

A wide-angle photograph of the Seattle skyline. The Space Needle is prominent on the left. The city is filled with various high-rise buildings. In the background, Mount Rainier is visible under a clear blue sky. The text 'Anticipez les besoins en consommation de bâtiments' is overlaid in white, bold, sans-serif font.

# Anticipez les besoins en consommation de bâtiments

# Objectifs

Seattle ville neutre en 2050, pour cela il faut anticiper la consommation d'énergie totale et les émissions en CO2 des bâtiments non destinés à l'habitation

Evaluer l'intérêt de l'ENERGY STAR Score sur la prédiction

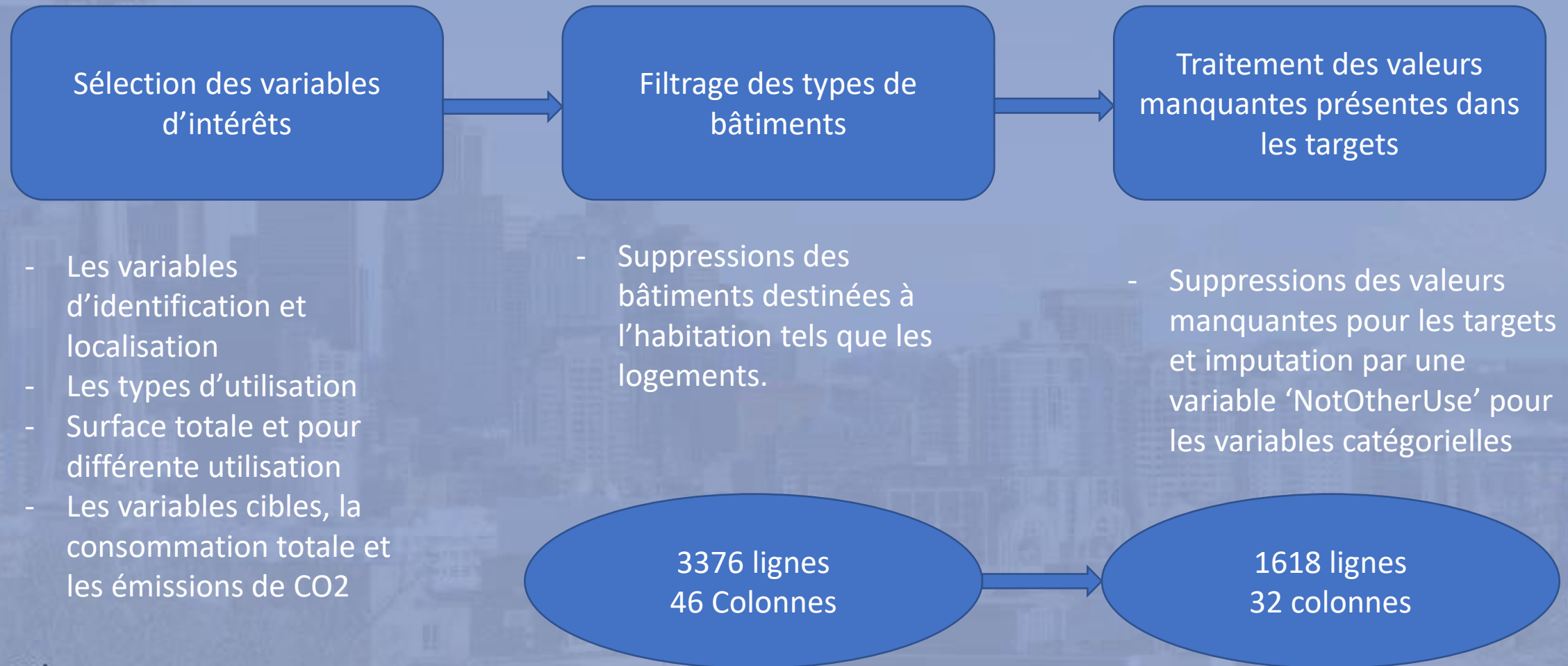
# Jeu de données

- Le jeu de données présente 3376 lignes et 46 colonnes

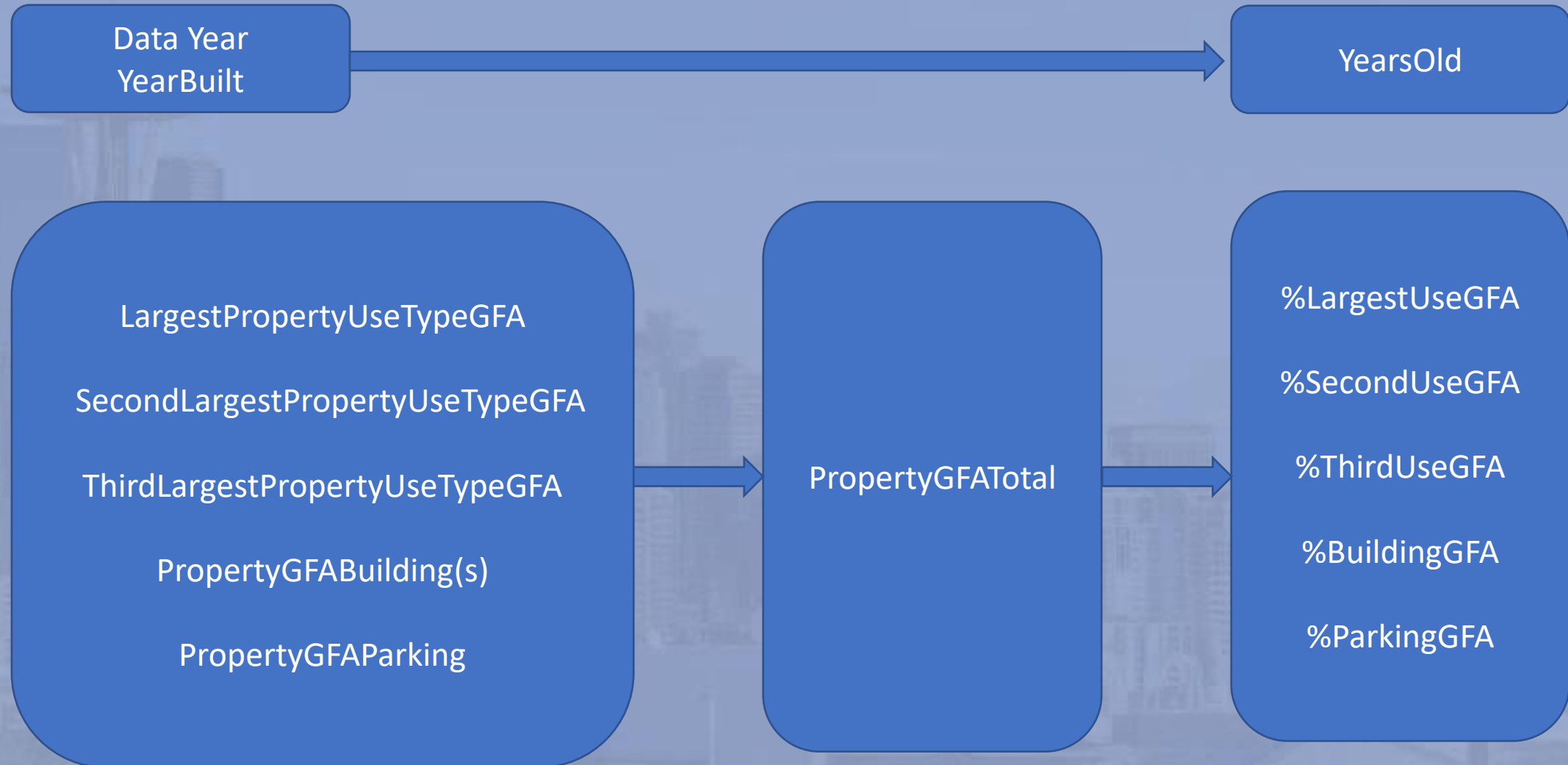
Soit 3376 bâtiments et 46 variables qui contiennent des informations sur les bâtiments:

- Variables d'identification et de localisation
- Variables de dimensions et d'utilisation
- Variables de consommations et d'émissions

