

---

# anticipez la consommation électrique de bâtiments

Soutenance Olivier Legrand  
Parcours Data Scientist  
Projet P4

---

---

# Interprétation de la problématique

- Jeu de données: Seattle energy benchmarking, pour les années 2015 et 2016
  - Identification des cibles
    - Emissions GES: TotalGHGEmissions
    - Consommation d'énergie totale: SiteEnergyUse
  - Prédiction de la consommation totale d'énergie et d'émissions de GES à partir des caractéristiques de bâtiments:
    - Ground floor area
    - number of buildings
    - building type
    - Largest property use type
    - etc.
  - Evaluation de l'ENERGYSTARScore comme prédicteur des émissions de GES:
    - On cherchera à évaluer si ce prédicteur est fortement associé aux émissions de GES
    - On cherchera à évaluer l'impact de cet indicateur dans la qualité des prédictions.
-

---

# Pistes

- Plusieurs problématiques associées au jeu de données:
    - potentielle fuite de données → On s'empêchera d'utiliser toutes les variables "dérivées" (Intensity).
    - pas de données issues des relevés annuels, mais possibilité d'utiliser les nature et proportion d'énergie utilisées → on devra créer de nouvelles variables, mais ne pas utiliser Electricity, NaturalGas, SteamUse
  - Comment les variables liées au permis d'exploitation commerciale sont-elles associées aux grandeurs cibles? Type et importance des corrélations
  - Exploiter également les associations entre les variables catégorielles et les variables cibles.
  - Envisager d'utiliser une prédiction sur une des cibles pour prédire l'autre.
-

# Analyse exploratoire

1. Nettoyage
2. Analyse exploratoire
3. Feature Engineering

---

# Nettoyage (1)

## 1. Fusion des deux tables:

- 1.1. Transformation de la colonne Location en 6 colonnes: Address, ZipCode, Latitude, Longitude, State, City

## 2. Sélection des colonnes pertinentes:

- 2.1. Variables du permis d'exploitation: GFA (PropertyGFAs, LargestPropertyUseTypeGFAs) NumberofFloors/Building, PrimaryPropertyType, LargestPropertyType (and Second-, Third-), DataYear, YearBuilt + Electricity, NaturalGas, Steam, OtherFuelUse

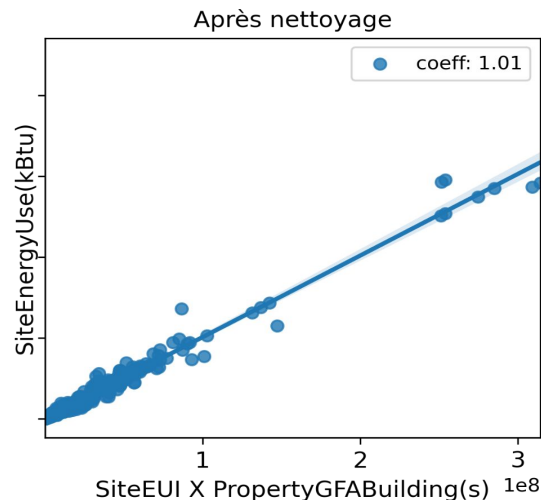
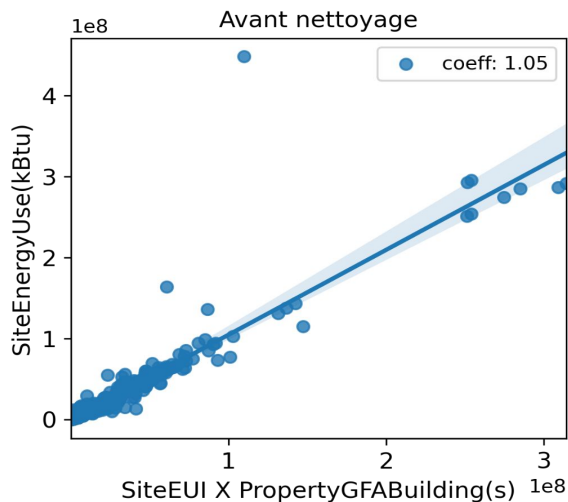
## 3. Traitement des valeurs manquantes:

- 3.1. SecondLargestPropertyUseType: on remplace par "None", car l'absence de valeur est cohérente, mais la présence de NaN peut empêcher certains traitements numériques. Idem pour ThirdLargestPropertyUseType.
  - 3.2. Pour les variables quantitatives, on supprime les lignes incomplètement renseignées - sauf pour ENERGYSTARScore: trop de valeurs manquantes.
-

# Nettoyage (2)

## Traitement des outliers

1. Utilisation de la colonne 'Outliers'
2. Sélection des individus pour lesquels SEU, GHG, NumberofBuildings > 0; 0 < NumberofFloors < 80
3. Correction des valeurs négatives de PropertyGFABuilding(s)
4. Utilisation de la relation linéaire entre SEU et SEUIntensity



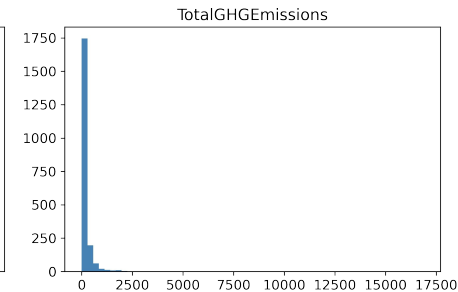
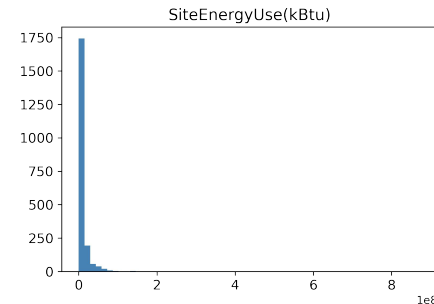
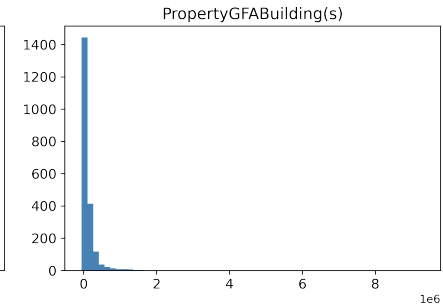
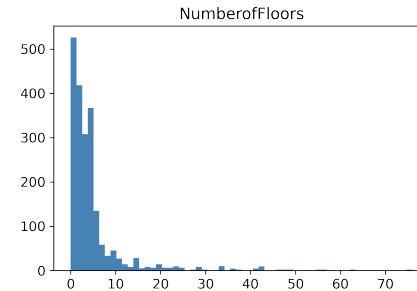
# Analyse exploratoire

