

# Anticipation des besoins en consommation des bâtiments de Seattle



**Seattle**

## Contexte :

La ville de Seattle souhaite devenir une ville neutre en émission de carbone en 2050.



## Mission :

- Analyser la consommation en énergie et les émissions de G.E.S. des **bâtiments non destinés à l'habitation**.
- Tenter de prédire ces deux valeurs à partir des propriétés géographiques, architecturales et d'usage des bâtiments.
- Tester la pertinence de l'indicateur ENERGYSTARScore, difficile et coûteux à acquérir.

## Données :

Relevé des différents indicateurs pour les bâtiments de Seattle pour l'année 2016.

**I) Présentation et préparation du jeu de données**

**II) Méthodes de prédiction**

**III) Prédiction des indicateurs de consommation et d'émission**

Résultats 1 : Prédictions de la consommation énergétique

Résultats 2 : Prédictions des émissions de G.E.S.



# Partie I : Présentation et préparation des données

3

# Présentation des données et des variables d'intérêts

4

- Données fournies le 15 mars 2018 par Seattle sur des relevés de 2016.
- 3,376 propriétés
- 46 variables décrivant les propriétés (géographiques, architecturales, usage, consommations, émissions)

GEOGRAPHIQUES	ARCHITECTURAUX	USAGE	EMISSIONS/CONSO
<ul style="list-style-type: none"><li>• Longitude</li><li>• Latitude</li><li>• Code de District (District Code)</li><li>• Quartier (Neighborhood)</li></ul>	<ul style="list-style-type: none"><li>• Nombre de bâtiments</li><li>• Nombre d'étages</li><li>• Année de construction</li><li>• Surface totale (propertyGFATotal)</li></ul>	<ul style="list-style-type: none"><li>• Type de propriété principale</li></ul>	<ul style="list-style-type: none"><li>• TotalGHGEmmissions</li><li>• SiteEnergyUse</li></ul>

## Modification et création de variable :

Type de propriété principale : 'Office', 'Small- and Mid-Sized Office', 'Large Office' = Office

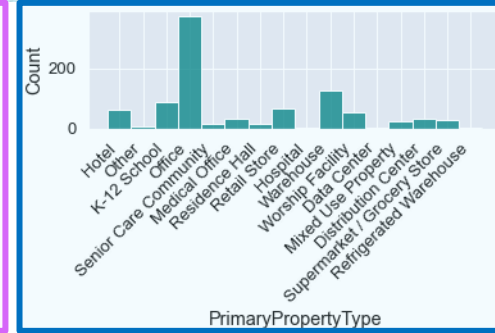
Parking\_per : Surface parking / Surface totale

Largest\_UseType\_per = Surface d'utilisation principale / Surface totale

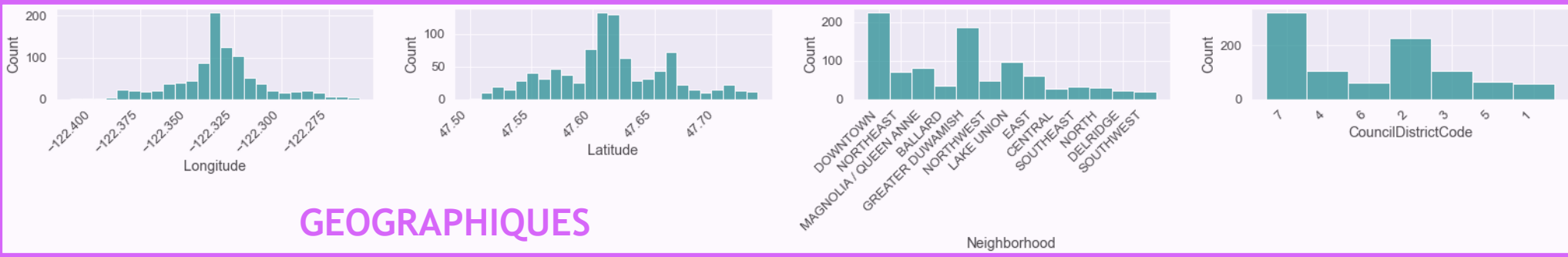
# Répartition empirique et modification des données

5

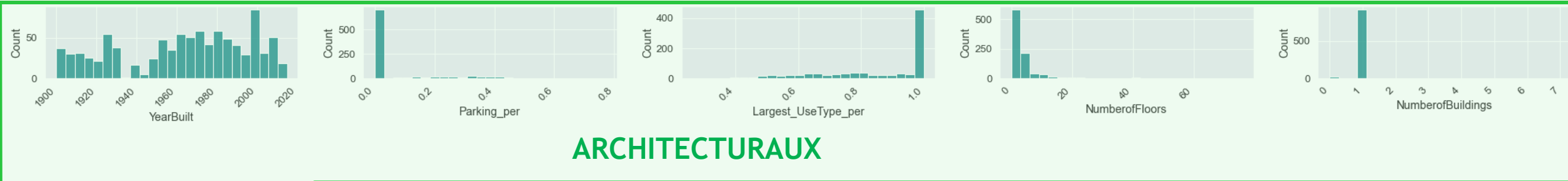
## USAGE



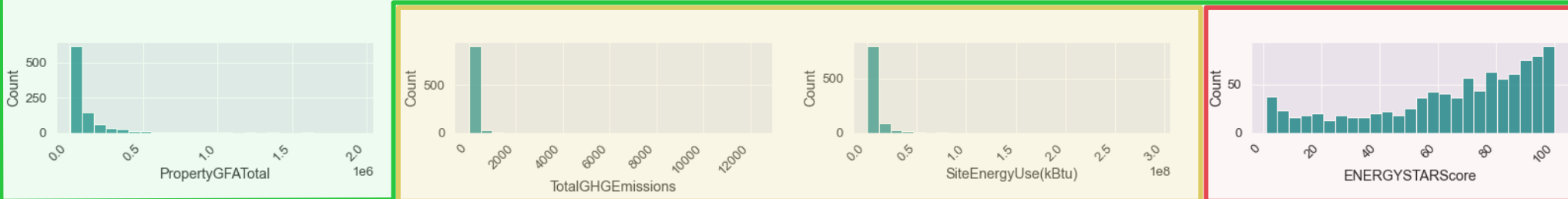
## GEOGRAPHIQUES



## ARCHITECTURAUX



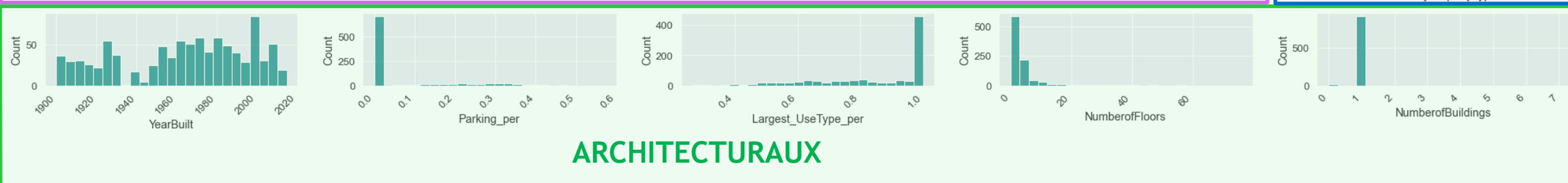
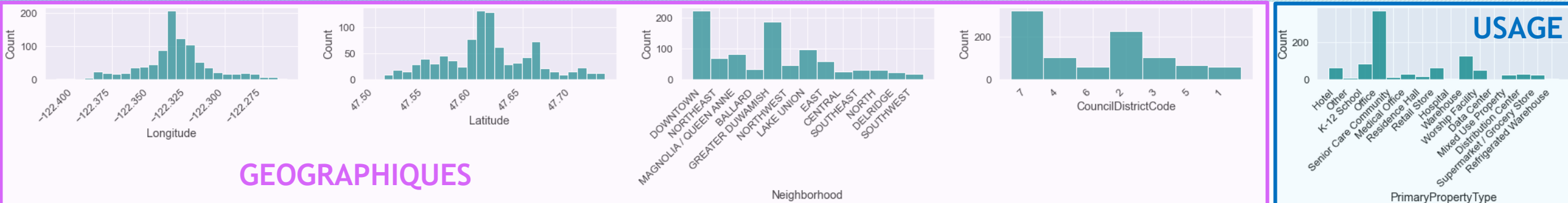
## EMISSION + CONSOMMATION



# Répartition empirique et modification des données

6

## Encodage des données catégorielles (OneHotEncoding)



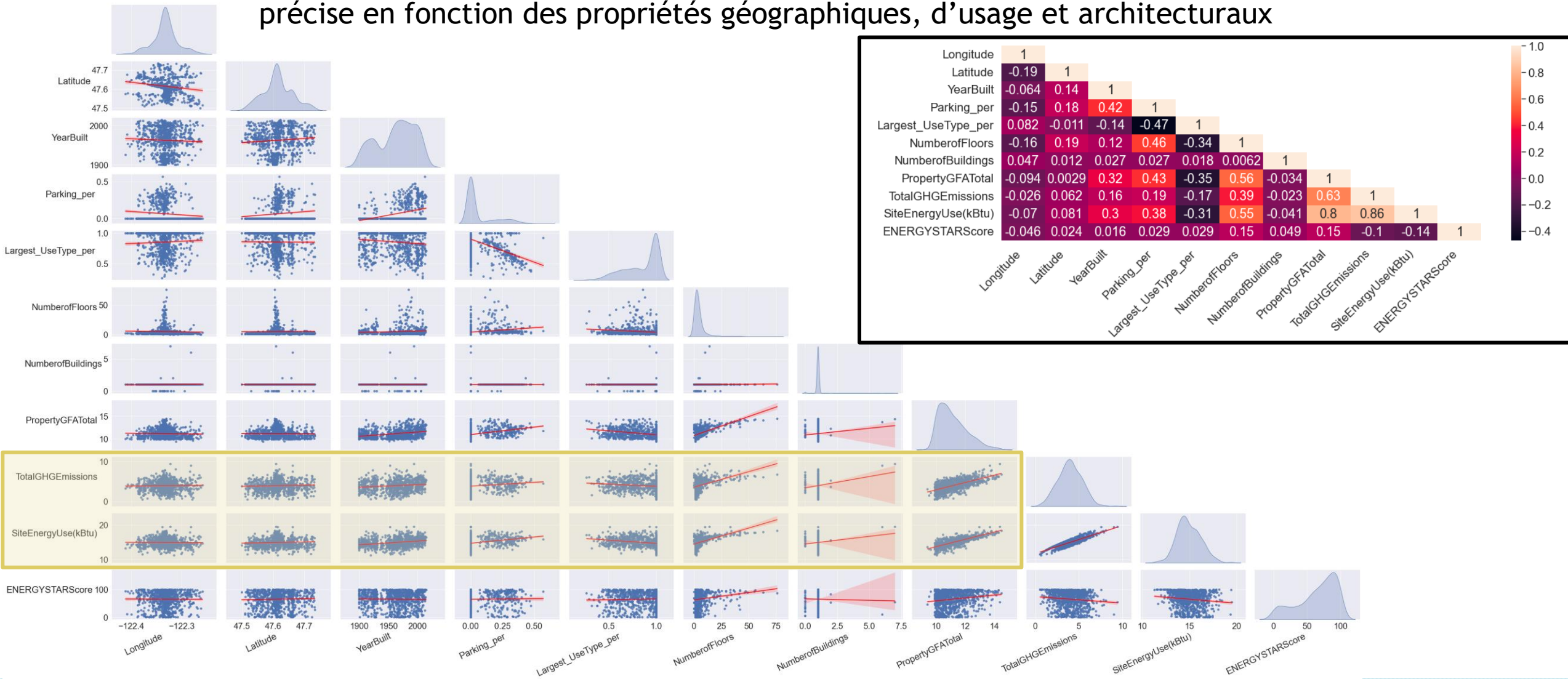
Transformation en logarithme naturel (+1)