# Nettoyage des données



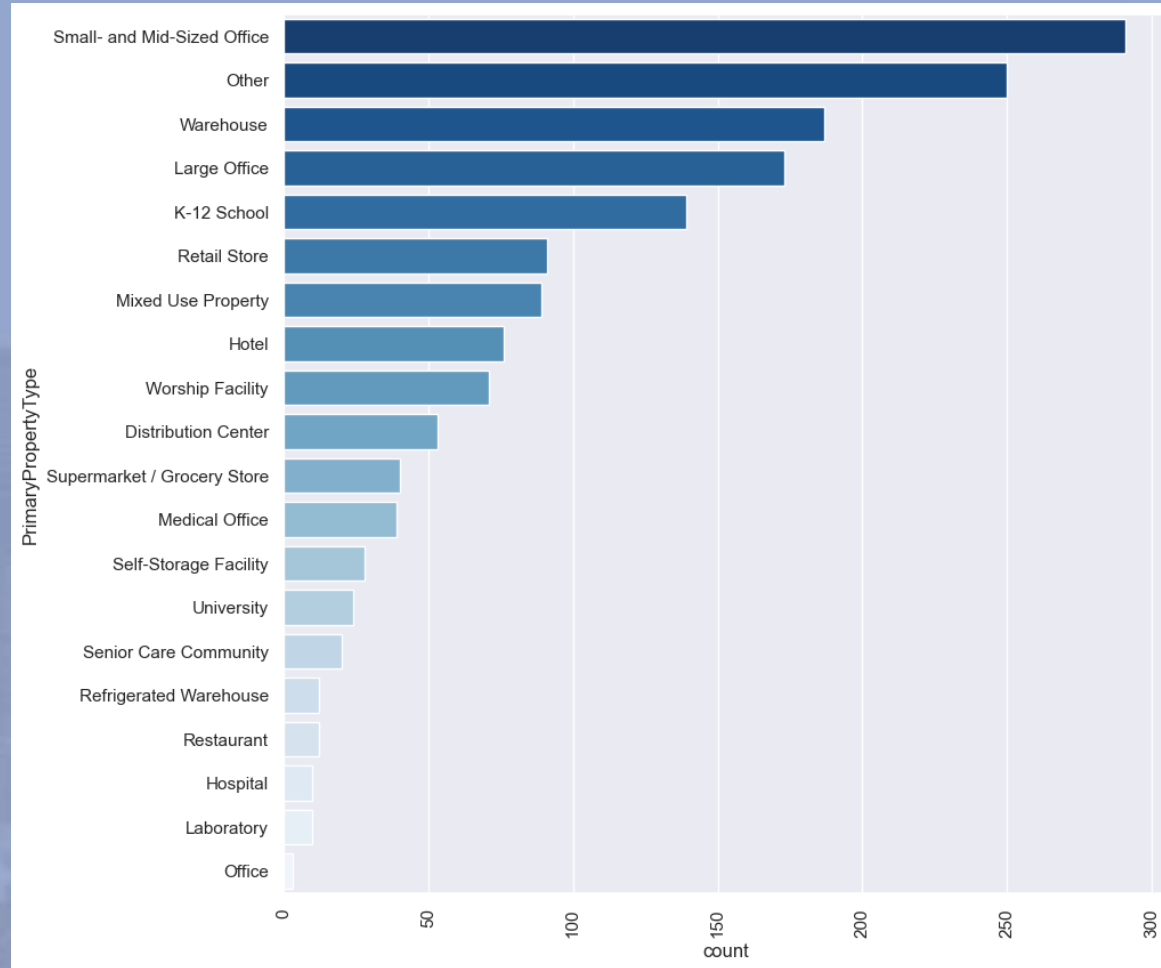
# Création de variables



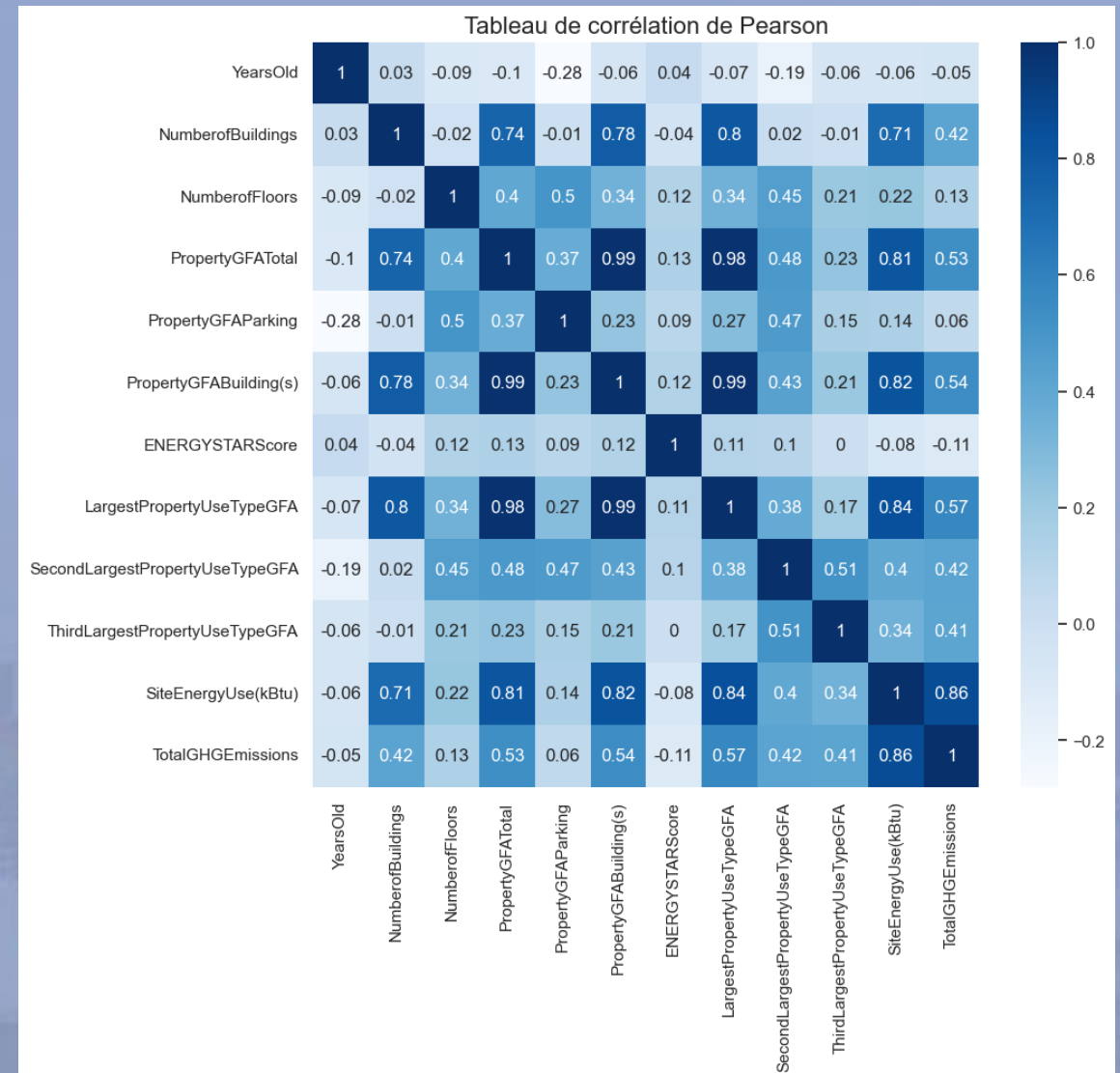
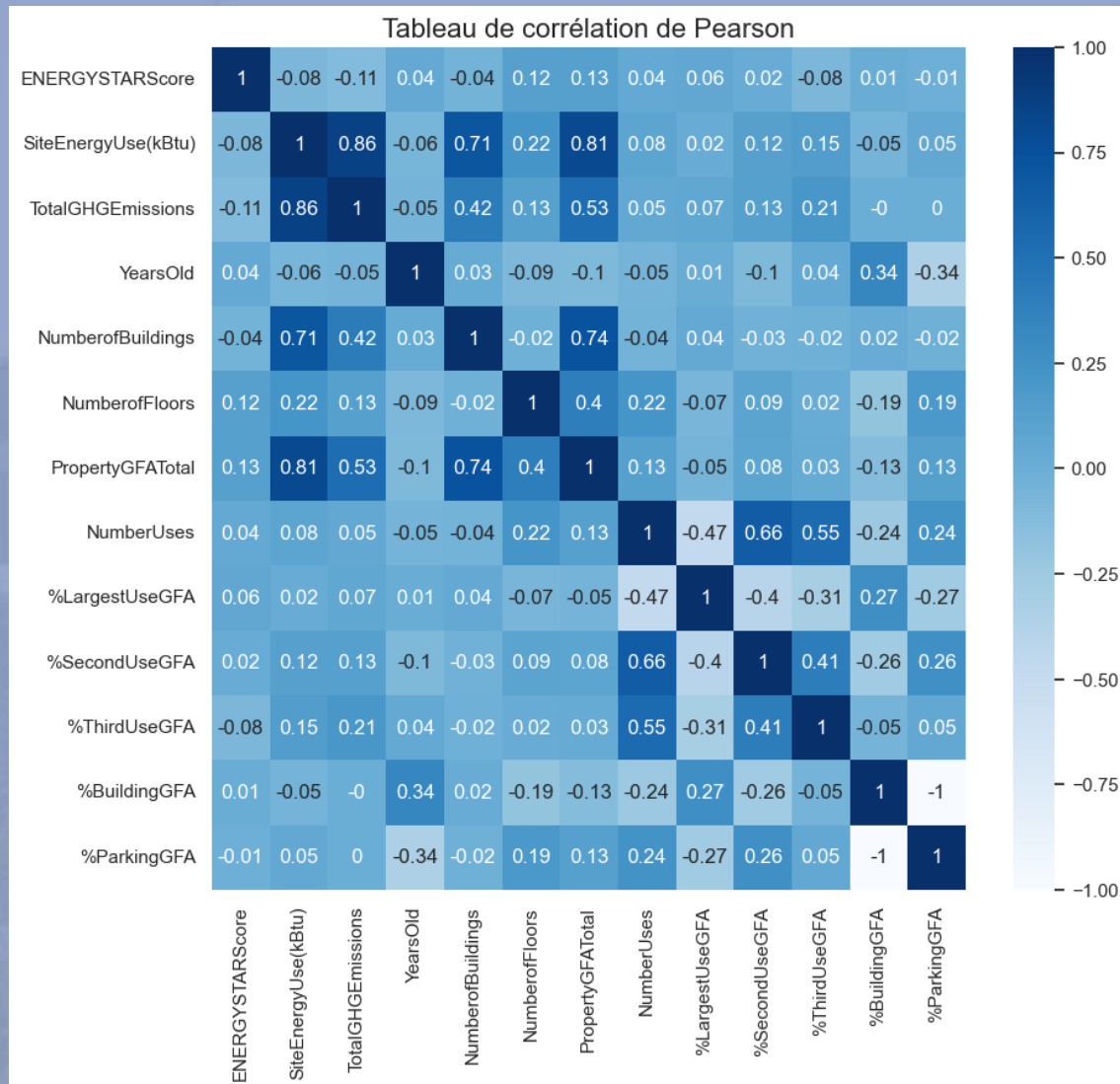