## Distributions des variables quantitatives

- Non-gaussiennes
- Grandes dispersions
- Grandes différences d'échelles

## Corrélations importantes pour certains groupes de variables

- 98% pour PropertyGFATotal, PropertyGFABuilding(s)
- 97% pour PropertyGFATotal, LargestPropertyUseTypeGFA
- 78% pour PropertyGFATotal, SecondLargestPropertyUseTypeGFA
- 94% pour Electricity et SiteEnergyUse, -92% pour Electricity et NaturalGas
- 89 % pour SiteEnergyUse et TotalGHGEmissions



# Analyse exploratoire

## Cas des variables catégorielles

- Grand nombre de modalités pour chaque variable (DataYear: 2 modalités, mais LargestPropertyUseType:53 et YearBuilt: 112 par exemple)
- Grand nombre de modalités presque vides: source potentielle de bruit
- Corrélations: ANOVA

	$\eta^2$	
	SiteEnergyUse	TotalGHGEmissions
PrimaryPropertyType	0.077	0.056
LargestPropertyUseType	0.060	0.067
SecondLargestPropertyUseType	0.042	0.019
BuildingType	0.023	0.013
YearBuilt	0.023	0.007



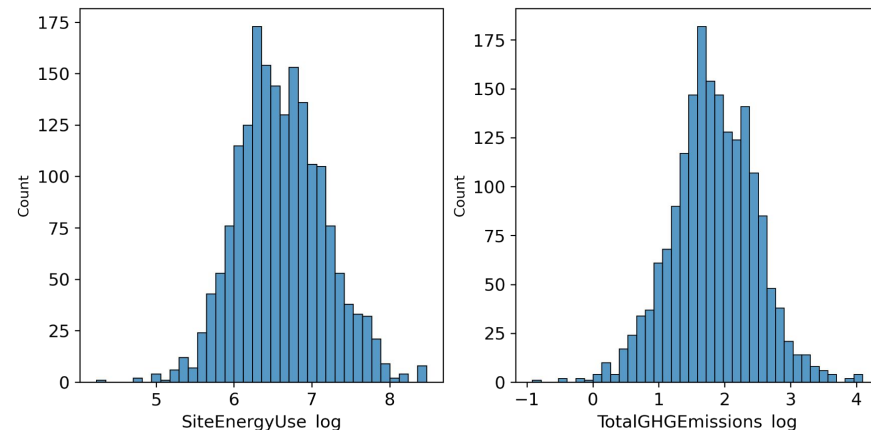
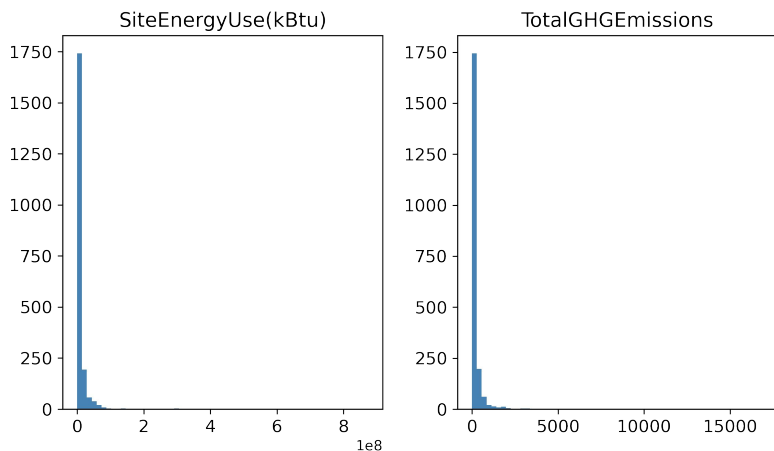
# Analyse exploratoire

## Corrélations et associations

- Intéressant d'envisager une PCA sur les variables quantitatives corrélées aux cibles, et très corrélées entre elles
  - ANOVA → mise à l'écart des v. catégorielles non associées, ou dont l'effet est très faible
  - ENERGYSTARScore très peu corrélée aux cibles et aux autres variables
  - TotalGHGEmissions et SiteEnergyUse très corrélées entre elles, SiteEnergyUse plus corrélée aux autres variables que TotalGHGEmissions: modèle séquentiel vs modèle non-séquentiel
-