# Brève analyse de la relation entre les variables

7

Objectif : Prédire la consommation en énergie et les émissions de GES de façon précise en fonction des propriétés géographiques, d'usage et architecturaux





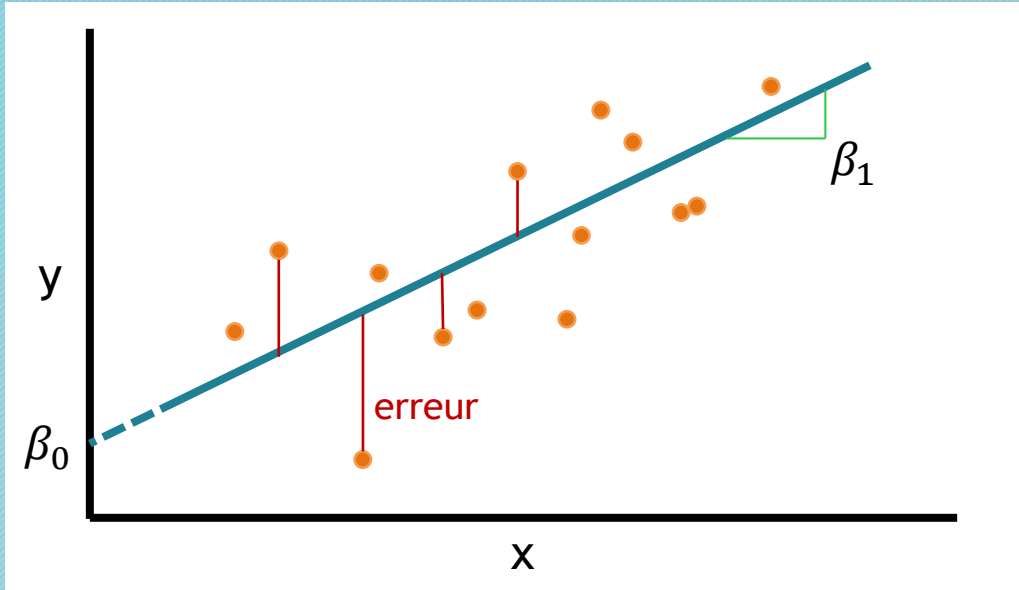
## Partie II : Méthodes de prédiction

8

- 1) Prédiction des valeurs de Consommation et des Emissions de GES
  - Comparaison : 10 modèles de Machine Learning succinctement optimisés par validation croisée
  - 4 critères (métriques)
  - Données standardisées ou non
  - Différents échantillons de départ (Random state de l'échantillonnage (split) des données d'entraînement/test)
- 2) Optimisation plus poussée du meilleur modèle + Courbe d'apprentissage
- 3) Importance des variables (dont ENERGYSTARScore)

# Les modèles utilisés : Modèles linéaires

10



$$y = \beta_0 + \beta_1 x_1 + \cdots + \beta_i x_i + \varepsilon$$

Objectif : Minimiser la somme des carrés des erreurs (ou résidus) = Méthode des Moindres Carrés (OLS)

Régression linéaire : Ordinary Least Squares

$$L_{OLS}(\hat{\beta}) = \sum_{i=1}^n (y_i - x_i' \hat{\beta})^2$$

Régression Ridge : OLS + L1 penalty

- Le poids des variables

$$L_{lasso}(\hat{\beta}) = \sum_{i=1}^n (y_i - x_i' \hat{\beta})^2 + \lambda \sum_{j=1}^m |\hat{\beta}_j|.$$

Régression Lasso : OLS + L2 penalty

- Le nombre de variables

$$L_{ridge}(\hat{\beta}) = \sum_{i=1}^n (y_i - x_i' \hat{\beta})^2 + \lambda \sum_{j=1}^m \hat{\beta}_j^2$$

Régression ElasticNet : OLS + L1 + L2 penalty

$$L_{enet}(\hat{\beta}) = \frac{\sum_{i=1}^n (y_i - x_i' \hat{\beta})^2}{2n} + \lambda \left( \frac{1-\alpha}{2} \sum_{j=1}^m \hat{\beta}_j^2 + \alpha \sum_{j=1}^m |\hat{\beta}_j| \right),$$



- Méthodes parallèles  
Bagging (tree) + Random Forest

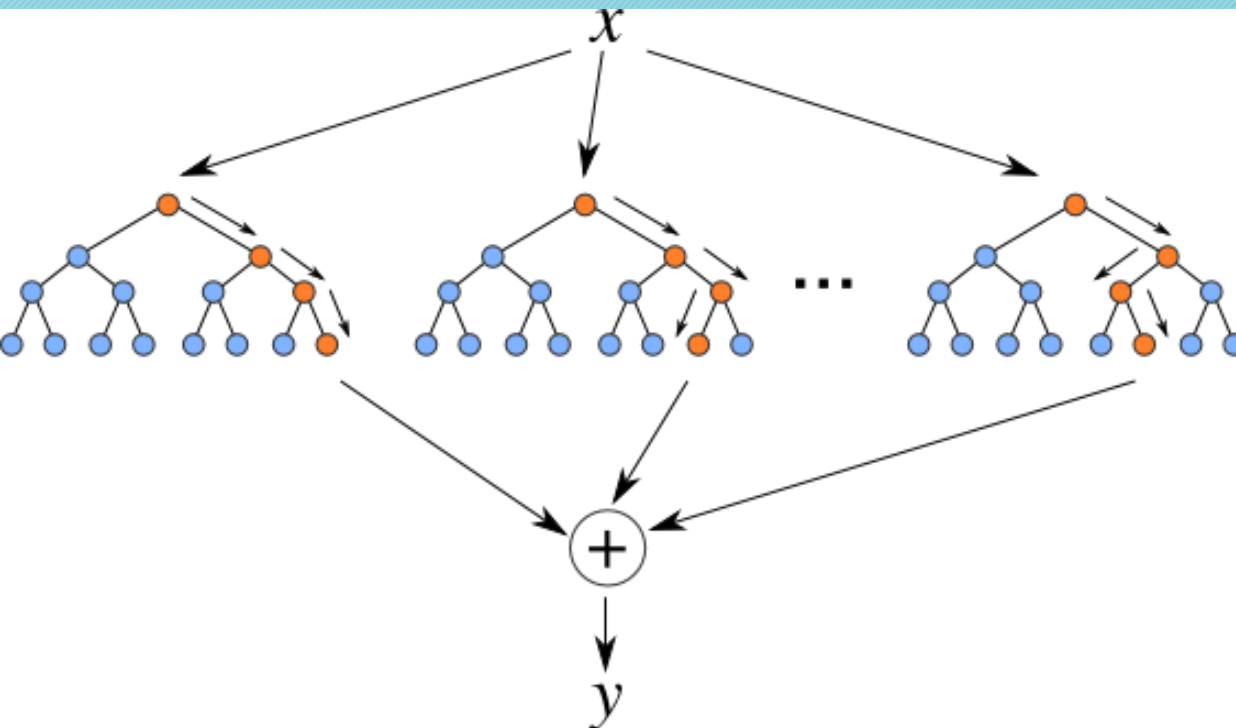


Image : Bakshi, Random Forest Regression. Medium, 2020

- Méthodes séquentielles  
Adaboost + (X)Gradient Boosting

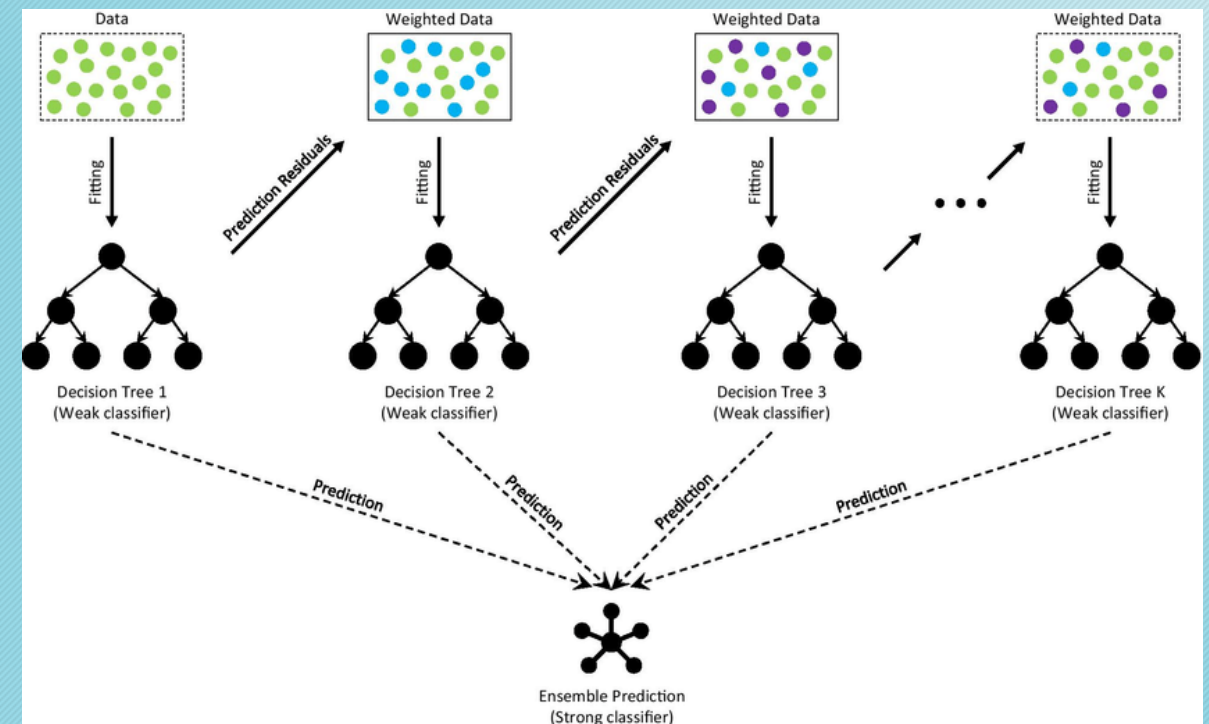


Image : Deng et al. BMC Medical Informatics and Decision Making 2021



- Méthode de support vector machine
- Méthode des K plus proches voisins
- Méthode de Perceptron multicouches
- Modèle nul (moyenne)