# Analyse

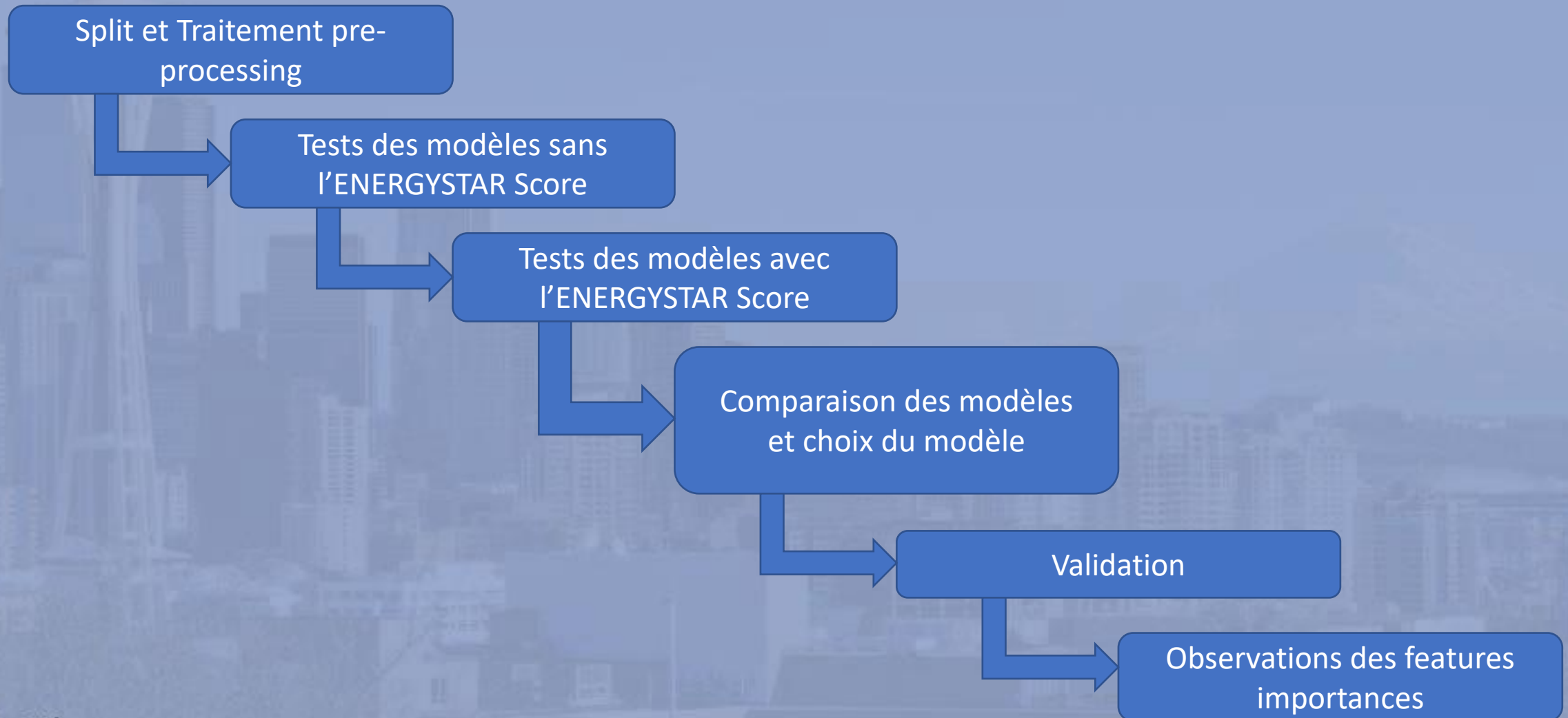
## Répartition de la catégorie 'PrimaryPropertyType'



- Corrélation de Pearson entre les variables quantitatives



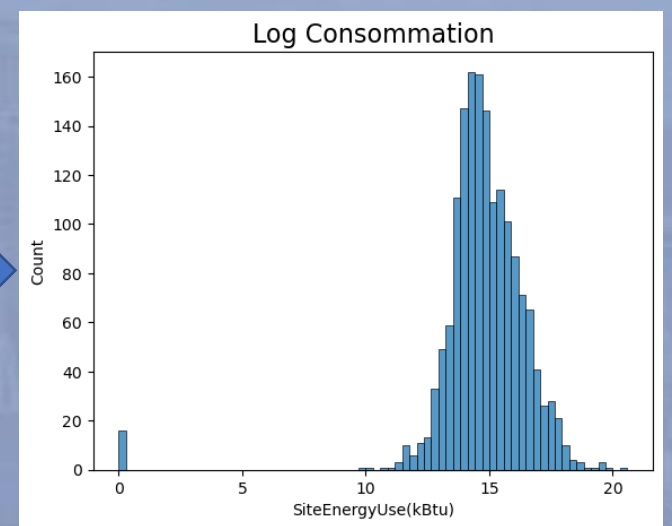
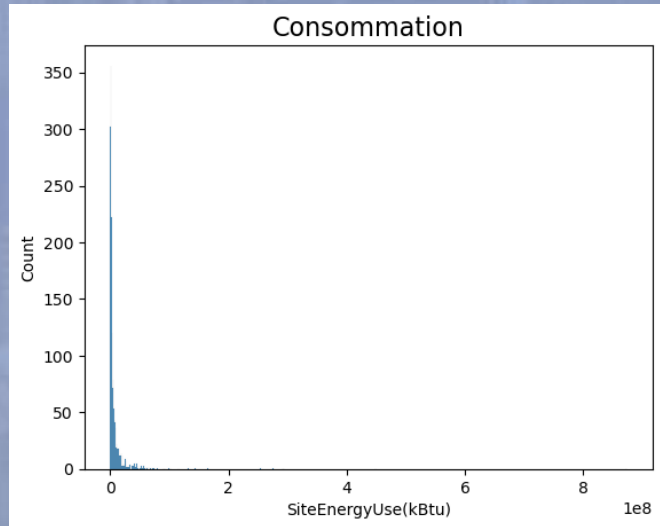
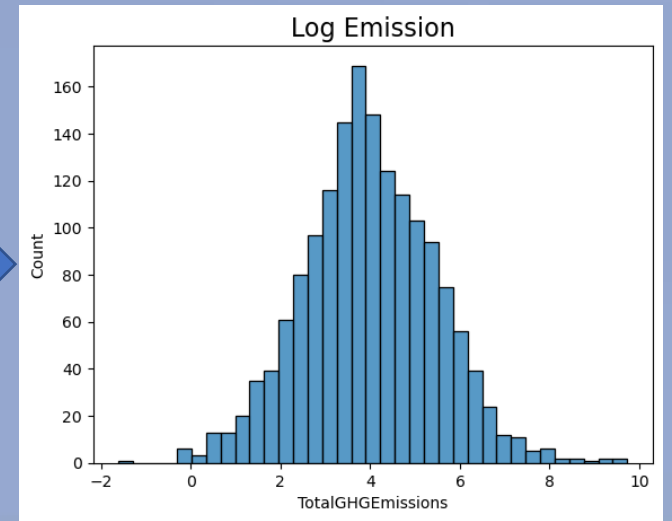
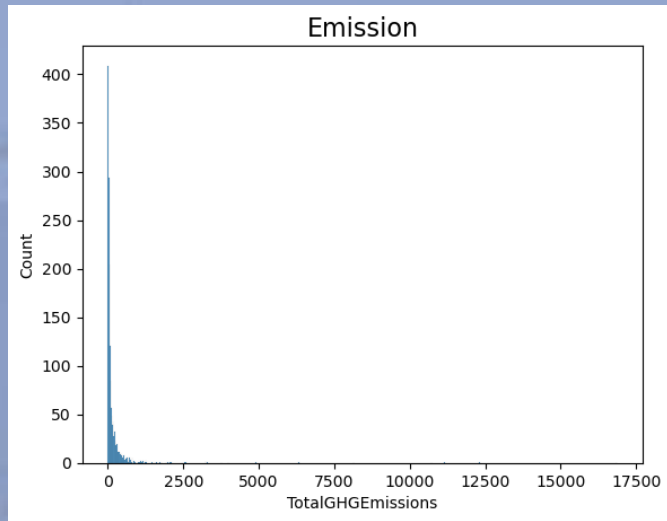
# Méthodologie