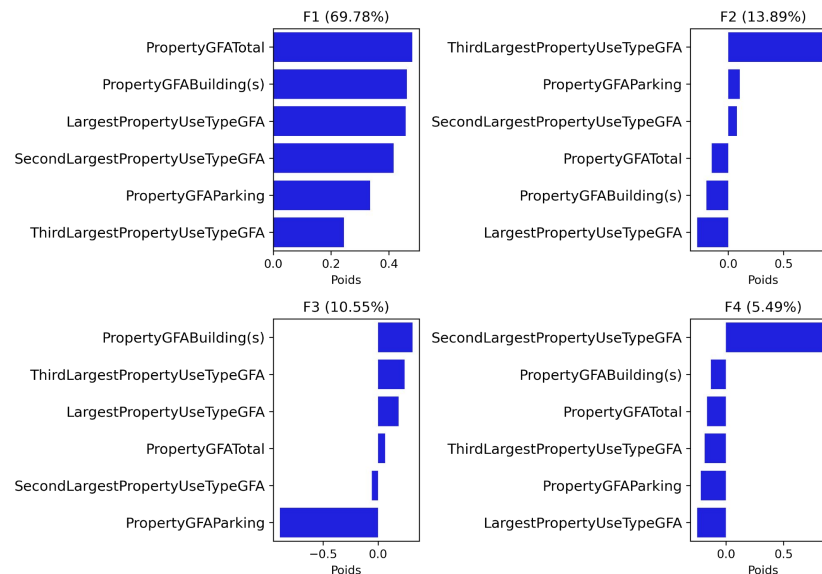
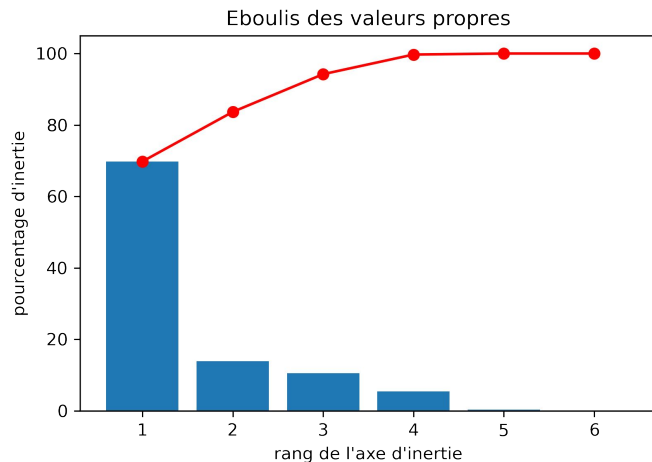
# Feature engineering

1. Très grande dispersion et non-gaussianité des données:
  - a.  $\log_{10}$  sur les cibles TotalGHGEmissions et SiteEnergyUse
  - b. standardisation des prédicteurs (sur jeu d'entraînement seulement)



# Feature engineering

1. Très grande dispersion et non-gaussianité des données:
  - a.  $\log_{10}$  sur les cibles TotalGHGEmissions et SiteEnergyUse
  - b. standardisation des prédicteurs (sur jeu d'entraînement seulement)
2. PCA sur le groupe de variables 'GFA' très corrélées entre elles



# Feature engineering

1. Très grande dispersion et non-gaussianité des données:
  - a.  $\log_{10}$  sur les cibles TotalGHGEmissions et SiteEnergyUse
  - b. standardisation des prédicteurs (sur jeu d'entraînement seulement)
2. PCA sur le groupe de variables 'GFA' très corrélées entre elles
3. Création des variables Energy\_ratio, NaturalGas\_ratio, Steam\_ratio

	Electricity(kBtu)	NaturalGas(kBtu)	SteamUse(kBtu)	OtherFuelUse(kBtu)	Electricity_ratio	NaturalGas_ratio	Steam_ratio	OtherFuel_ratio
0	3686160.0	1272388.0	2023032.0	0.0	0.527995	0.182253	0.289773	0.0
1	3905411.0	4448985.0	0.0	0.0	0.467477	0.532542	0.000000	0.0
2	49762435.0	3709900.0	19660404.0	0.0	0.680459	0.050730	0.268839	0.0
4	6066245.0	8763105.0	0.0	0.0	0.409077	0.590940	0.000000	0.0
5	7271004.0	4781283.0	0.0	0.0	0.603303	0.396722	0.000000	0.0

# Feature engineering

1. Très grande dispersion et non-gaussianité des données:
    - a.  $\log_{10}$  sur les cibles TotalGHGEmissions et SiteEnergyUse
    - b. standardisation des prédicteurs (sur jeu d'entraînement seulement)
  2. PCA sur le groupe de variables 'GFA' très corrélées entre elles
  3. Création des variables Energy\_ratio, NaturalGas\_ratio, Steam\_ratio
  4. One-hot encoding des variables catégorielles, puis réduction du nombre de modalités:
    - a. par groupements basés sur des seuils de population et/ou règles métiers et/ou des considérations portant sur les dépendances entre les cibles et les diverses modalités.
    - b. YearBuilt groupé en deux catégories: avant 1980, après 1980
-

# Modèles

1. Structure générale des modèles
2. Baseline: Régression linéaire
3. Régression polynomiale avec Lasso
4. K-NN
5. Random Forest
6. Comparaison des modèles
7. ENERGYSTARScore

---

# Structure des modèles

Modèle: pipeline de prétraitement + estimateur.

- Pipeline de prétraitement:
    - Standardisation des variables et PCA
    - one-hot encoding des variables catégorielles
    - optionnel: Transformation des features pour la régression polynomiale
  
  - Estimateurs:
    - Régression Linéaire: baseline
    - Régression Polynomiale: PolynomialFeatures + Lasso pour réduire la complexité du modèle polynomial,
    - KNN,
    - RandomForest pour réduire la complexité et prédire.
-

## Baseline: régression linéaire

- Jeu de données: On écarte 'DataYear' (SEU) ou 'DataYear' et ThirdLargestPropertyUseType
  - Pipeline: Standardisation, PCA, One-hot encoding
  - Méthode: validation croisée 5 folds
-