# Critères de comparaison : les différentes métriques

13

- Erreur quadratique moyenne (Mean Squared error MSE)  $MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$
- Erreur quadratique moyenne (Mean Squared error)  $RMSE = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}}$
- Coefficient de détermination  $R^2$  
$$R^2 = 1 - \frac{\overset{\text{erreur quadratique}}{\sum_{i=1}^n (y_i - \hat{y}_i)^2}}{\underset{\text{variance}}{\sum_{i=1}^n (y_i - \bar{y})^2}}$$
- Erreur Moyenne Absolue (Mean Absolute Error MAE)  $MAE = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i|$

- $y_i$  valeur vraie
- $\hat{y}_i$  valeur prédite
- $\bar{y}$  valeur moyenne

# Importance des Variables : Valeurs de Shapley

Lundberg & Lee 2017 :

Algorithme basé sur la théorie des jeux coopératifs (*Lloyd S. Shapley*)

La prédiction d'un individu = somme de contributions de chacune des variables.



## Partie III : Résultats

15

Résultats 1 : Prédiction de la consommation énergétique

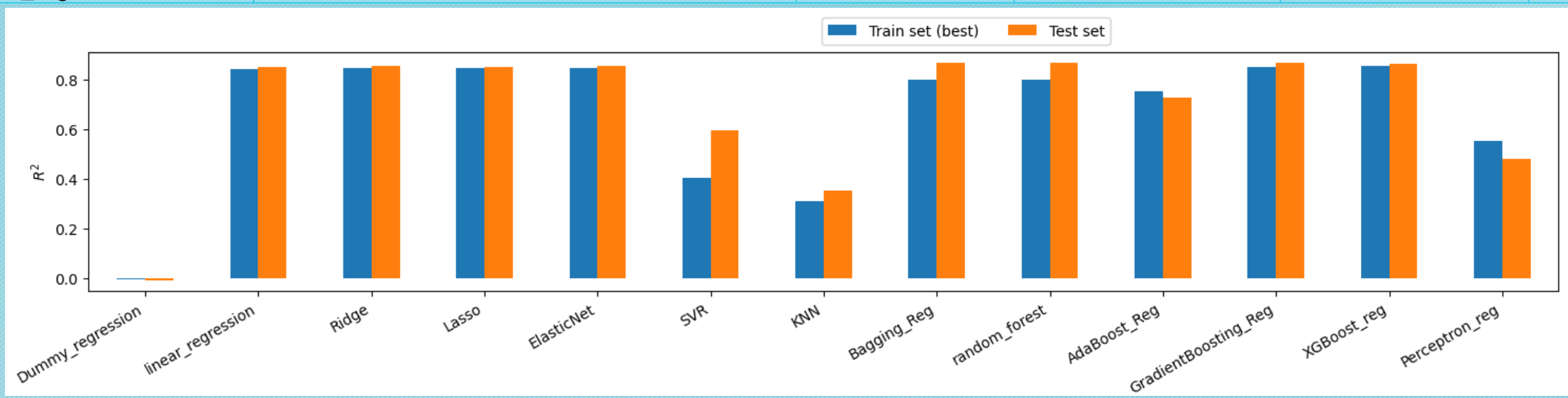
Résultats 2 : Prédiction des émissions de G.E.S.



# Résultats 1 : Prédictions de la consommation énergétique

16

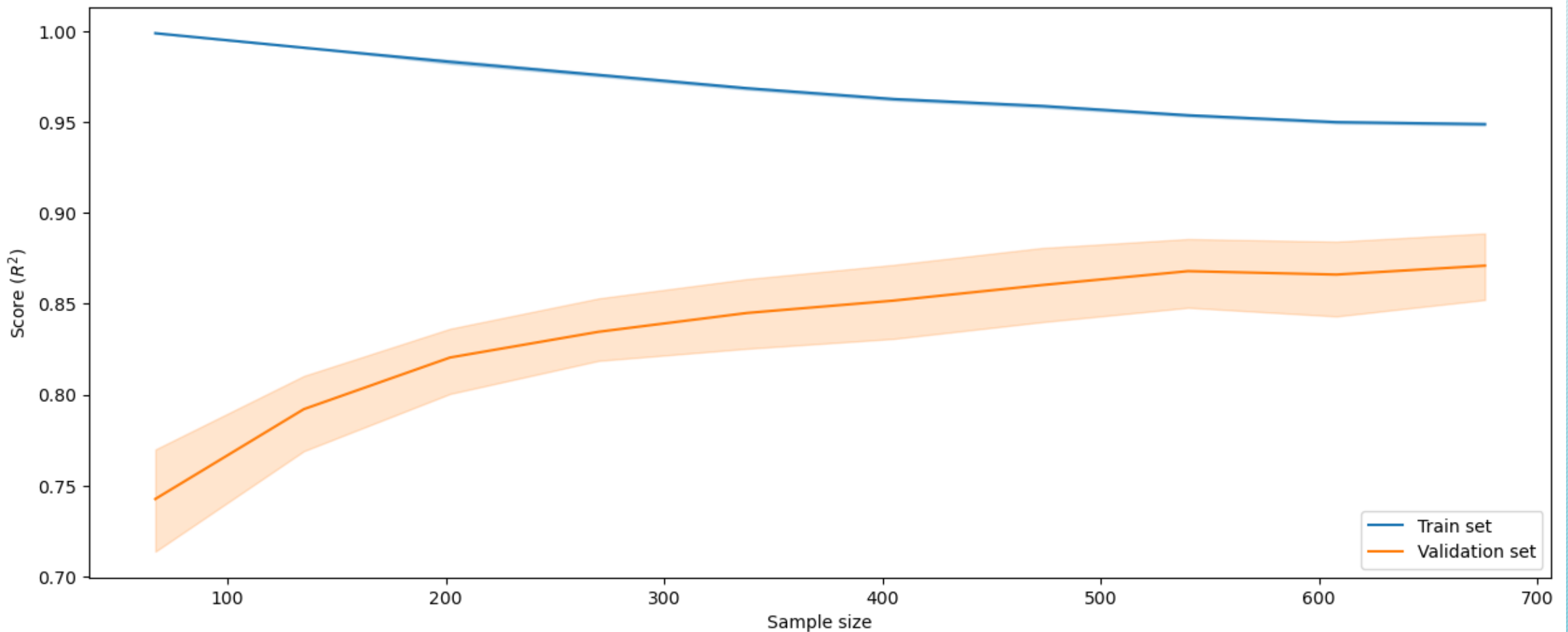
Modèle	$R^2$ : entraînement (meilleurs param)	$R^2$ : test	MSE : test	RMSE : test	MAE : test
Dummy_regression	-0.006	-0.000	1.178	1.387	0.957
linear_regression	0.857	0.871	0.423	0.179	0.331
Ridge	0.858	0.872	0.422	0.178	0.330
Lasso	0.864	0.871	0.423	0.179	0.329
ElasticNet	0.861	0.872	0.421	0.178	0.328
SVR	0.834	0.868	0.429	0.184	0.319
KNN	0.607	0.658	0.689	0.474	0.533
Bagging_Reg	0.809	0.847	0.460	0.212	0.344
random_forest	0.808	0.841	0.469	0.220	0.350
AdaBoost_Reg	0.740	0.776	0.557	0.310	0.458
GradientBoosting_Reg	0.848	0.884	0.401	0.161	0.300
XGBoost_reg	0.843	0.874	0.419	0.175	0.304
Perceptron_reg	0.713	0.872	0.421	0.177	0.321