# Pré-traitement des variables

# Transformation log pour les valeurs cibles

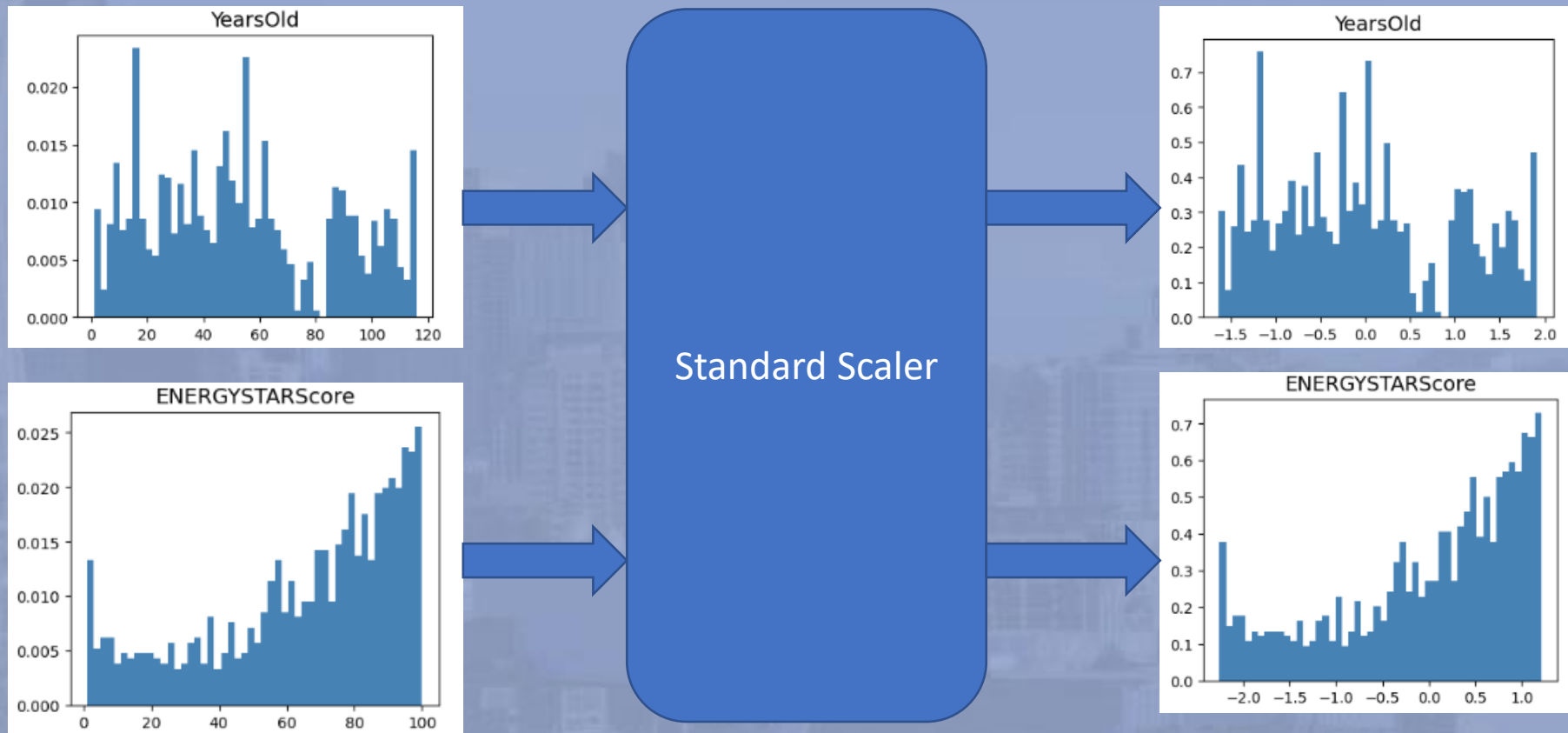


Transformation  
Logarithme

# Normalisation et encodage

- Normalisation avec Standard Scaler

Permet de réaliser une mise à l'échelle afin que les données soient centrées à 0 avec un écart type de 1



- Encodage avec One Hot Encoder

PrimaryPropertyType
Warehouse
Other
Hotel
Large Office
.....

One Hot Encoder

Warehouse	Other	Hotel	Large Office	....
1	0	0	0	
0	1	0	0	
0	0	1	0	
0	0	0	1	

# MSE , RMSE et $R^2$

- MSE :

Moyenne des distances euclidiennes entre les valeurs réelles et les valeurs prédites.

- RMSE :

Racine carrée de la moyenne des distances euclidiennes entre les valeurs réelles et les valeurs prédites.

L'erreur moyenne sur le log de la valeur cible

- $R^2$  (Coefficient de Corrélation):

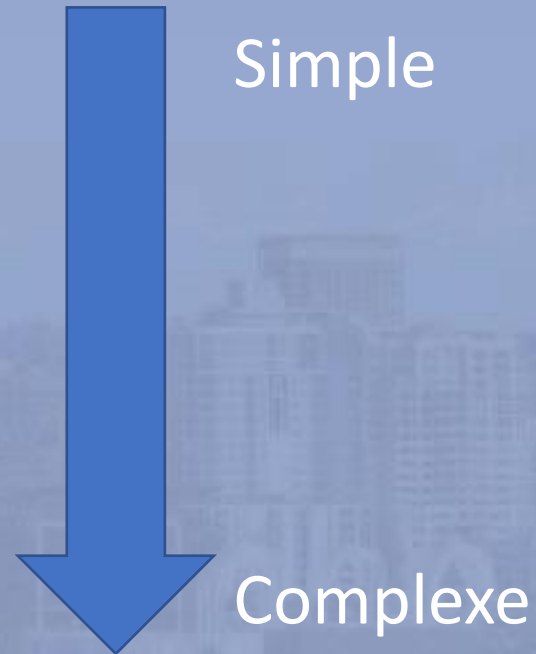
somme des distances euclidiennes entre valeurs prédites et valeurs réelles divisé par la somme des distances euclidiennes entre valeurs réelles et moyenne.

La capacité de la prédiction à expliquer la variance réelle (sur le log).