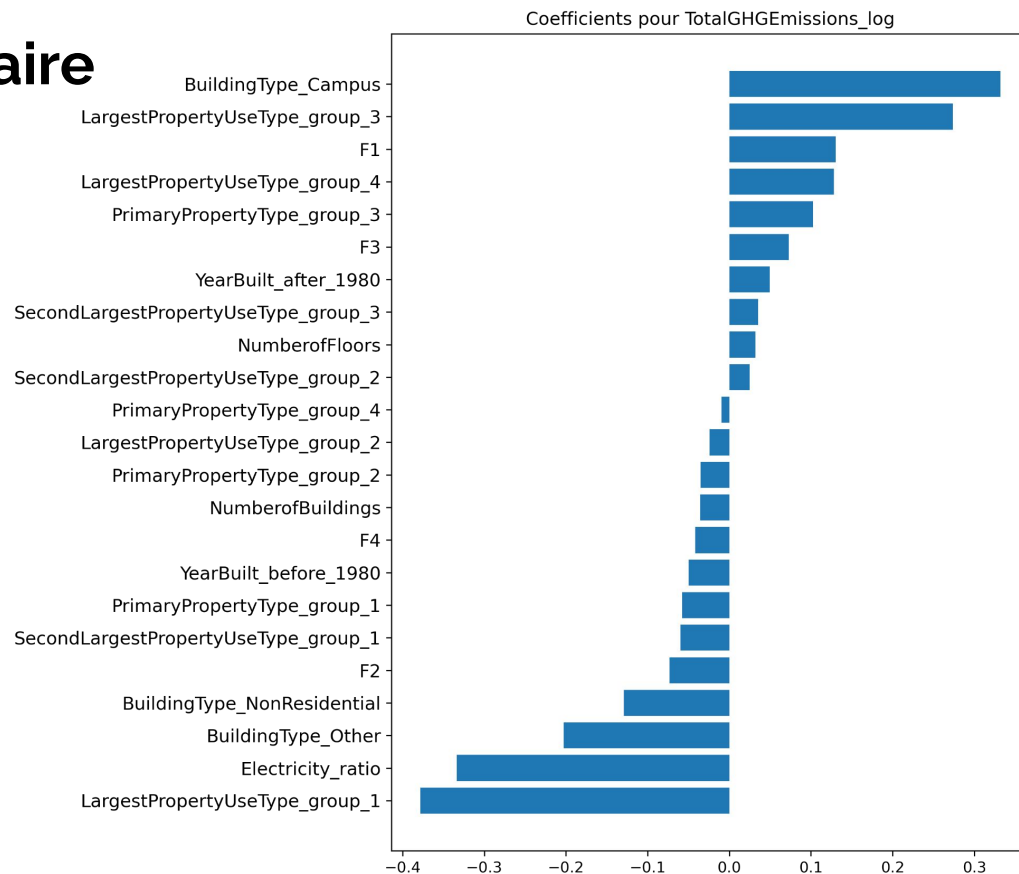
# Baseline: régression linéaire

	SiteEnergyUse	TotalGHGEmissions
R2 (entraînement)	0.59 +/- 0.01	0.65 +/- 0.1
R2 (test)	0.56 +/- 0.04	0.62 +/- 0.04

→ Le modèle semble stable, mais le score indique un possible sous-apprentissage.  
Régression polynomiale pour prendre en compte les interactions entre variables.

# Baseline: régression linéaire

Principaux prédicteurs:

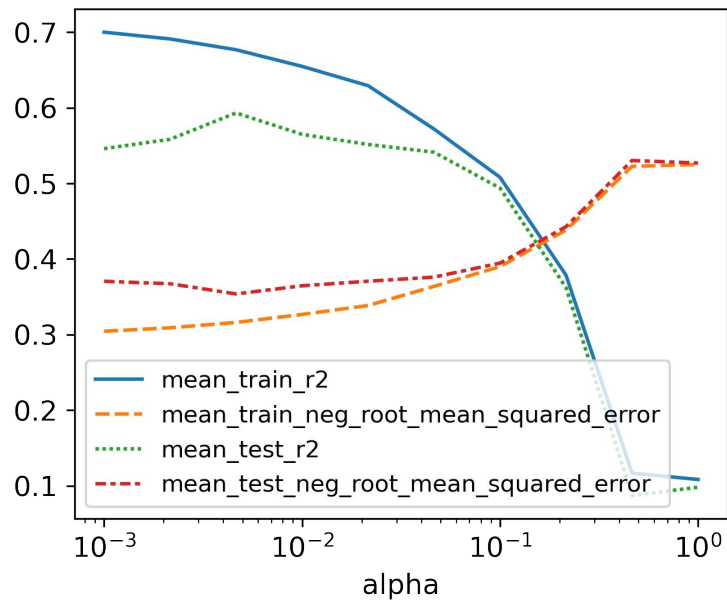


## Régression polynomiale avec lasso

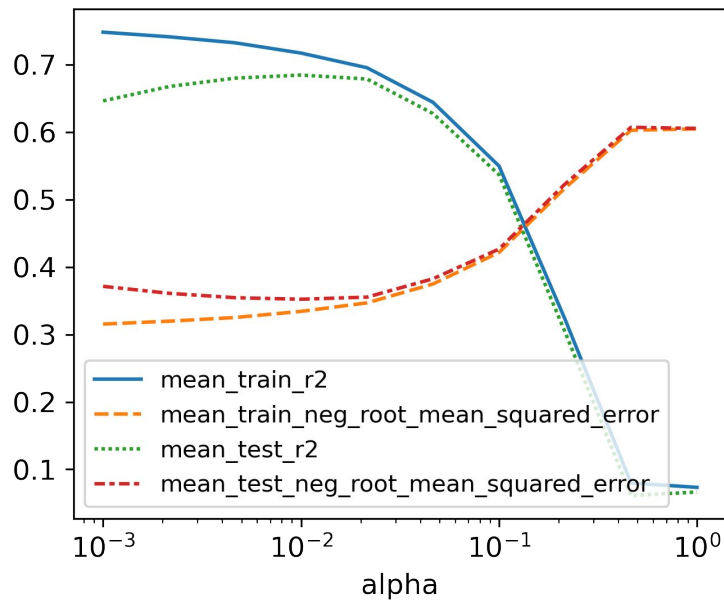
- Jeu de données: On écarte 'DataYear', 'ThirdLargestPropertyUseType', 'BuildingType' car risque de surapprentissage.
  - Modèle exploré seulement sur la cible SiteEnergyUse.
  - Pipeline: Standardisation, PCA, One-hot encoding, PolynomialFeatures de degré 2, Lasso
  - Méthode: GridSearch (5 folds) sur le set d'entraînement pour l'évaluation de alpha
-