# Résultats 1 : Optimisation des meilleurs modèles

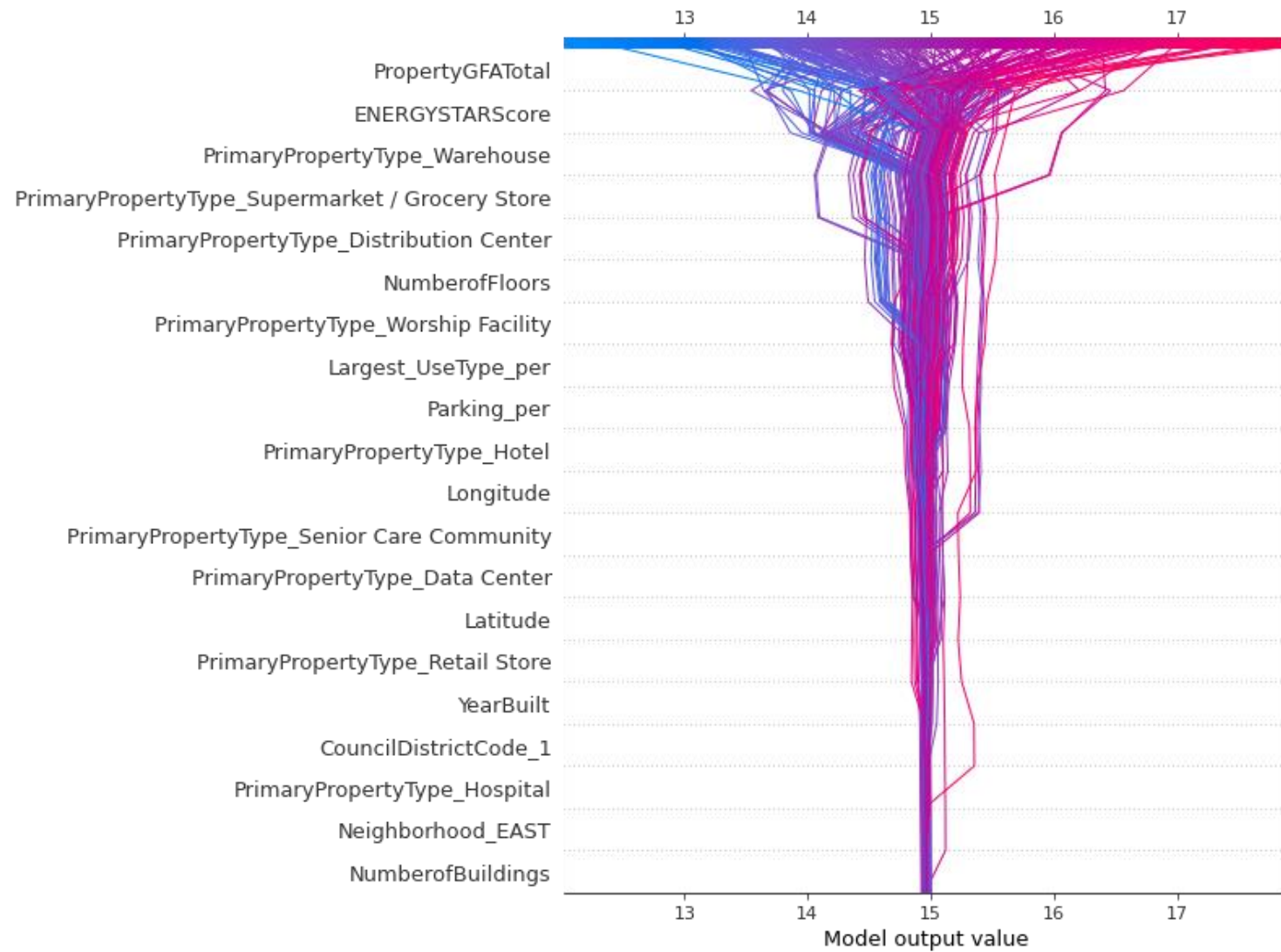
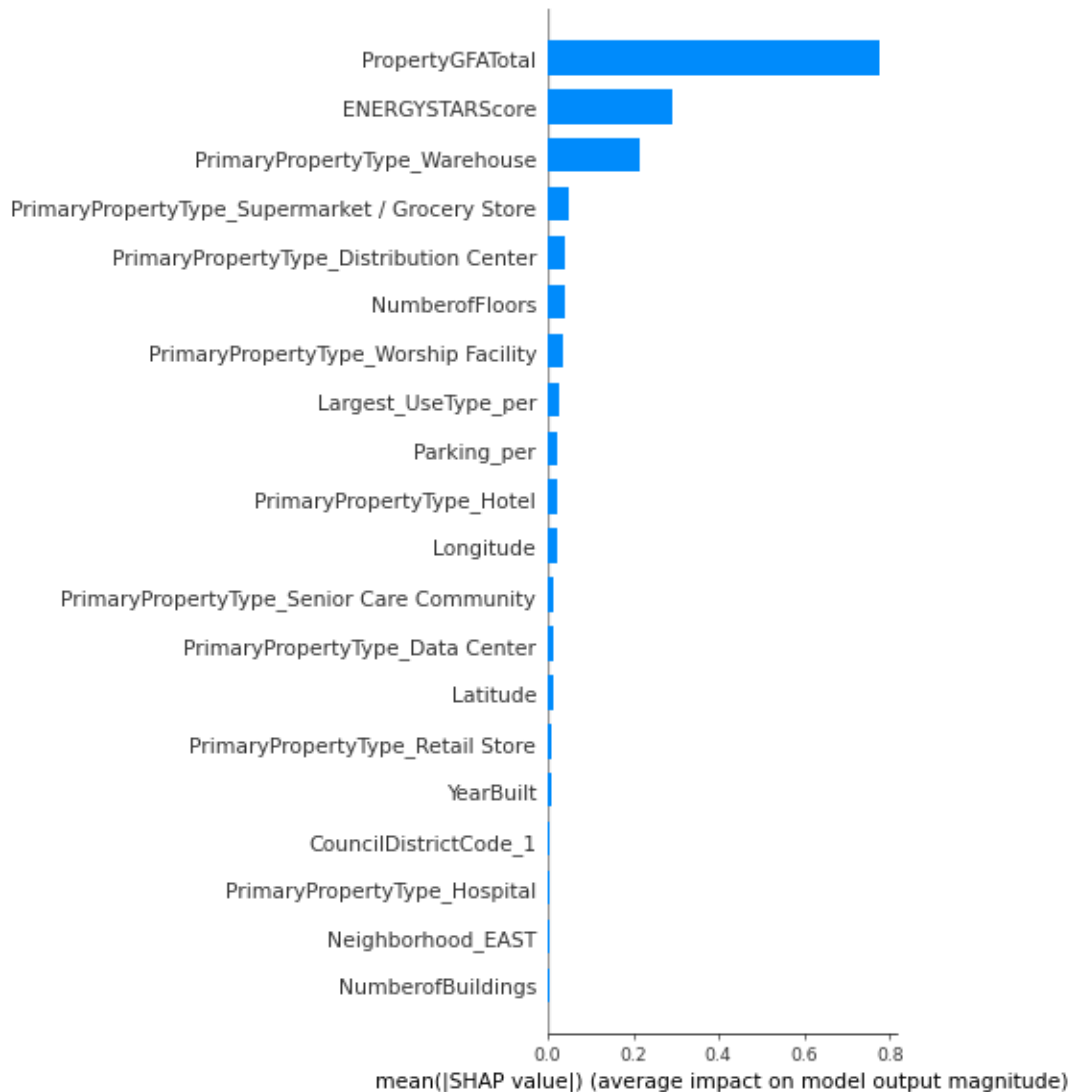
17

Données numériques	$R^2$ : test	MSE : test	RMSE : test	MAE : test
Non Standardisées	0.888	0.395	0.156	0.303



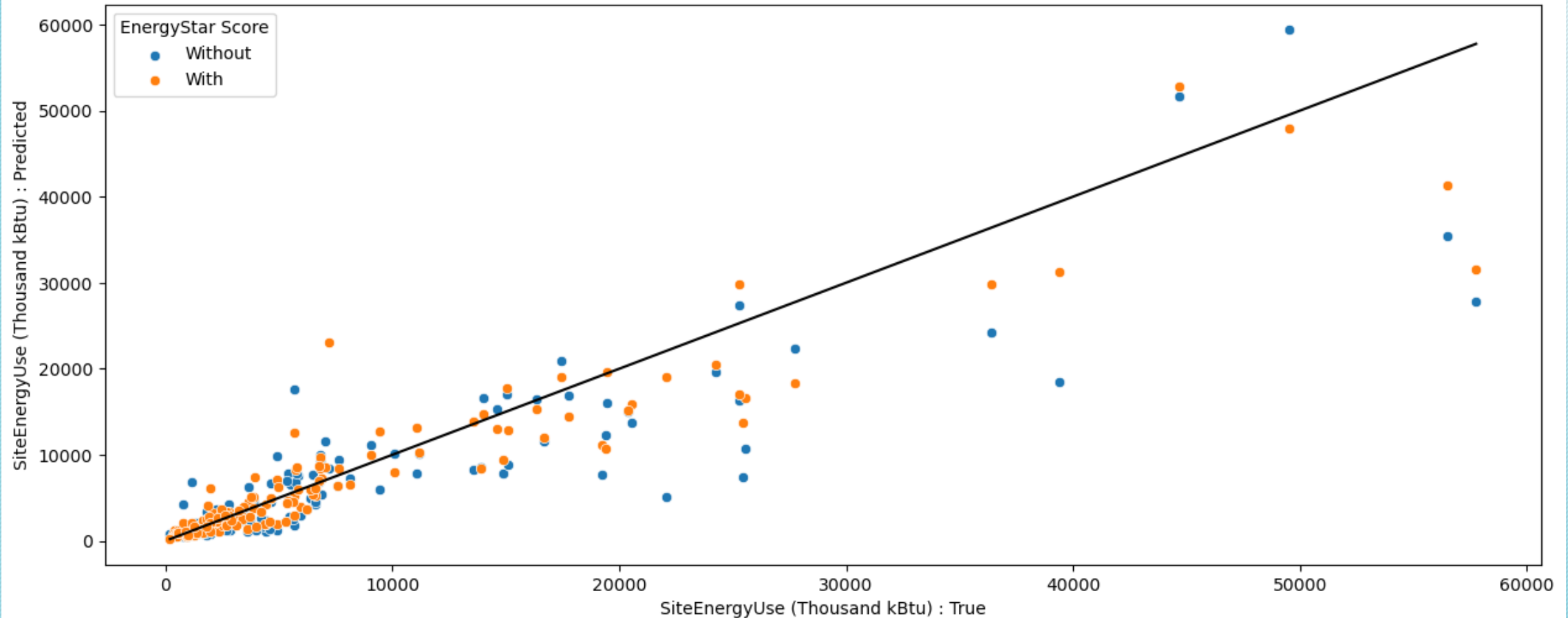
# Résultats 1 : Importances des variables

18



# Résultats 1 : Comparaison des valeurs prédites/vraies

19



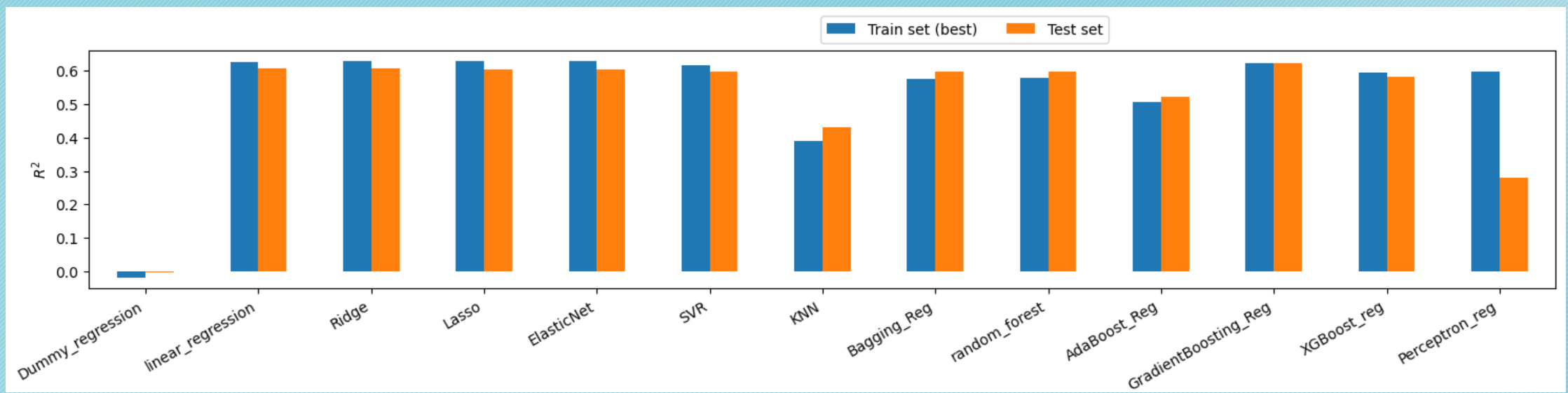
ENERGYSTARScore	$R^2$ : test	MSE : test	RMSE : test	MAE : test
Avec	0.888	0.395	0.156	0.303
Sans	0.793	0.536	0.287	0.406
Ecart	-10,7 %	+35,7 %	+84,0 %	+34,0 %