# Test des différents modèles

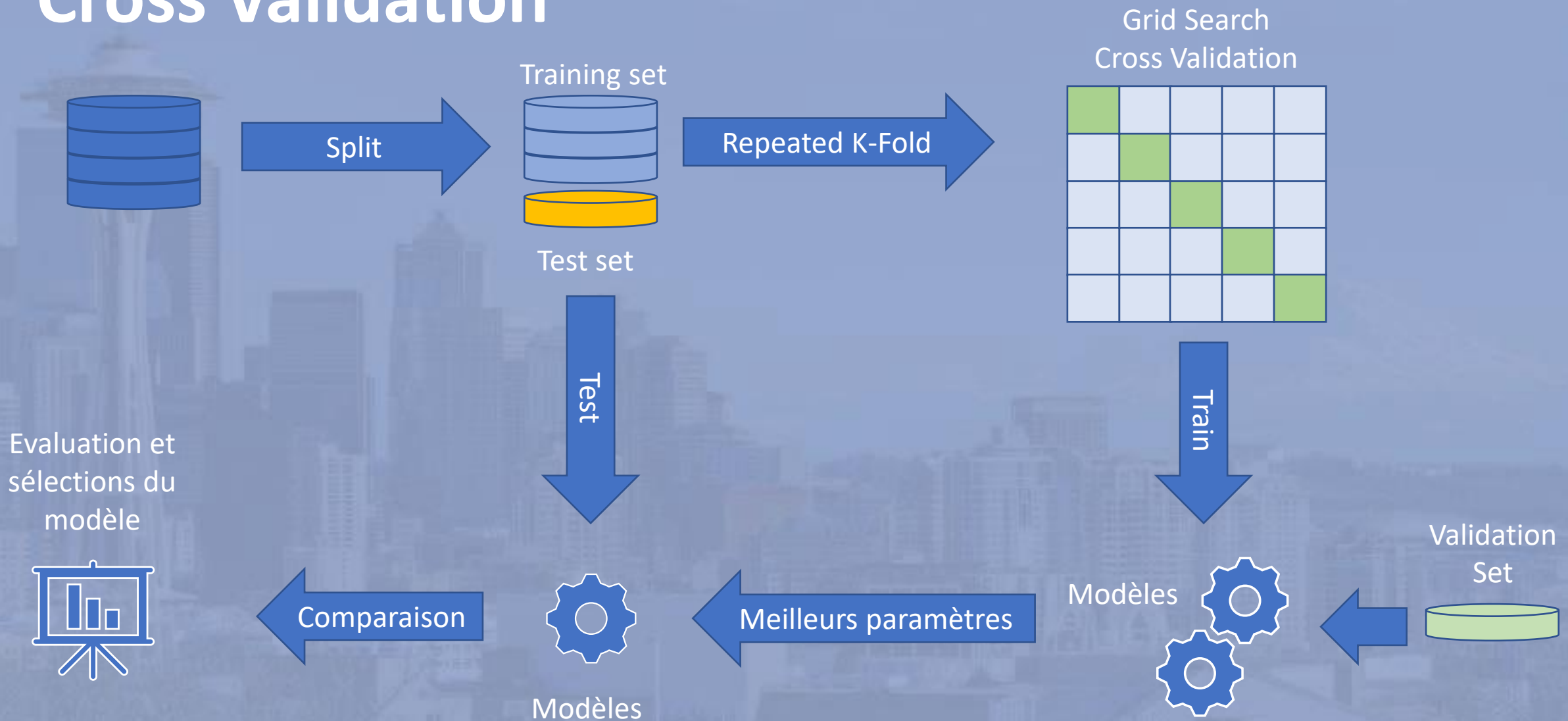
- Les modèles qui vont être testés:

- Régression Ridge
- Régression Lasso
- Régression Elastic Net
- kNN Regressor
- Support Vector Regressor
- XGBoost Regressor
- LightGBM Regressor
- Random Forest Regressor

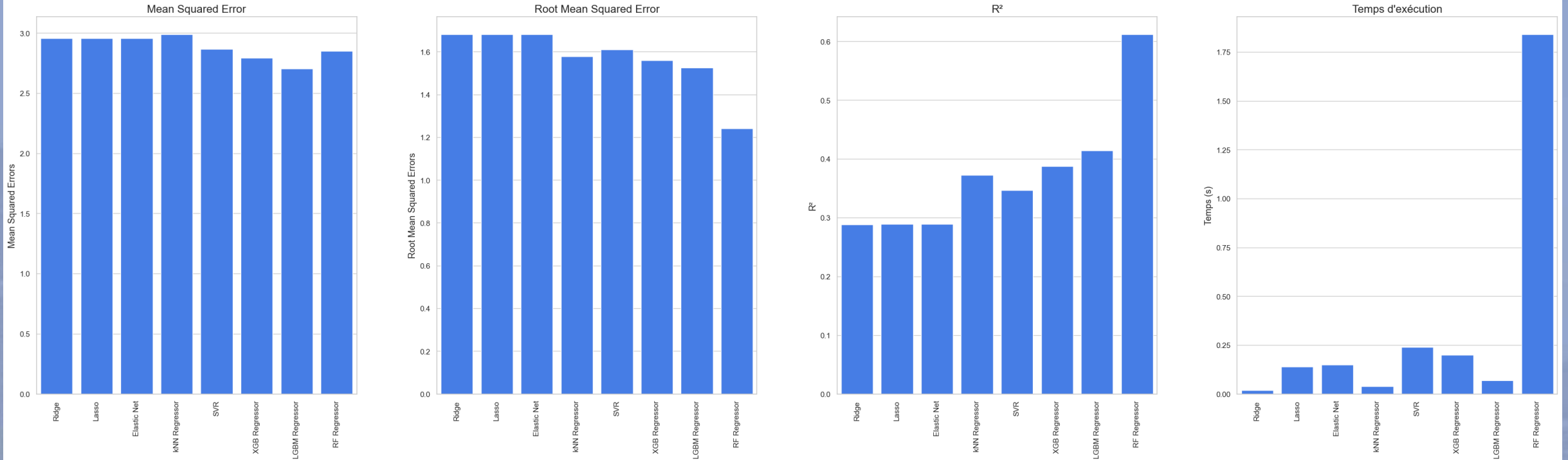




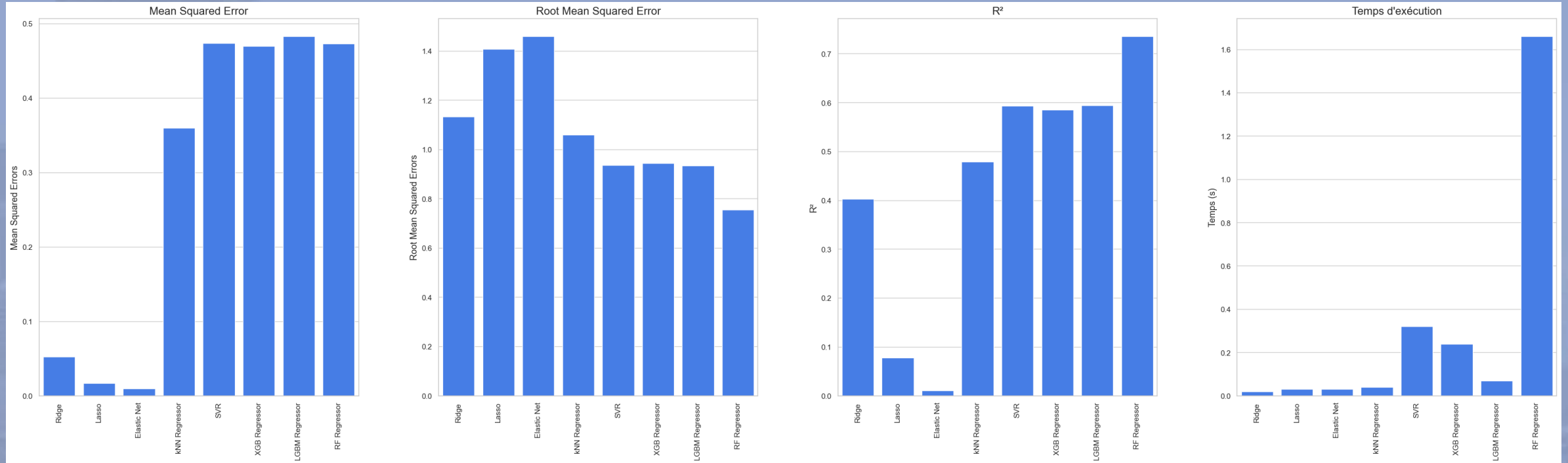
# Cross Validation



# Comparaison des modèles pour la consommation d'énergie



# Comparaison des modèles pour l'émission de CO2



## Consommation d'énergie

	R <sup>2</sup> Train	R <sup>2</sup> Test
<b>Ridge</b>	0.303815	-0.728460
<b>Lasso</b>	0.304169	-0.700649
<b>Elastic Net</b>	0.303641	-0.735701
<b>kNN Regressor</b>	0.422469	0.362411
<b>SVR</b>	0.360211	0.322745
<b>XGB Regressor</b>	0.445751	0.373037
<b>LGBM Regressor</b>	0.405397	0.354933
<b>RF Regressor</b>	0.592719	0.297099