# Régression polynomiale avec lasso

SiteEnergyUSE



TotalGHGEmissions



## Régression polynomiale avec lasso

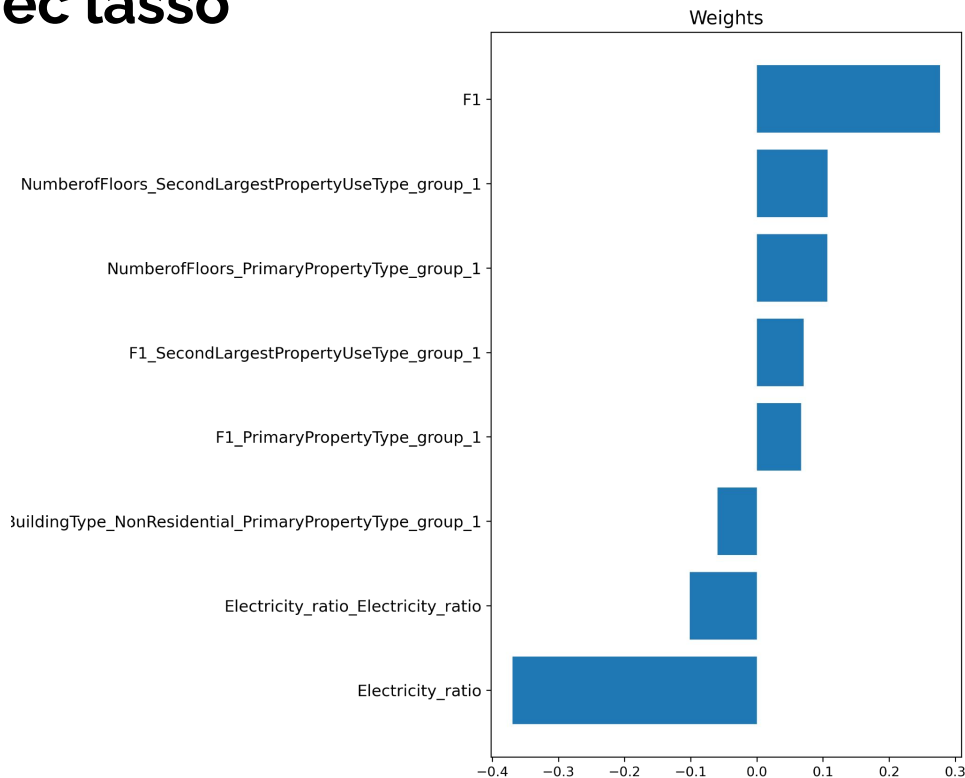
	SiteEnergyUse	TotalGHGEmissions
alpha	0.00717	0.00717
R2 (entraînement)	0.65 +/- 0.01	0.730 +/- 0.002
<b>R2 (test)</b>	<b>0.59 +/- 0.03</b>	<b>0.705 +/- 0.012</b>

→ Amélioration par rapport à la régression linéaire. Toujours en sous-apprentissage.  
Pour aller plus loin: knn, randomforest

# Régression polynomiale avec lasso

Principaux prédicteurs:

→ Prise en compte des interactions.

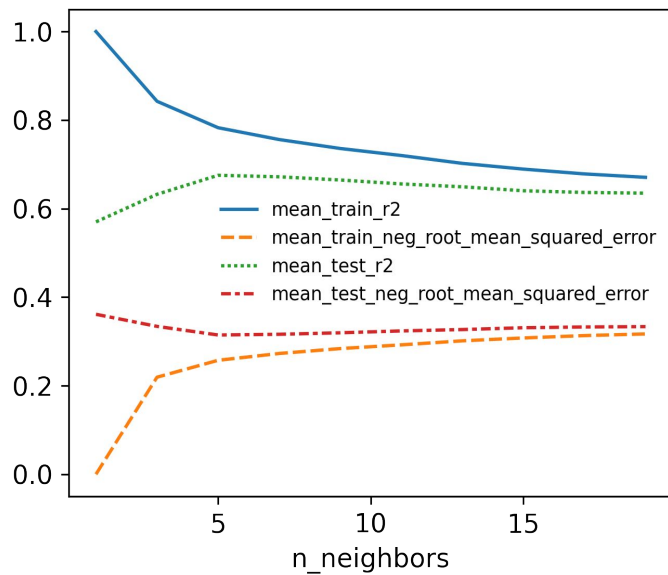


## k-NearestNeighbors

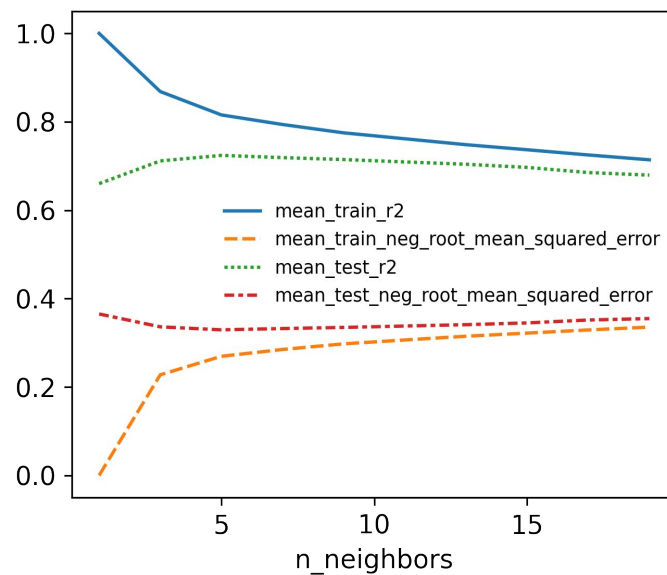
- Jeu de données: On écarte 'DataYear', 'ThirdLargestPropertyUseType', 'BuildingType' (SEU), et 'DataYear', 'ThirdLargestPropertyUseType', 'YearBuilt'
  - Pipeline: Standardisation, PCA, One-hot encoding
  - Méthode: GridSearch pour la détermination du nombre optimal de p.p. voisins
-