# Résultats 2 : Prédictions des émissions de G.E.S.

20

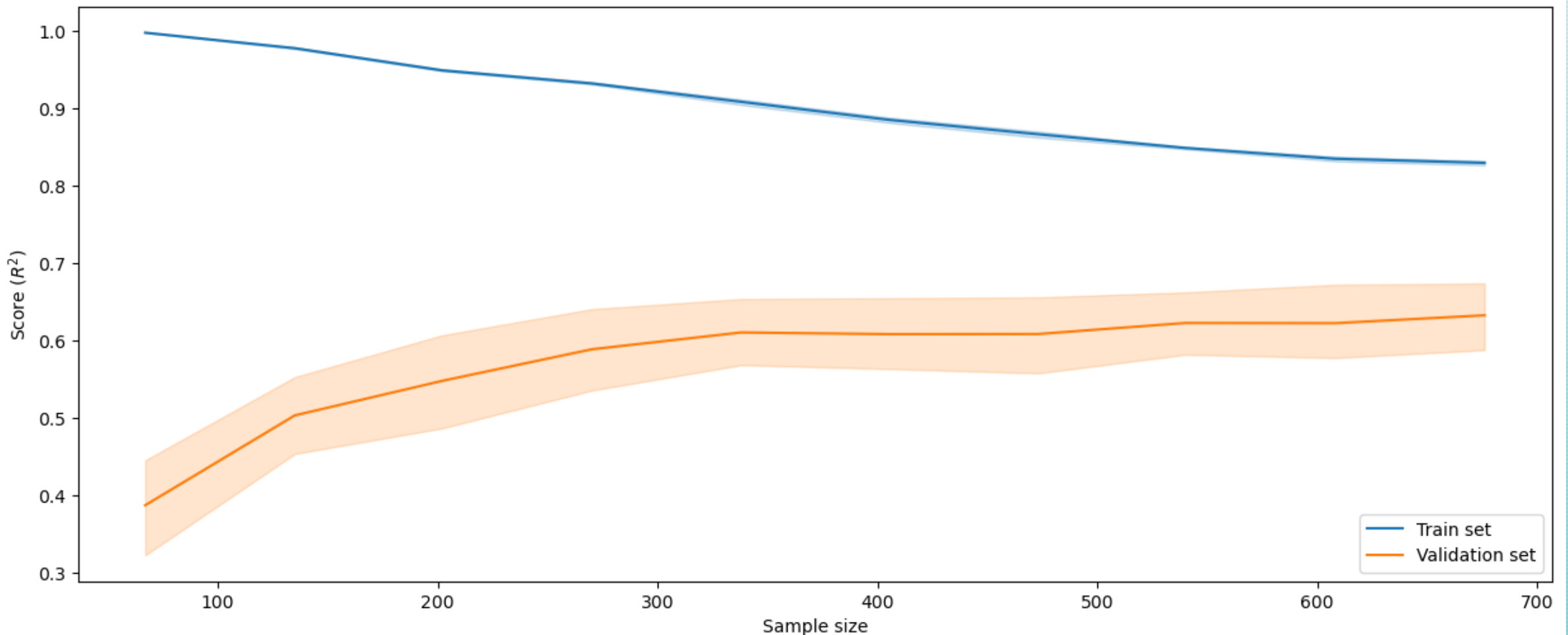
Modèle	$R^2$ : entraînement (meilleurs param)	$R^2$ : test	MSE : test	RMSE : test	MAE : test
Dummy_regression	-0.018	-0.003	1.339	1.792	1.062
linear_regression	0.626	0.607	0.839	0.703	0.685
Ridge	0.630	0.606	0.839	0.704	0.685
Lasso	0.630	0.603	0.843	0.710	0.686
ElasticNet	0.630	0.604	0.841	0.707	0.685
SVR	0.616	0.597	0.849	0.720	0.666
KNN	0.390	0.432	1.008	1.015	0.801
Bagging_Reg	0.577	0.598	0.848	0.718	0.676
random_forest	0.578	0.597	0.849	0.721	0.677
AdaBoost_Reg	0.508	0.523	0.923	0.853	0.778
GradientBoosting_Reg	0.622	0.623	0.821	0.674	0.653
XGBoost_reg	0.594	0.582	0.864	0.747	0.683
Perceptron_reg	0.597	0.281	1.134	1.285	0.901



# Résultats 2 : Optimisation des meilleurs modèles

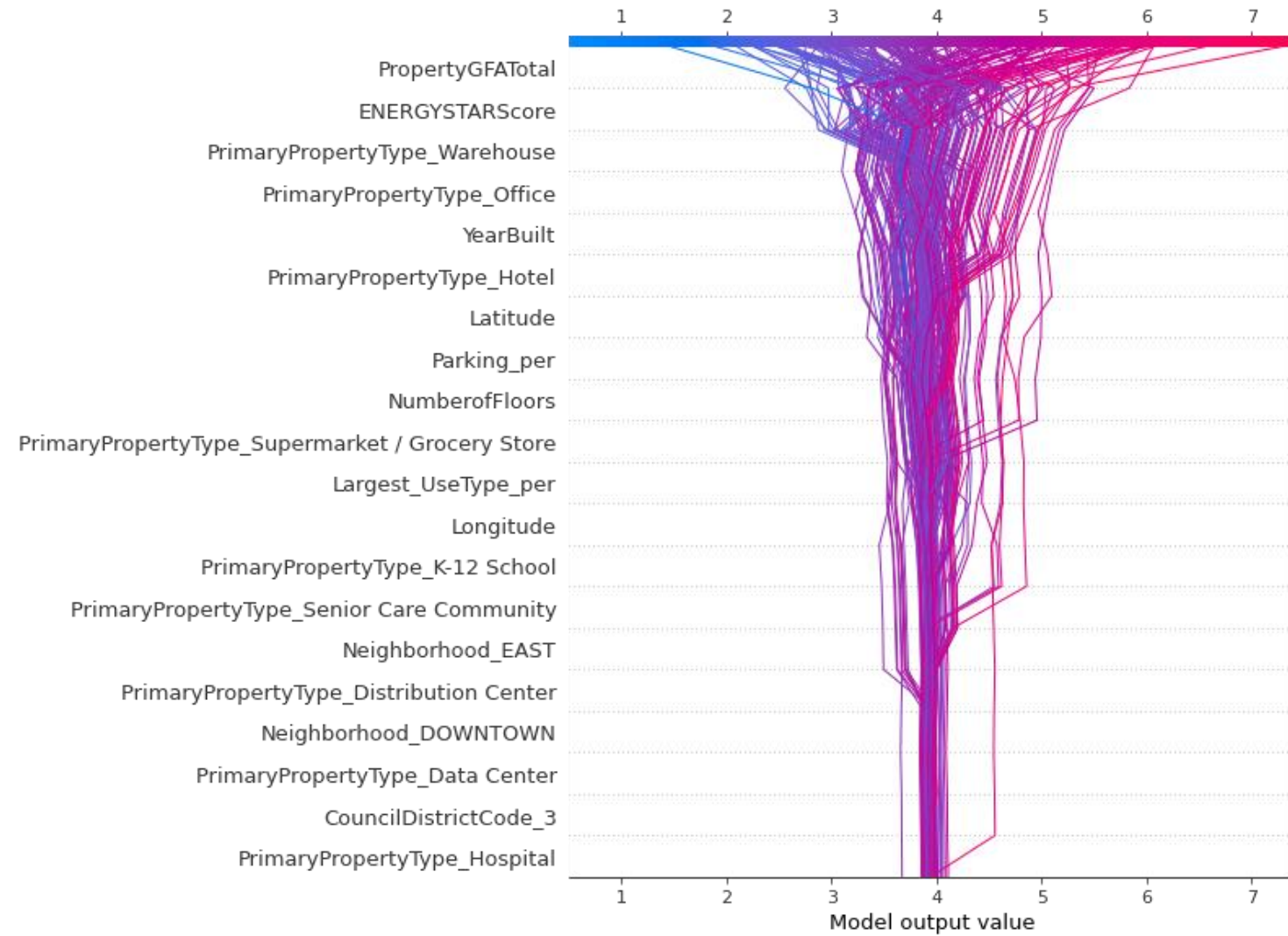
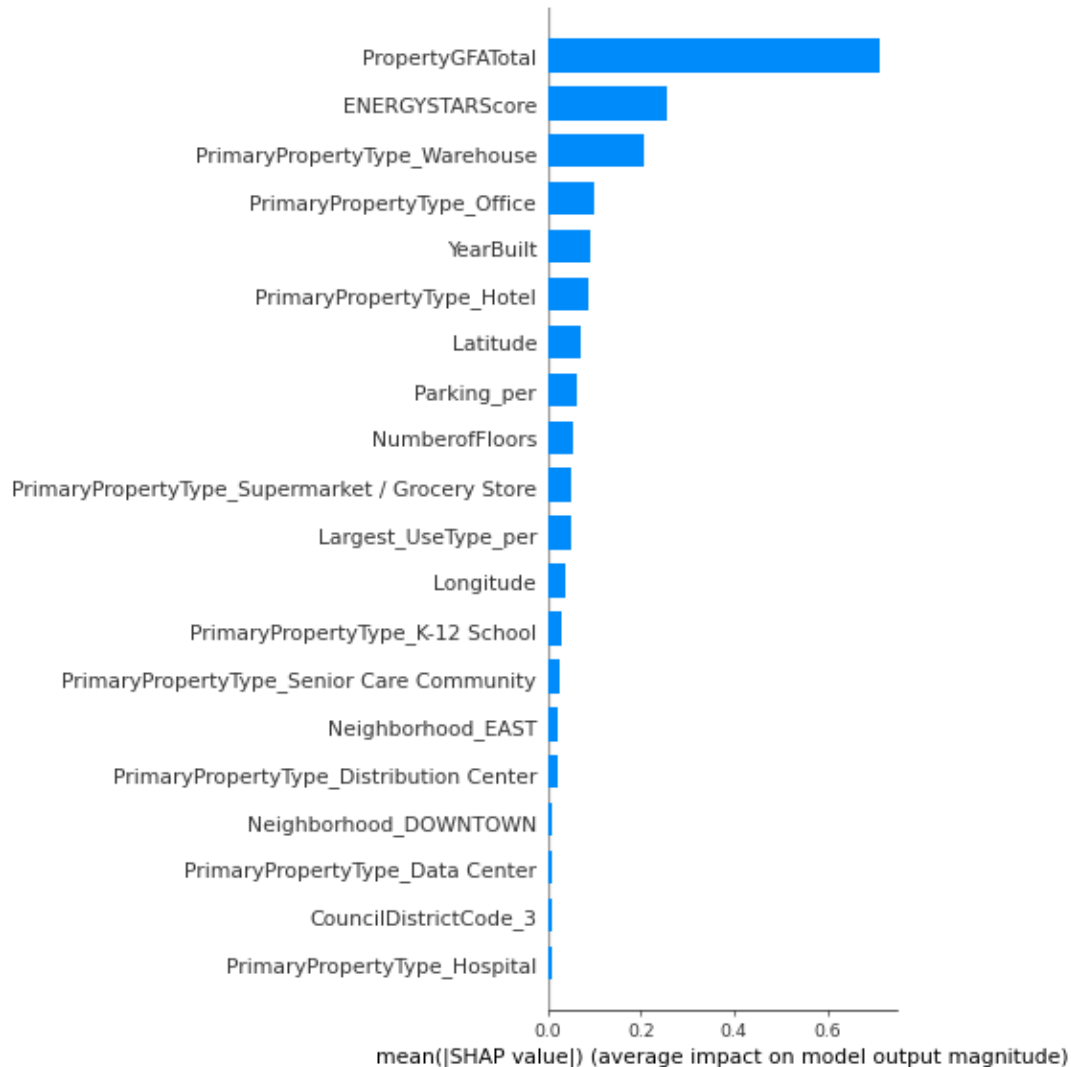
21

