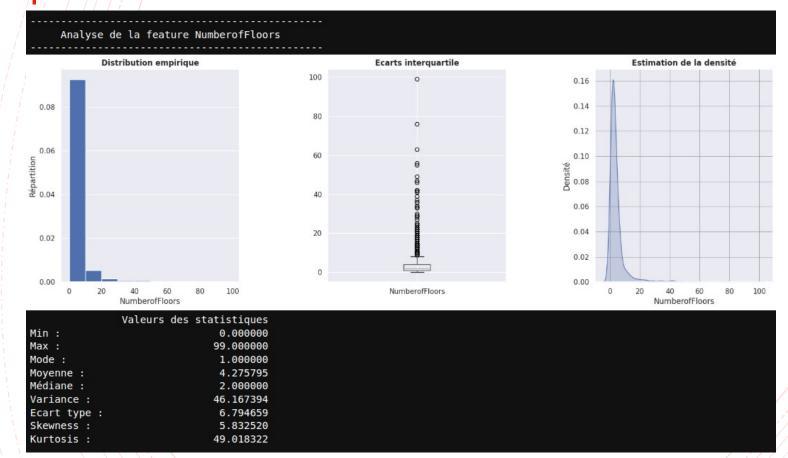
Variables quantitatives

- Analyse des corrélations entre features
 - Suppression des variables peu informatives



Variables quantitatives

- Analyse statistiques
 - Correction des valeurs aberrantes

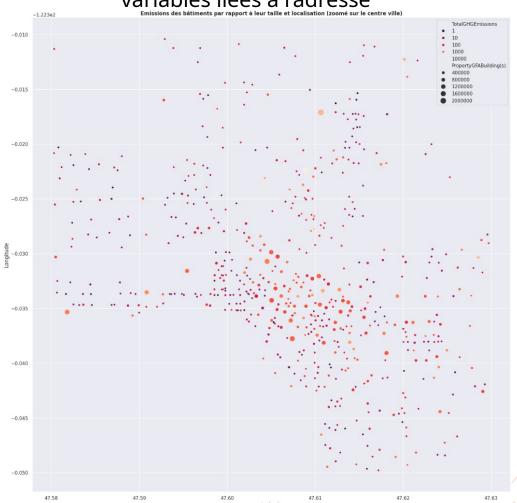


Variables quantitatives

- Création de nouvelles features
 - Répartition des énergies en pourcentage de la consommation totale
 - Pourcentage d'utilisation du bâtiment pour :
 - Bâtiment hors parking
 - Parking
 - Plus grande utilisation
 - Deuxième plus grande utilisation
 - Troisième plus grande utilisation

Variables qualitatives -Localisation

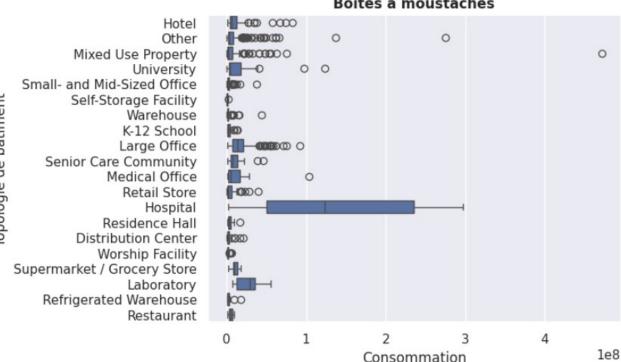
 Carte des bâtiments pour apprécier les variables liées à l'adresse



Variables qualitatives – Type de bâtiment

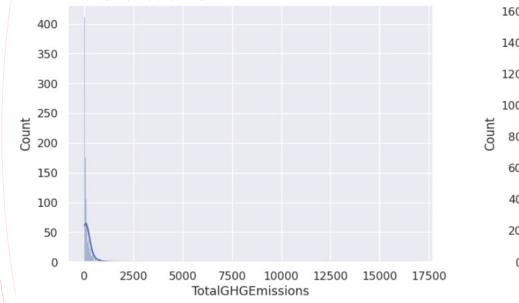
- feature
- Différents types de traitement
 - OneHotEncoder si peu de catégories
 - TargetEncoder sinon (pipeline)

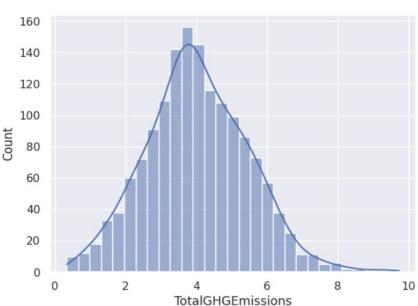




Normalisation

- Passage au log
 - Plusieurs variables et la cible ont des distributions très étalées