- Problème d'overfitting pour certains modèles
- Comportement sur le test set

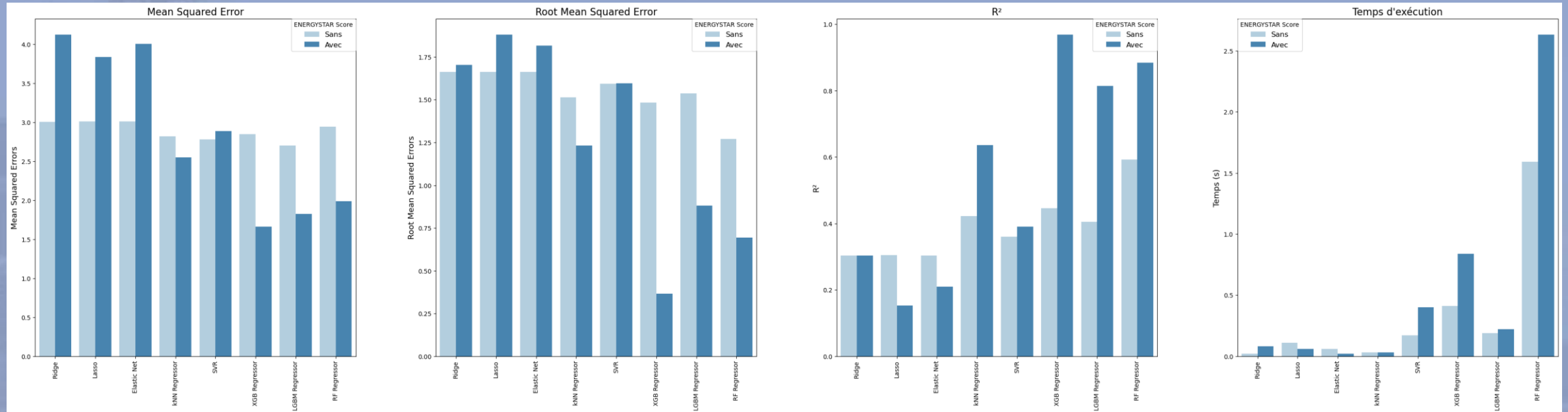
## Emission CO2

	R <sup>2</sup> Train	R <sup>2</sup> Test
<b>Ridge</b>	0.402866	0.361437
<b>Lasso</b>	0.078094	0.061417
<b>Elastic Net</b>	0.010933	0.007445
<b>kNN Regressor</b>	0.479059	0.412488
<b>SVR</b>	0.593057	0.524511
<b>XGB Regressor</b>	0.585509	0.538961
<b>LGBM Regressor</b>	0.594466	0.508146
<b>RF Regressor</b>	0.734652	0.476599

- Pas d'overfitting
- Comportement sur le test set

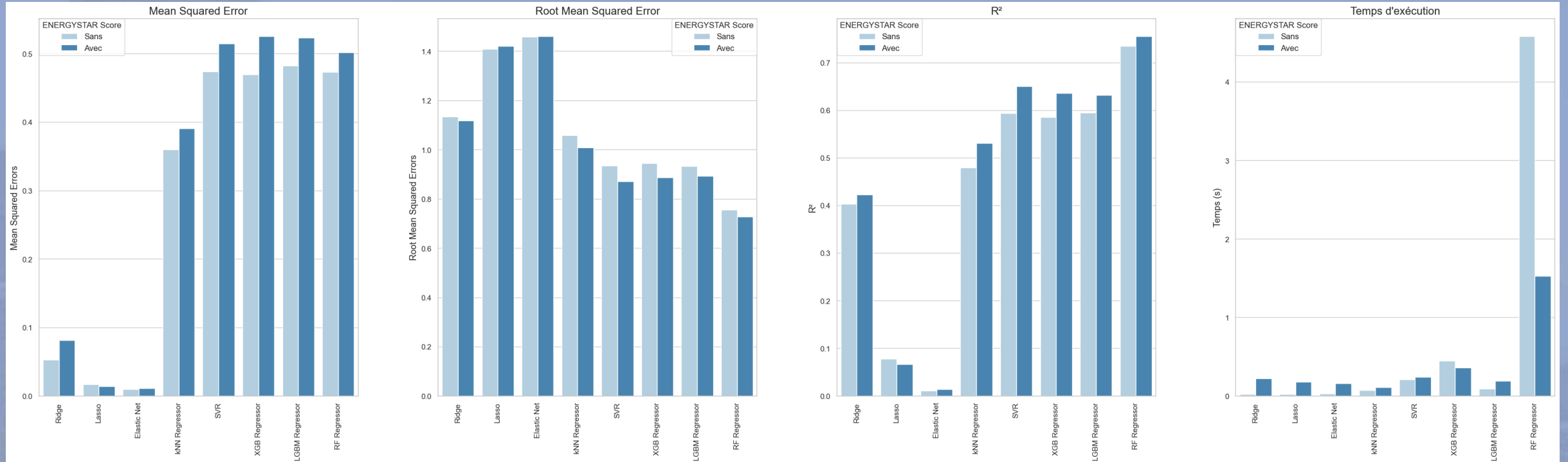
# Nouveaux tests avec l'ENERGYSTAR Score

# Consommation énergie





# Emission CO2



## Consommation d'énergie

	<b>R<sup>2</sup> Train</b>	<b>R<sup>2</sup> Test</b>
<b>Ridge</b>	0.303395	0.244833
<b>Lasso</b>	0.153382	0.141042
<b>Elastic Net</b>	0.209925	0.172316
<b>kNN Regressor</b>	0.635795	0.658304
<b>SVR</b>	0.390350	0.410527
<b>XGB Regressor</b>	0.967926	0.823031
<b>LGBM Regressor</b>	0.814276	0.761886
<b>RF Regressor</b>	0.884349	0.847370

- Pas d'overfitting
- Comportement sur le test set

## Emission CO2

	<b>R<sup>2</sup> Train</b>	<b>R<sup>2</sup> Test</b>
<b>Ridge</b>	0.423086	0.372486
<b>Lasso</b>	0.066071	0.046992
<b>Elastic Net</b>	0.014402	0.008995
<b>kNN Regressor</b>	0.530612	0.403665
<b>SVR</b>	0.649610	0.535243
<b>XGB Regressor</b>	0.635468	0.557803
<b>LGBM Regressor</b>	0.631510	0.533800
<b>RF Regressor</b>	0.754852	0.533090