Données numériques	$R^2$ : test	MSE : test	RMSE : test	MAE : test
Non standardisées	0.636	0.806	0.65	0.652



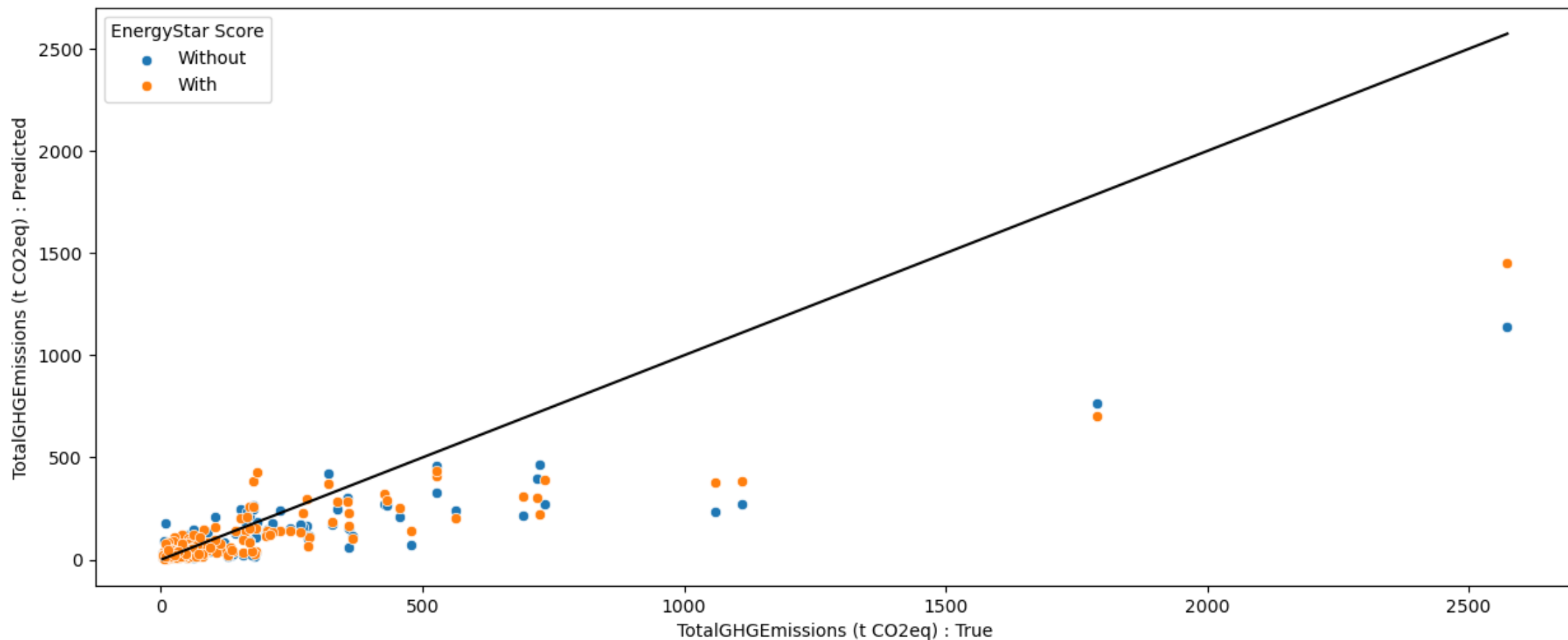
# Résultats 2 : Importances des variables

22



# Résultats 2 : Comparaison des valeurs prédites/vraies

23



ENERGYSTARScore	$R^2$ : test	MSE : test	RMSE : test	MAE : test
Avec	0.636	0.806	0.65	0.652
Sans	0.594	0.852	0.726	0.67
Ecart	-6,6 %	+5,7 %	+11,7 %	+2,8 %



- 1) Le modèle prédit avec précision la consommation énergétique des bâtiments non destinées à l'habitation.
- 2) Les prédictions des émissions de G.E.S. sont moins précises.
- 3) Dans les deux cas, ENERGYSTARScore augmentent la précision des prédictions.