# k-NearestNeighbors

SiteEnergyUse



TotalGHGEmissions





## k-NearestNeighbors

	SiteEnergyUse	TotalGHGEmissions
n_neighbors	5	5
R2 (entraînement)	0.783 +/- 0.002	0.816 +/- 0.004
<b>R2 (test)</b>	<b>0.675 +/- 0.02</b>	<b>0.724 +/- 0.018</b>

# Random Forest

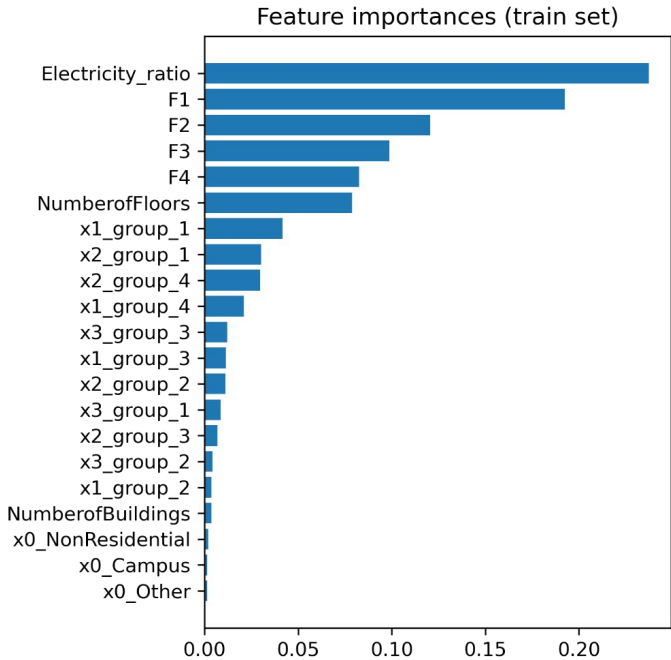
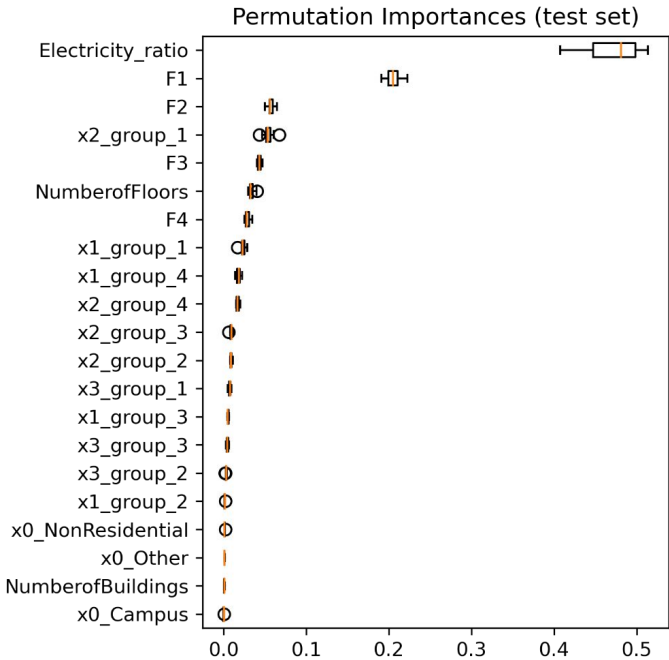
- Jeu de données: Même jeu de données que pour la régression linéaire simple
  - Pipeline: Standardisation, PCA, One-hot encoding
  - Méthode: GridSearch pour déterminer les valeurs optimales de max\_features, max\_depth et n\_estimators
  - + On utilise deux méthodes (feature permutation et feature importance) pour évaluer l'importance des différents prédicteurs et créer un modèle plus simple.
-

# Random Forest

	SiteEnergyUse	TotalGHGEmissions
n_estimators, max_depth, max_features	500, 'None', 'sqrt'	500, 'None', 'auto'
R2 (entraînement)	0.971 +/- 0.000	0.975 +/- 0.001
<b>R2 (test)</b>	<b>0.790 +/- 0.023</b>	<b>0.816 +/- 0.027</b>

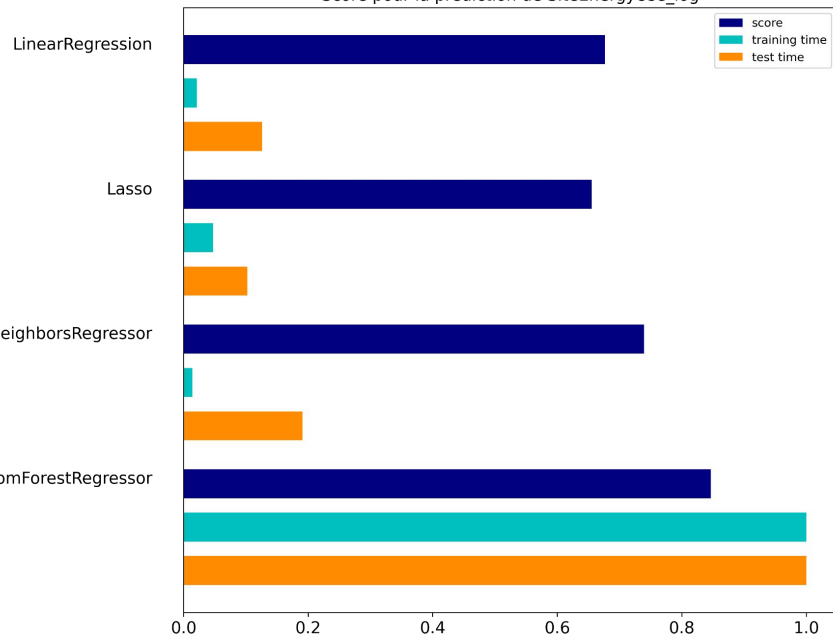
# Feature importance avec Random Forest

x0: BuildingType  
x1: PrimaryPropertyType  
x2: LargestPropertyUseType  
x3: SecondLargestPropertyUseType