- Pas d'overfitting
- Comportement sur le test set

# Choix du modèle

Pour la prédiction de la consommation d'énergie:  
Random Forest Regressor

Pour la prédiction de l'émission de CO<sub>2</sub>:  
XGB Regressor

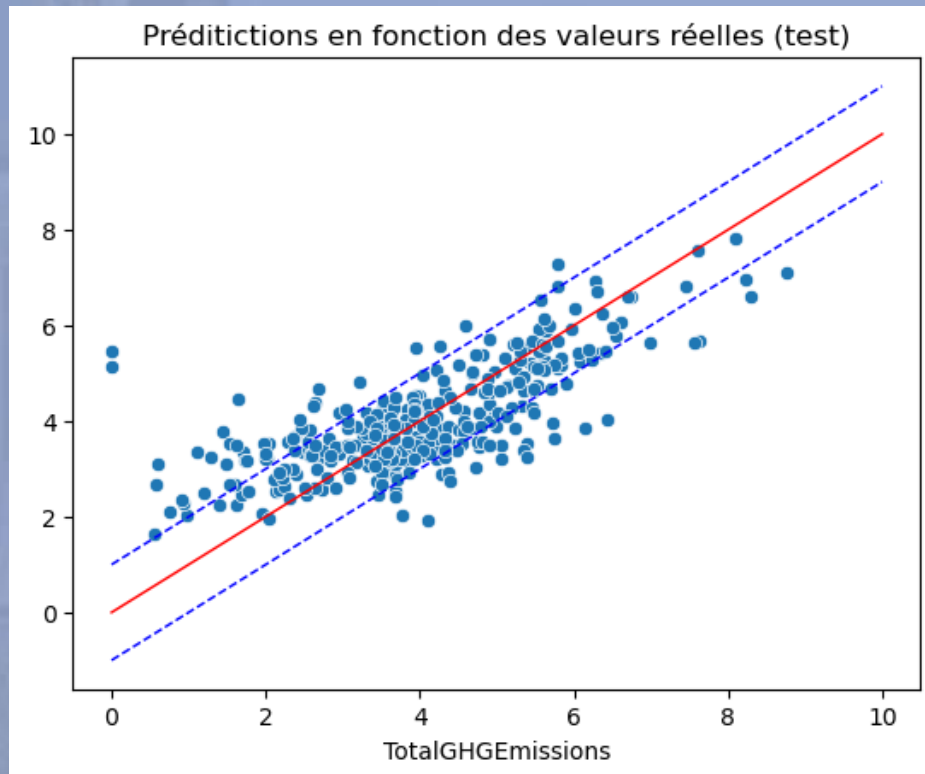
# Validation et Prédictions

XGB Regressor

$$R^2 = 0,56$$

$$\text{MSE} = 1,01$$

$$\text{RMSE} = 1$$

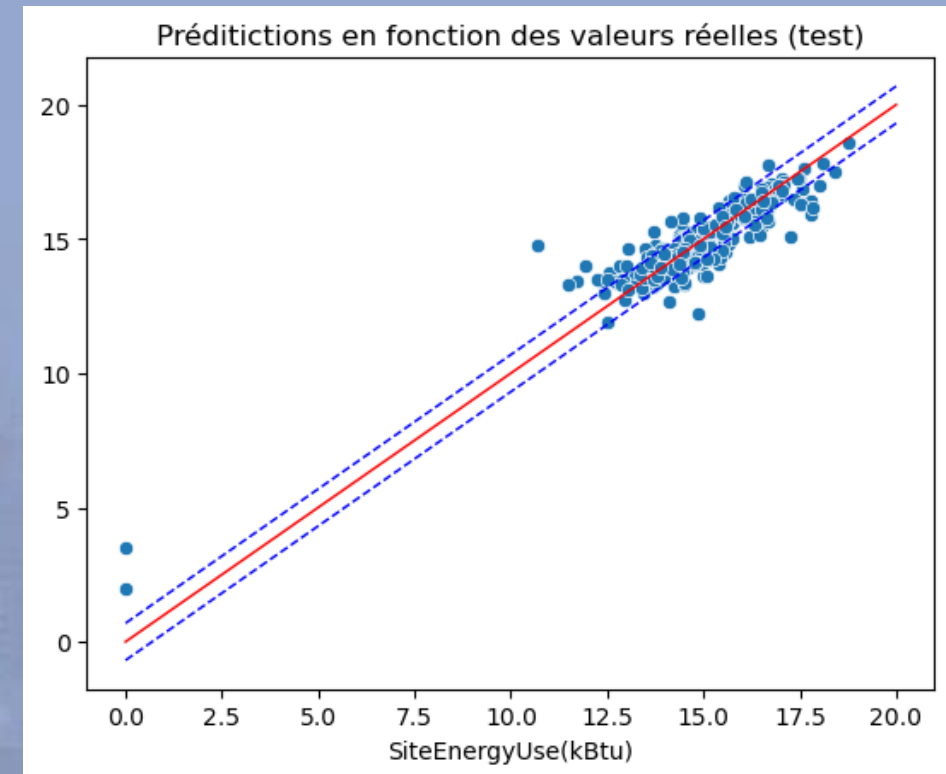


Random Forest Regressor

$$R^2 = 0,84$$

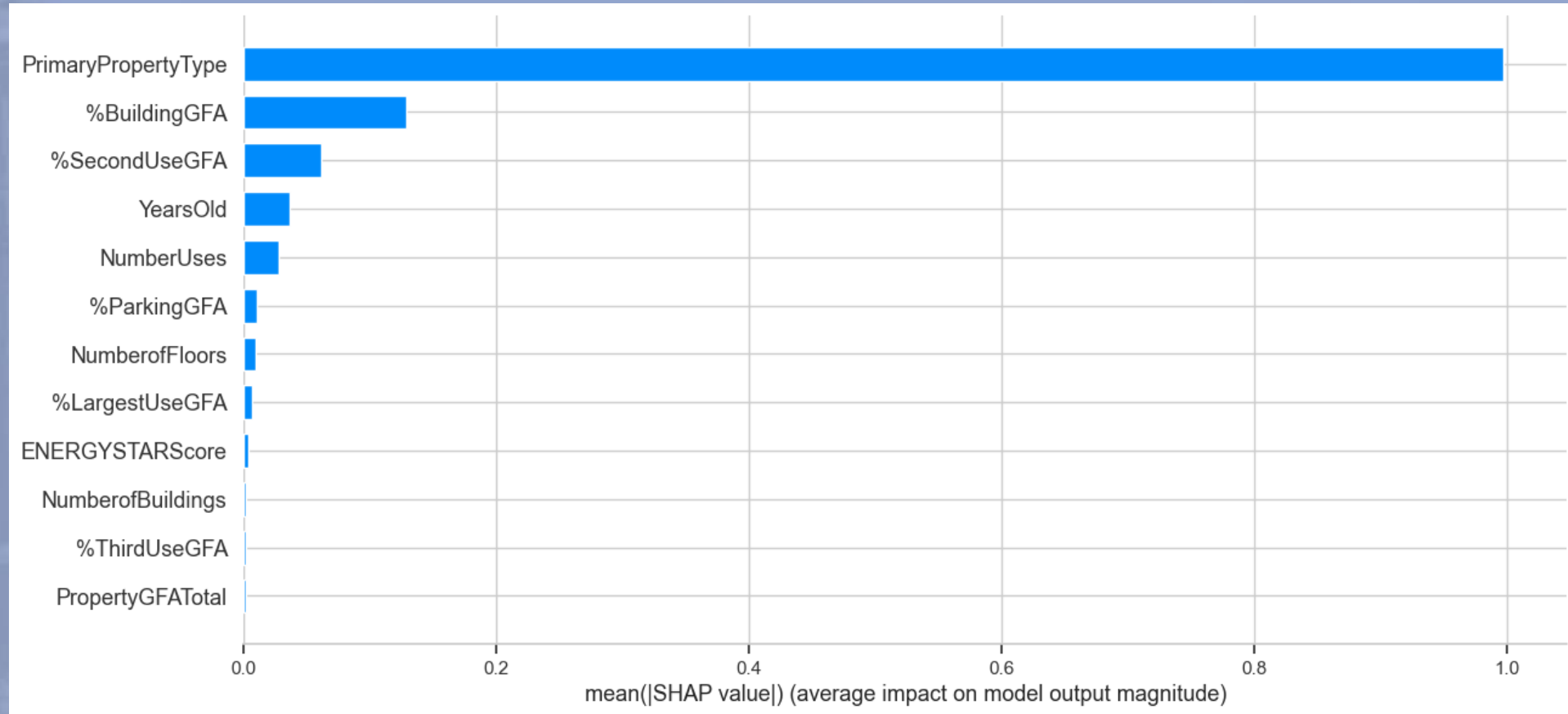
$$\text{MSE} = 0,48$$

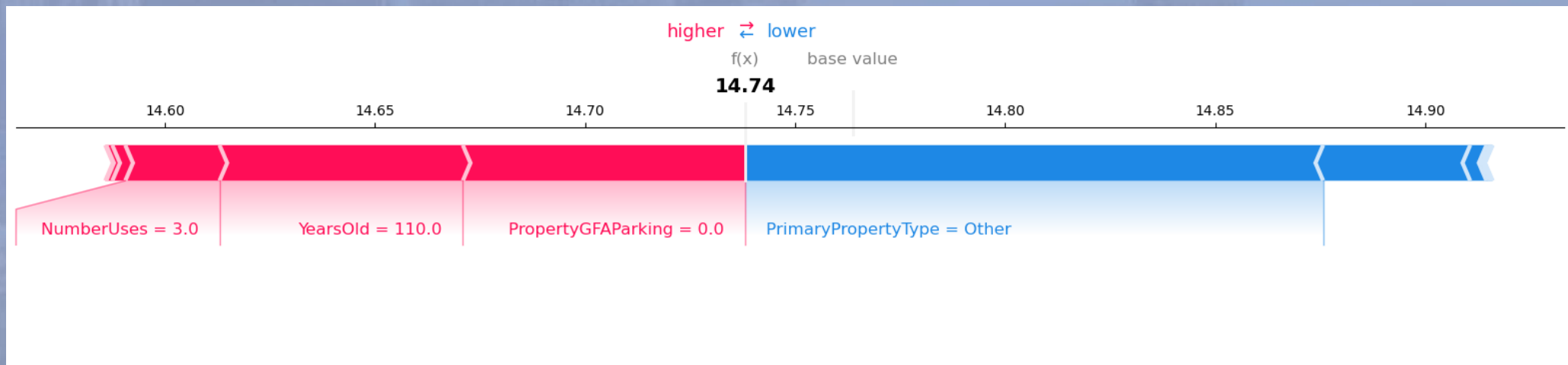
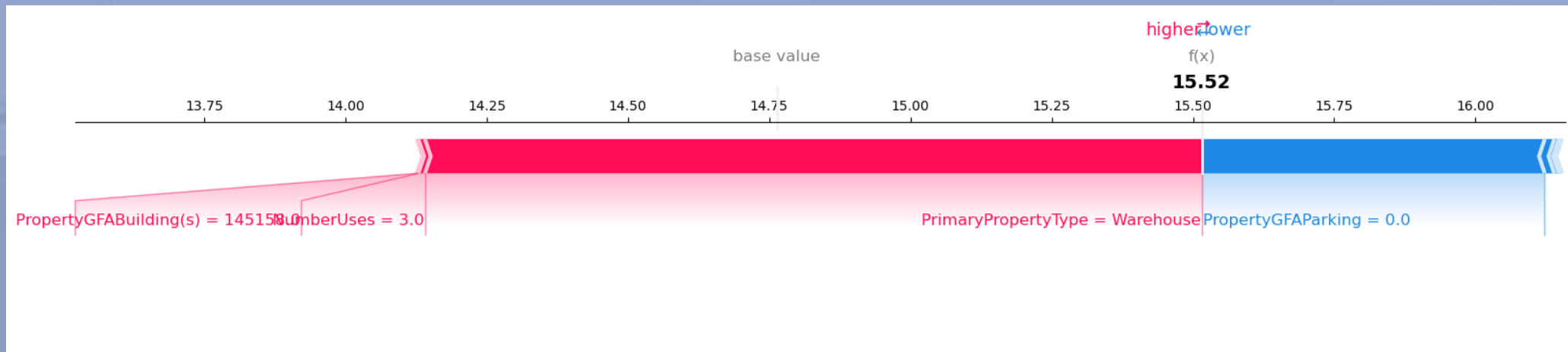
$$\text{RMSE} = 0,69$$