Distribution avant et après passage au log

Normalisation

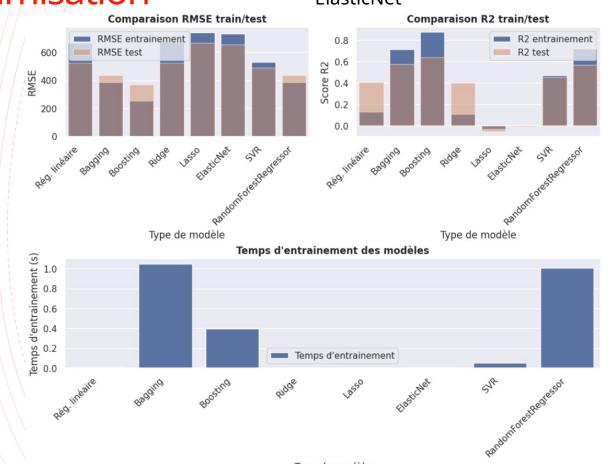
- Préprocesseur
 - Normalisation selon le type de données
 - Quantitative = StandardScaler
 - Qualitative = TargetEncoder
 - Booléenne = MinMaxScaler
 - Sera appliqué dans un pipeline

Métriques d'évaluation

- RMSE
 - Grandeur de l'erreur dans l'unité de la cible
 - Les valeurs prédites seront passées à l'exponentiel avant calcul du score
- Coefficient de détermination R2
 - Performance du modèle
 - Les valeurs prédites seront passées à l'exponentiel avant calcul du score
- Temps d'entrainement

Scores avant optimisation

- Sélection algorithmes
 - SVR
 - ElasticNet

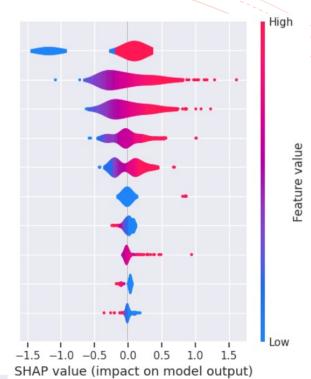


Scores après optimisation

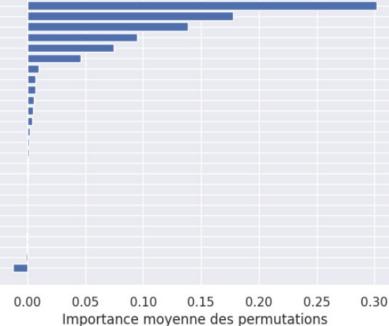
- Scores GridSearchCV :
 - RMSE moyenne entraînement : 327,27
 - R2 moyen entraînement : 0,797
- Scores set de test
 - RMSE: 322,26
 - R2:0,706

Feature importance

ThirdLargestPropertyUseTypeGFAPct
PropertyGFAParking
PropertyGFABuilding(s)
NaturalGas
Electricity
LargestPropertyUseTypeGFAPct
NumberofBuildings
SteamUse
LargestPropertyUseTypeGFA
ThirdLargestPropertyUseType







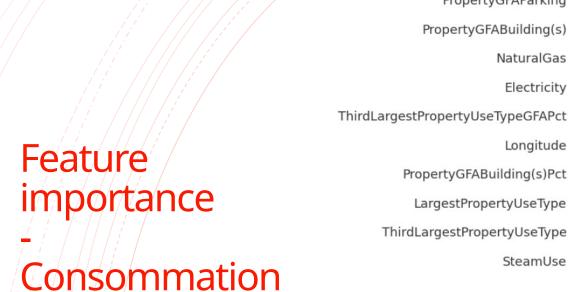
Scores avec ENERGYSTAR Score

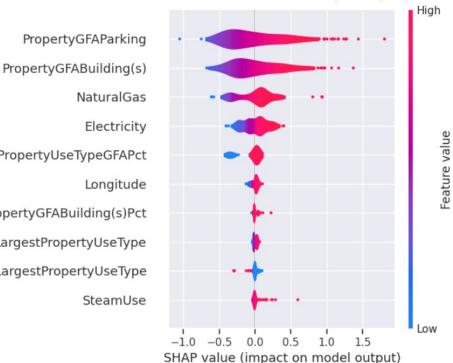
- Scores GridSearchCV:
 - RMSE moyenne entraînement : 340,95
 - R2 moyen entraînement : 0,842
- Scores set de test
 - RMSE: 286,17
 - R2:0,743

Scores après optimisation

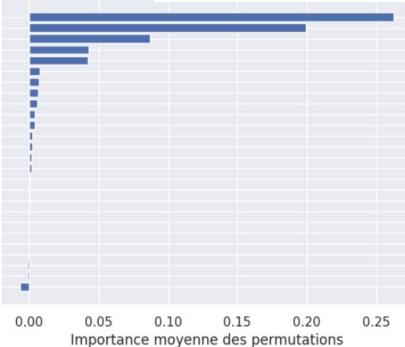
Consommation

- Scores GridSearchCV:
 - RMSE moyenne entraînement : 9,8M
 - R2 moyen entraînement : 0,821
- Scores set de test
 - RMSE: 8,4M
 - R2:0,817









Scores avec ENERGYSTAR Score

Consommation

- Scores GridSearchCV:
 - RMSE moyenne entraînement : 10,4M
 - R2 moyen entraînement : 0,839
- Scores set de test
 - RMSE: 10,7M
 - R2:0,862