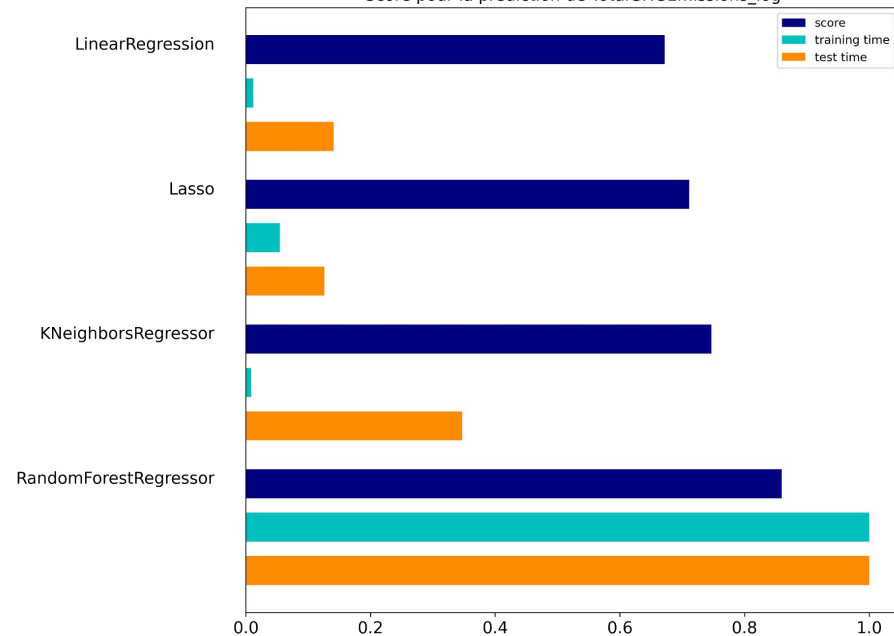


# Comparaison des modèles

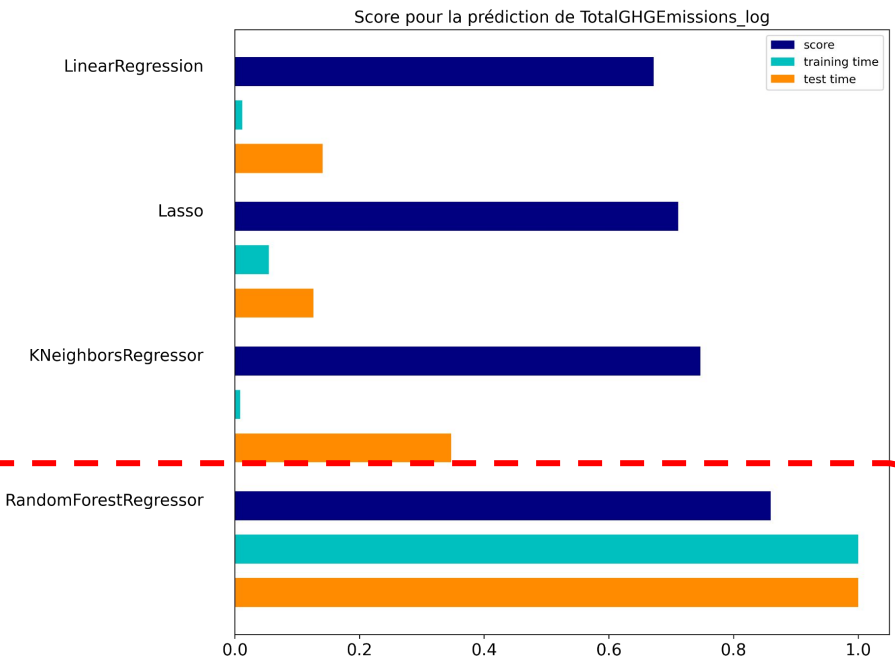
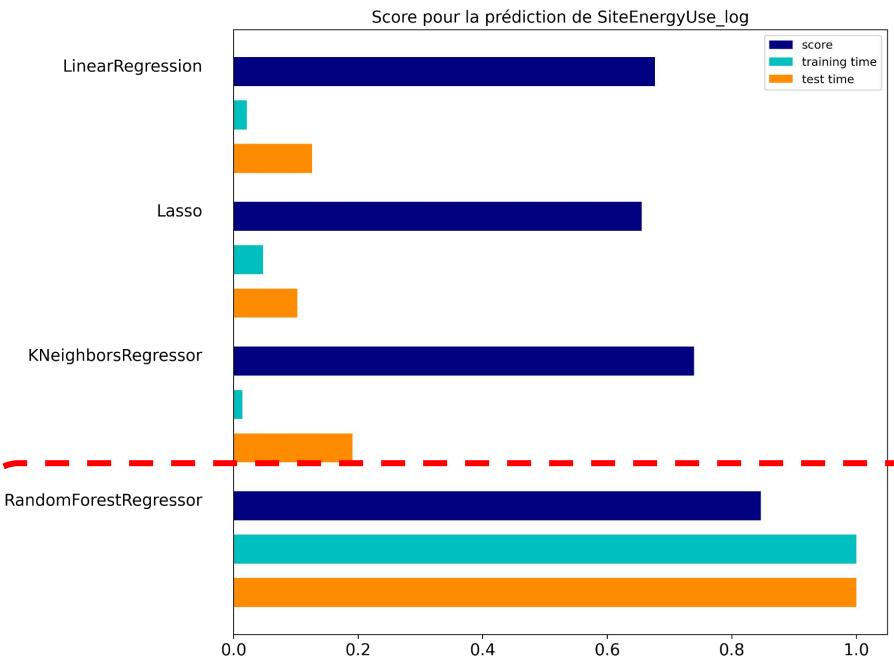
Score pour la prédiction de SiteEnergyUse\_log