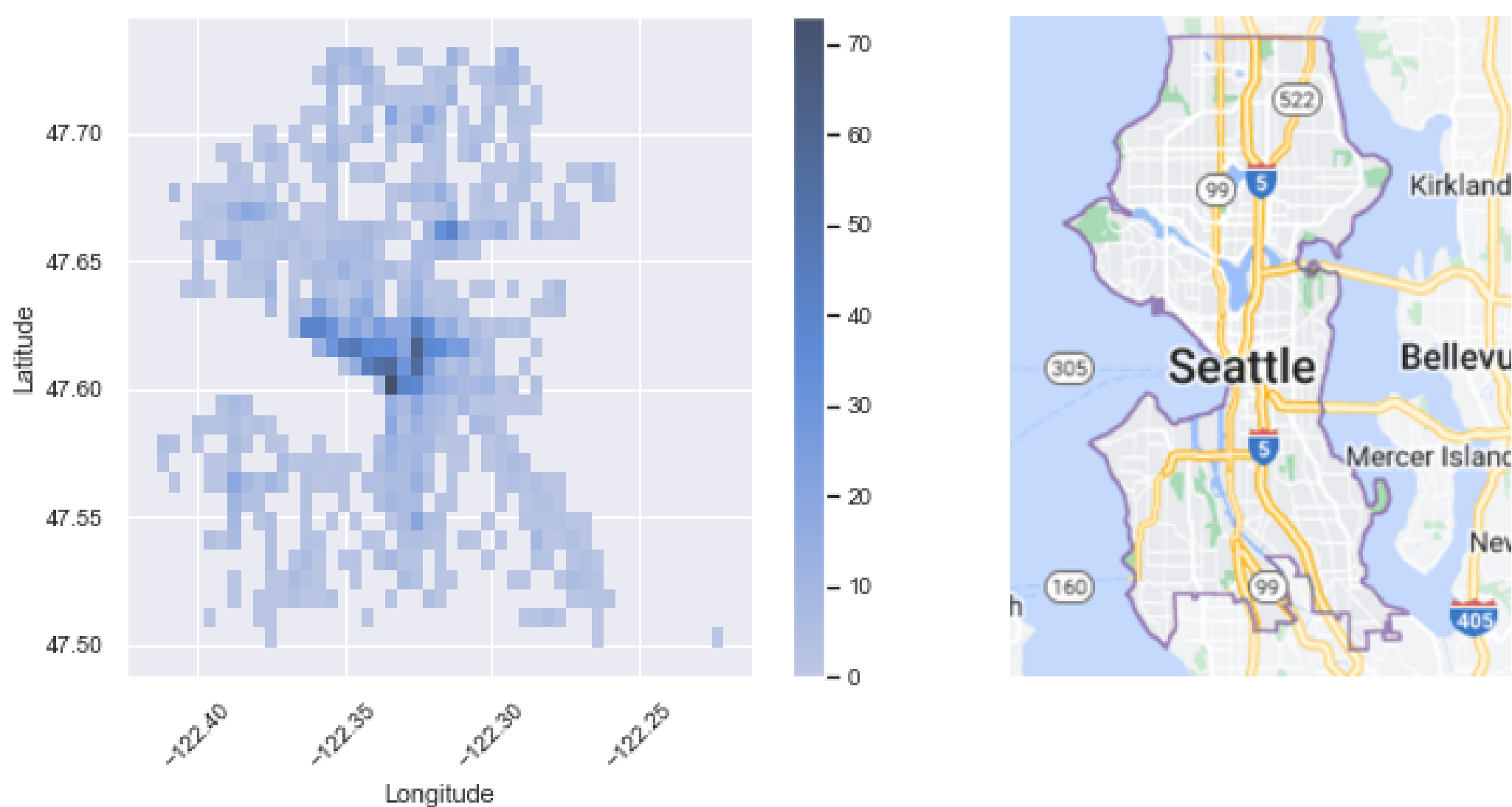


## Annexe : Analyse exploratoire des indicateurs

15

# Indicateurs géographiques : Longitude latitude

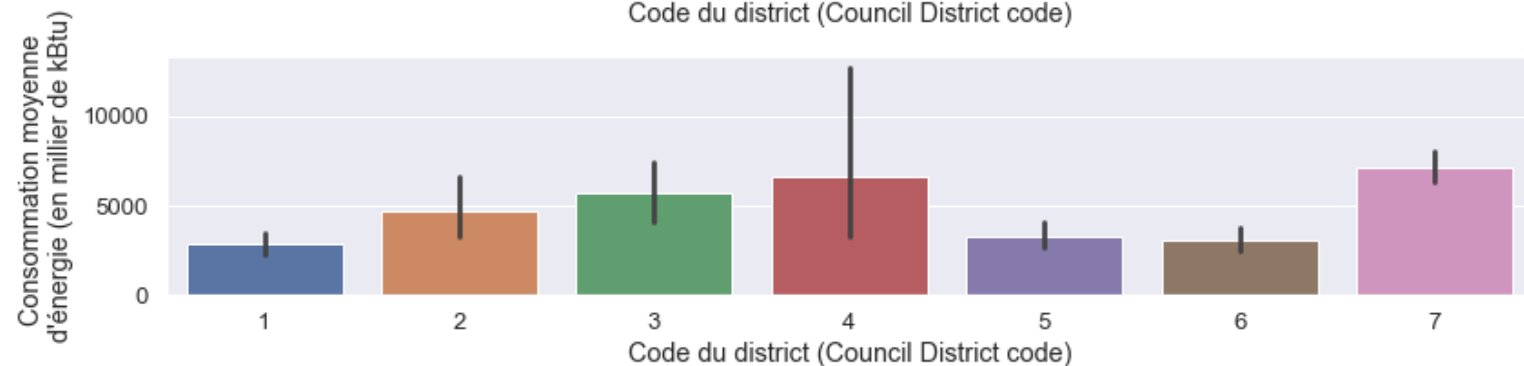
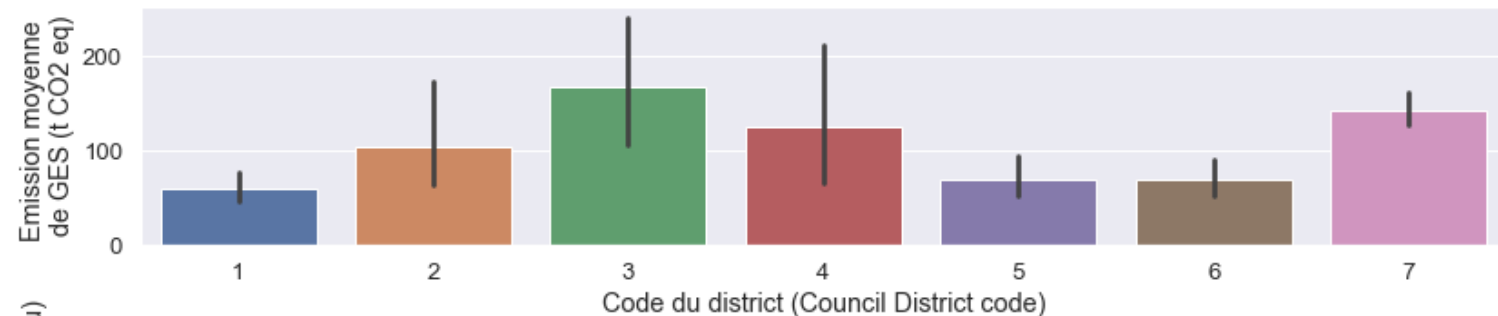
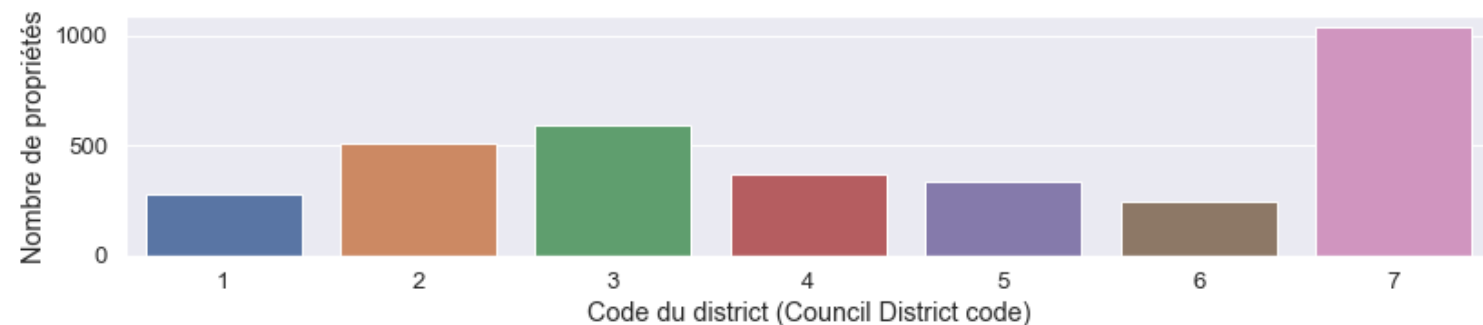
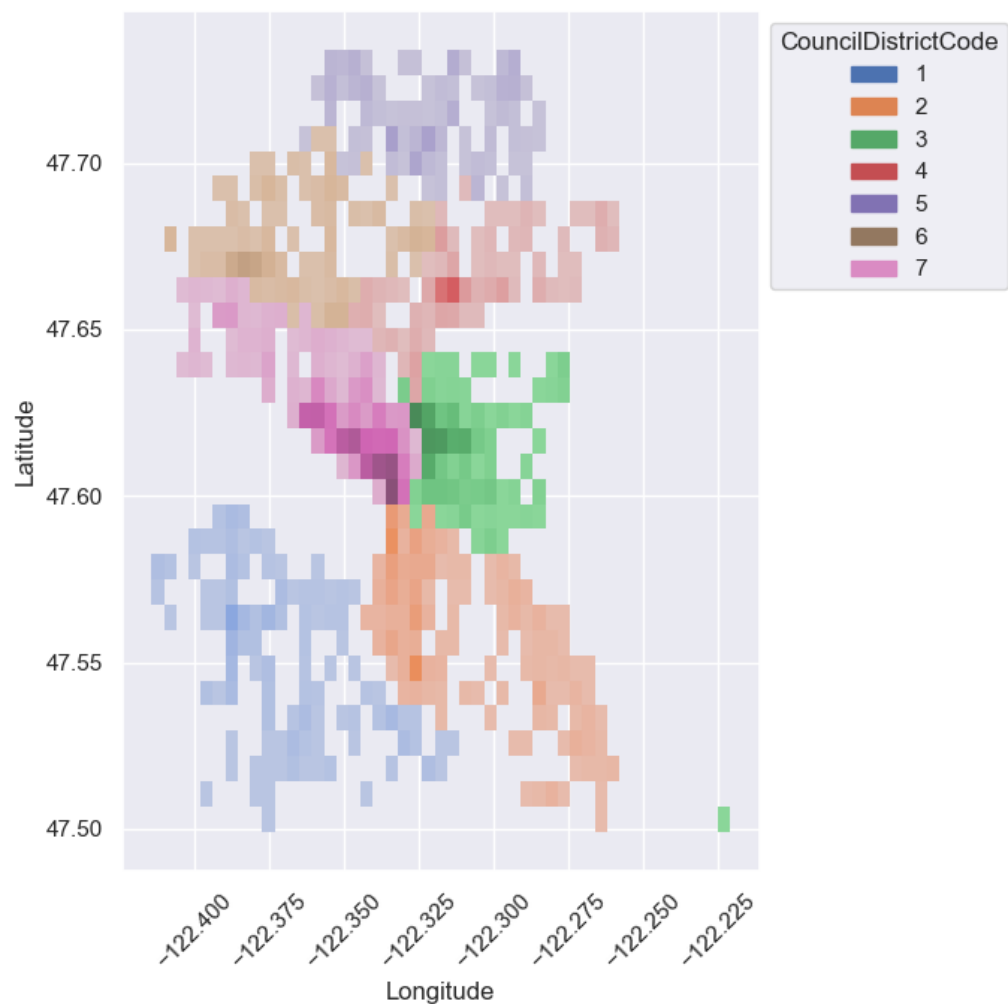
28