# FEATURES IMPORTANCES

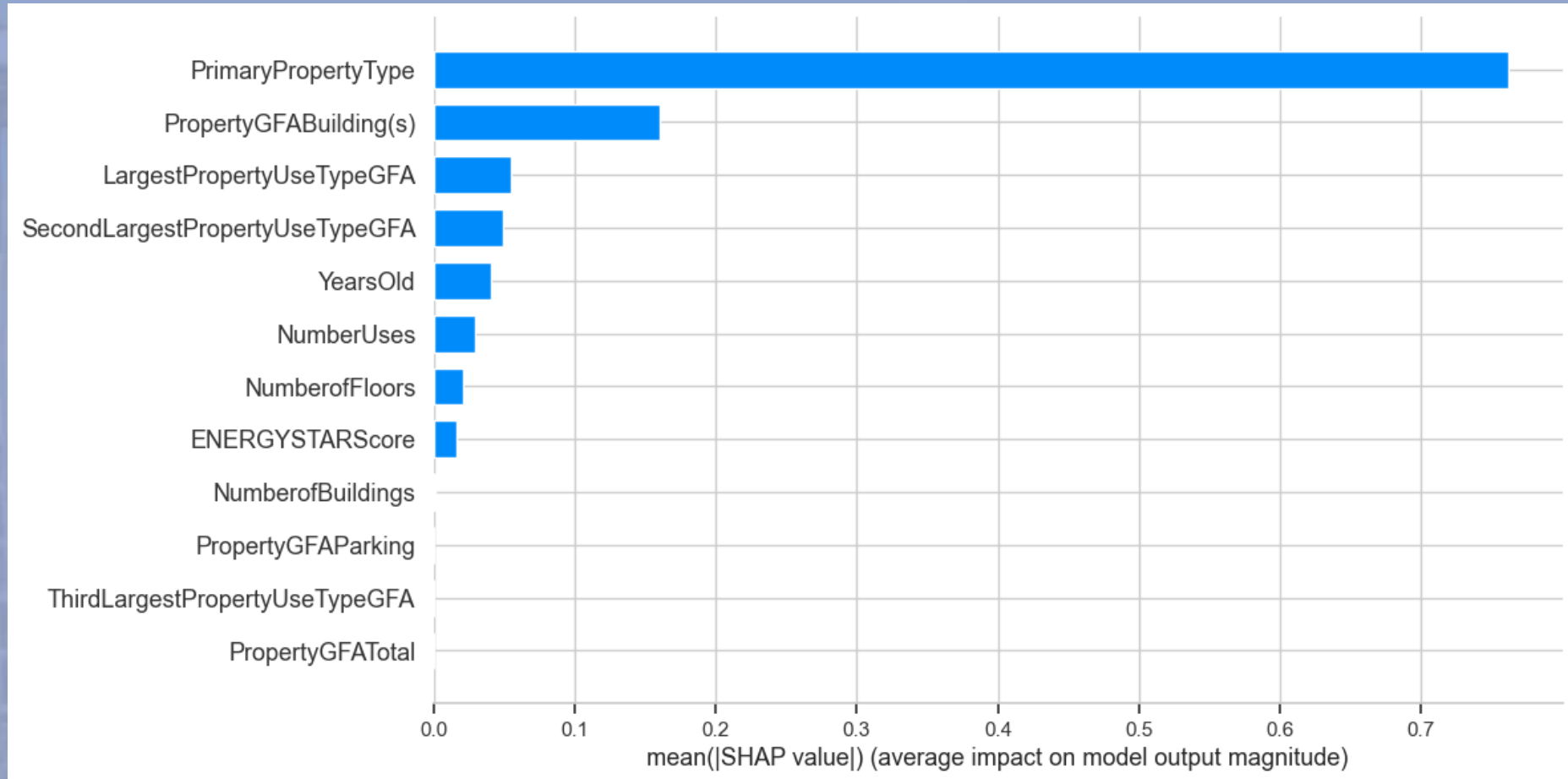
## Consommation énergie

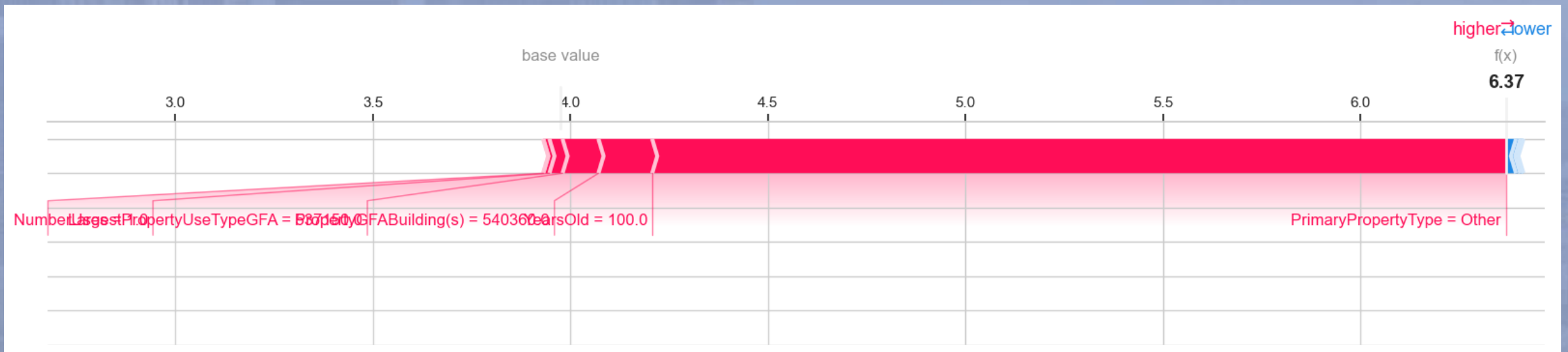
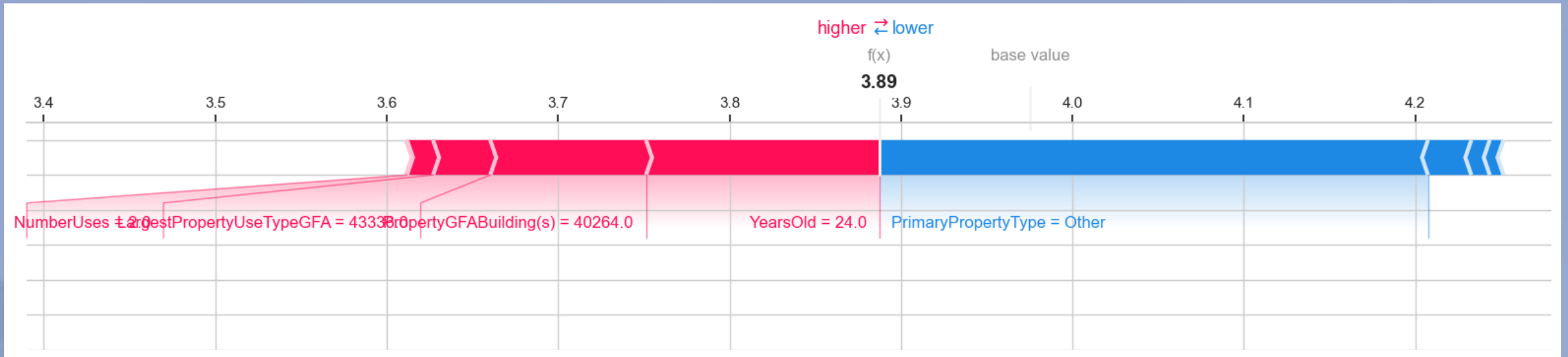






# Emission CO2





# Conclusion

- Prédiction satisfaisante mais lente pour la consommation d'énergie
- Prédiction décevante concernant l'émissions CO2
- L'ENERGYSTAR Score améliore les prédictions

MERCI DE VOTRE ATTENTION