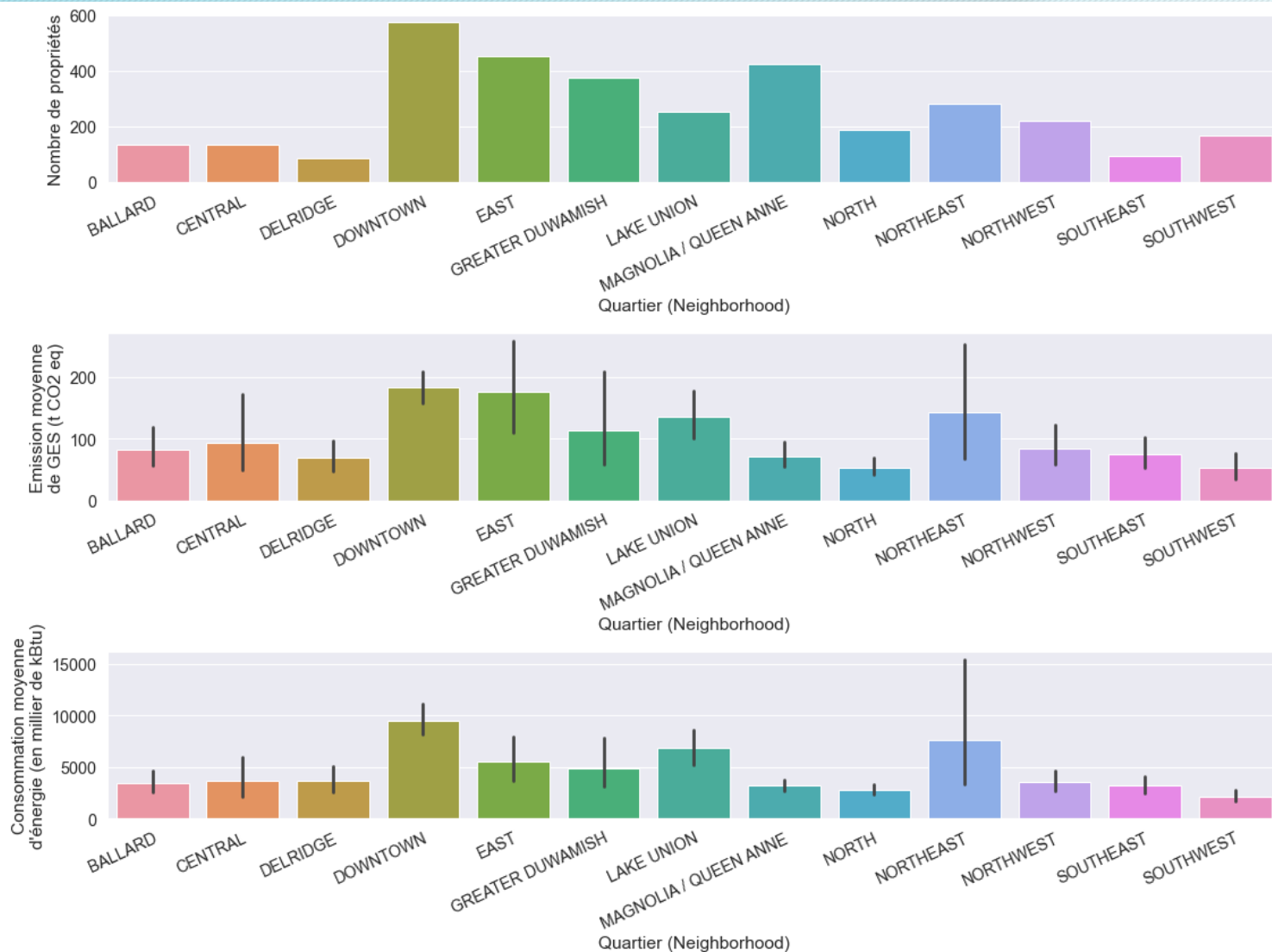
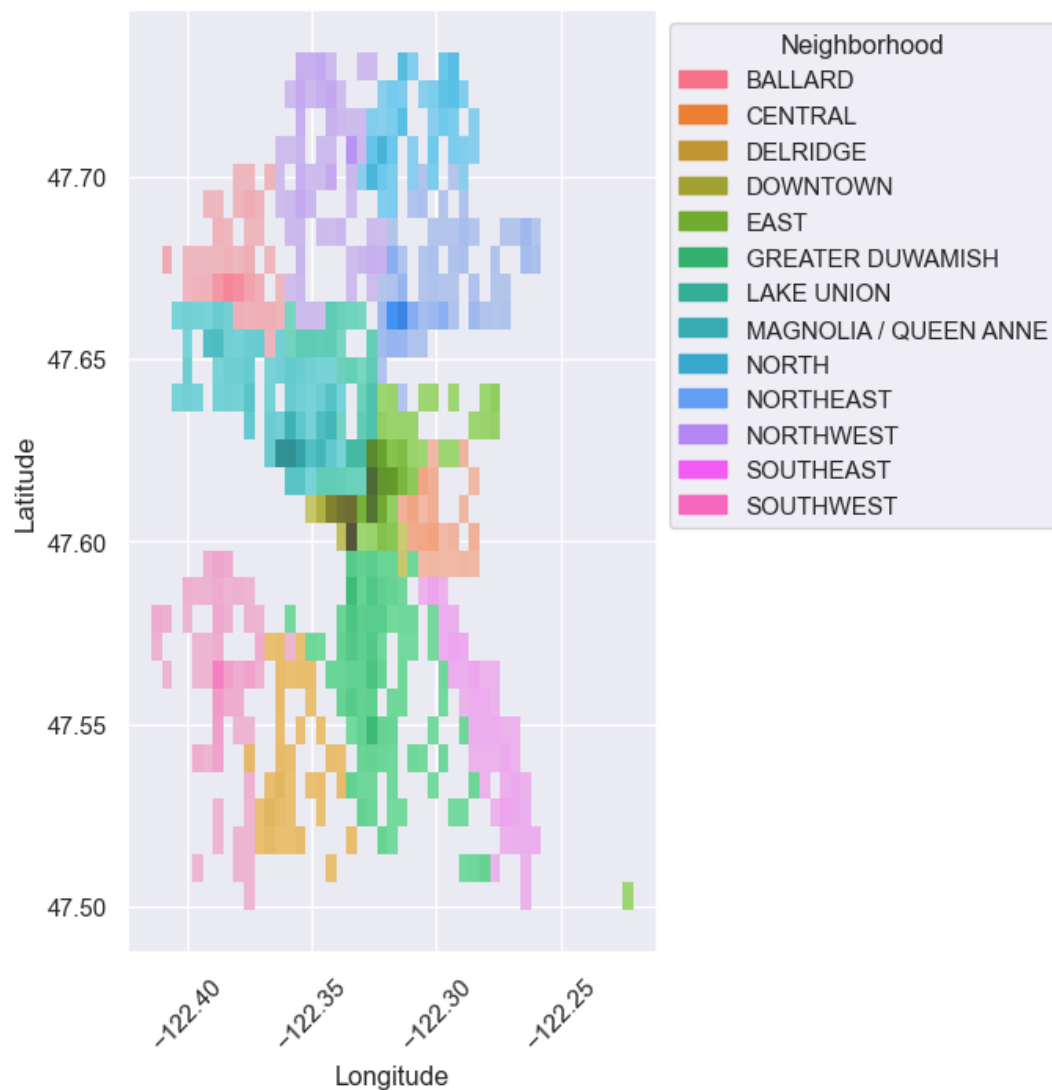
# Indicateurs géographiques : District

29

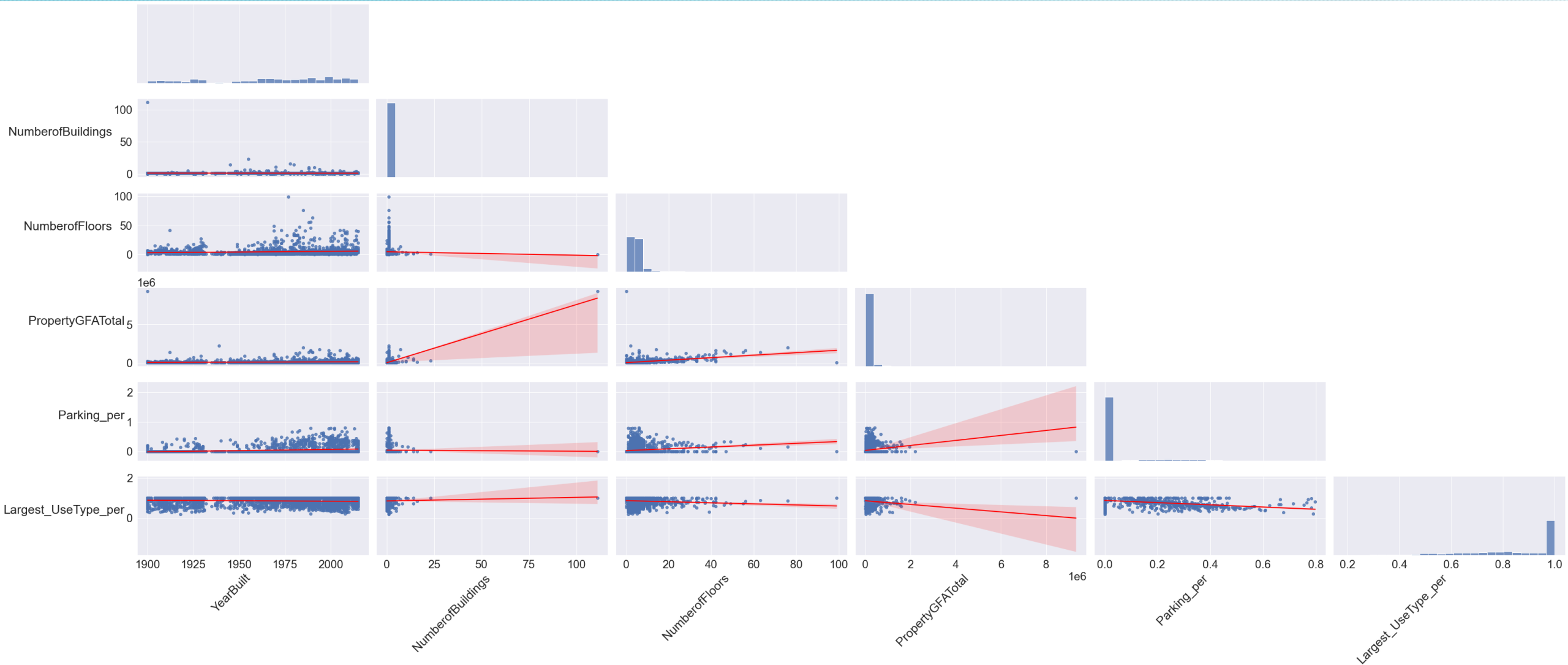


# Indicateurs géographiques : Quartier

30

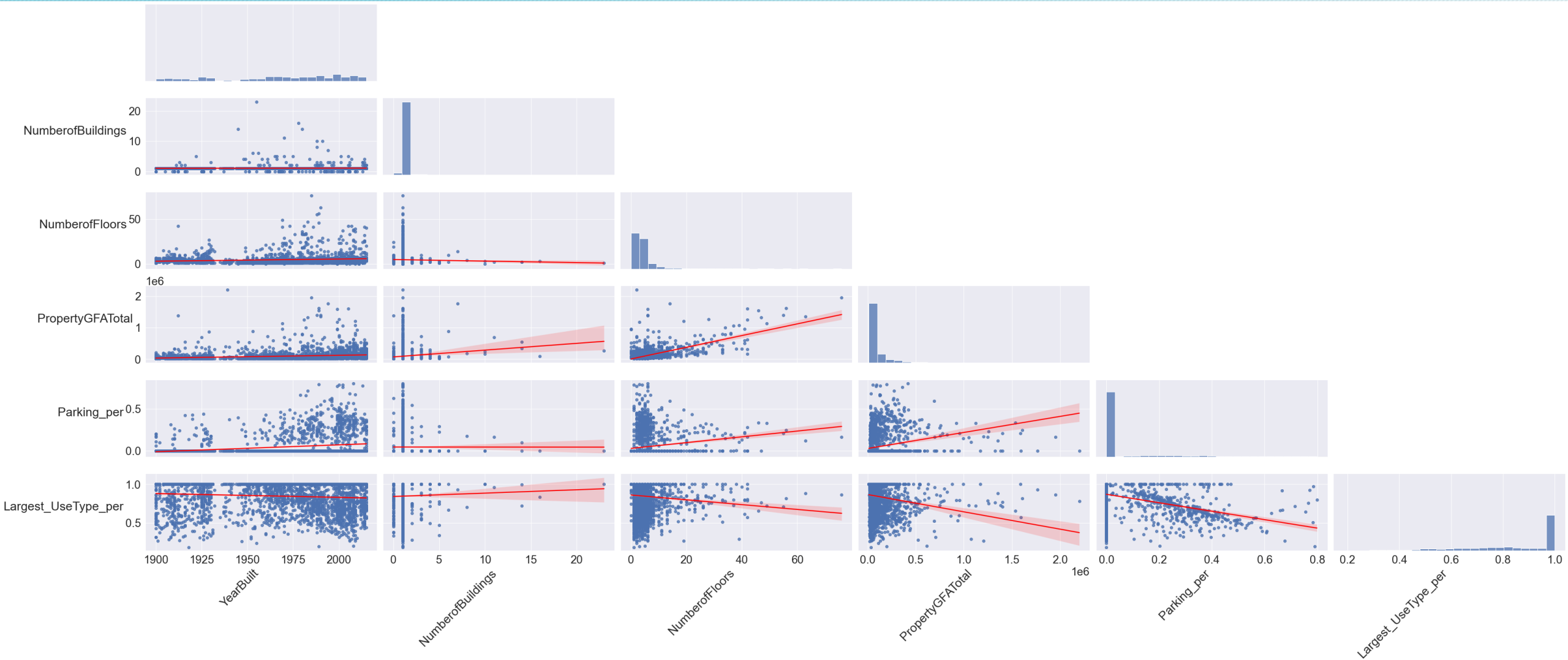


## 31



# Indicateurs architecturaux :

32





# Indicateurs d'usages : Type de propriété principale