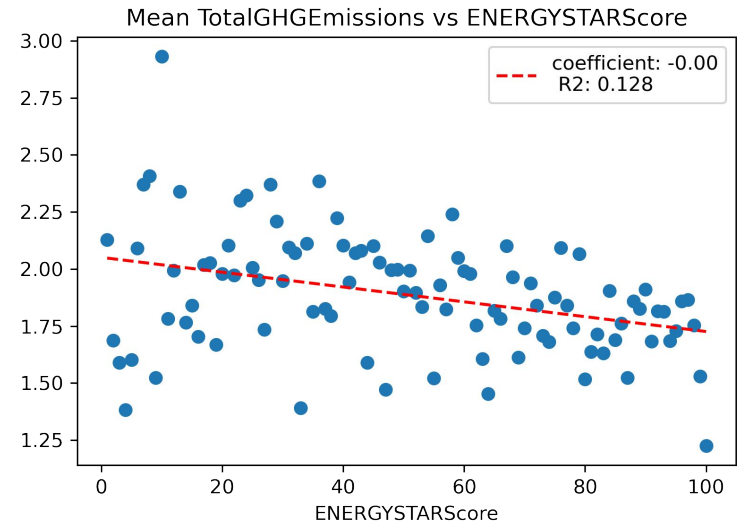
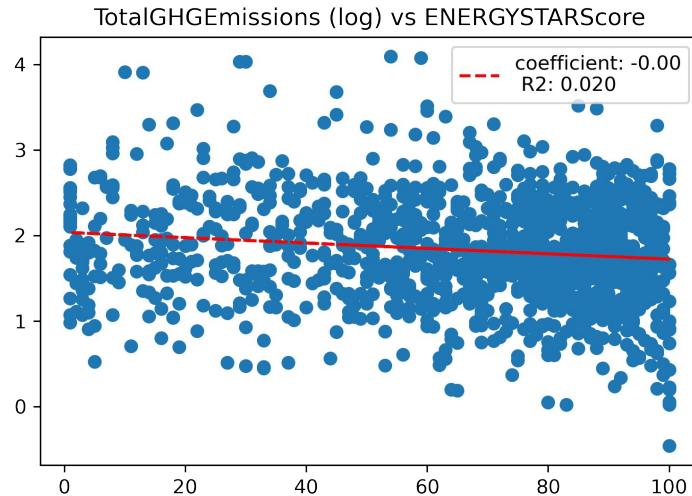
Score pour la prédiction de TotalGHGEmissions\_log



# Comparaison des modèles



# Pertinence d'ENERGYSTARScore pour la prédiction de TotalGHGEmissions



→ Corrélation présente, mais faible entre ENERGYSTARScore et TotalGHGEmissions.  
Evaluation réalisée par la comparaison de la qualité des prédictions avec et sans ESS.

## Pertinence d'ENERGYSTARScore pour la prédiction de TotalGHGEmissions

Comparaison de trois cas:

- prédiction s'appuyant sur le jeu de données complet
- prédiction s'appuyant sur le jeu de données restreint aux seuls individus pour lesquels ENERGYSTARScore est renseigné
- ENERGYSTARScore inclus dans l'ensemble des prédicteurs.

	TotalGHGEmissions
jeu de données complet	0.84 +/- 0.02
Jeu de données restreint	0.88 +/- 0.02
<b>Prise en compte d'ENERGYSTARScore</b>	<b>0.92 +/- 0.01</b>

---



# Modèle final

1. Modèle final
2. Pistes d'amélioration (1): Sélection des features
3. Pistes d'amélioration (2): modèles séquentiel vs non-séquentiel

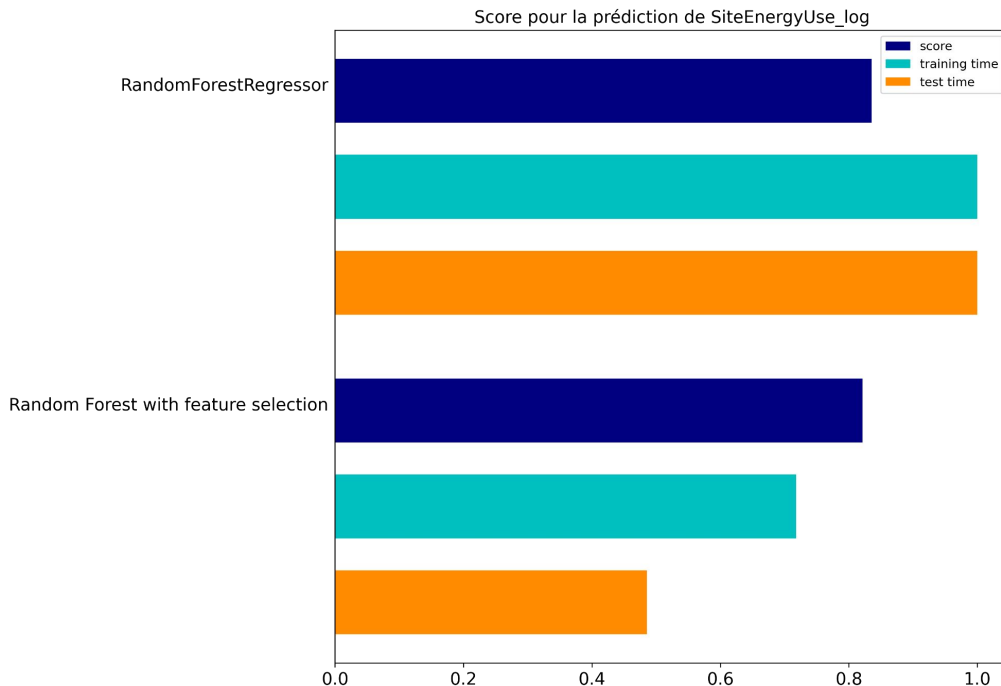
# Modèle final

1. Sélection des inputs
  2. PCA, normalisation, Onehot encoding
  3. Prédictions avec **RandomForest**
  4. Pistes d'améliorations:
    - a. Sélection des features
    - b. Modèle séquentiel vs modèle non-séquentiel
-

## Modèle final - Améliorations (1): Sélection des features

Sélection des features pour lesquelles feature importance > 2.5%

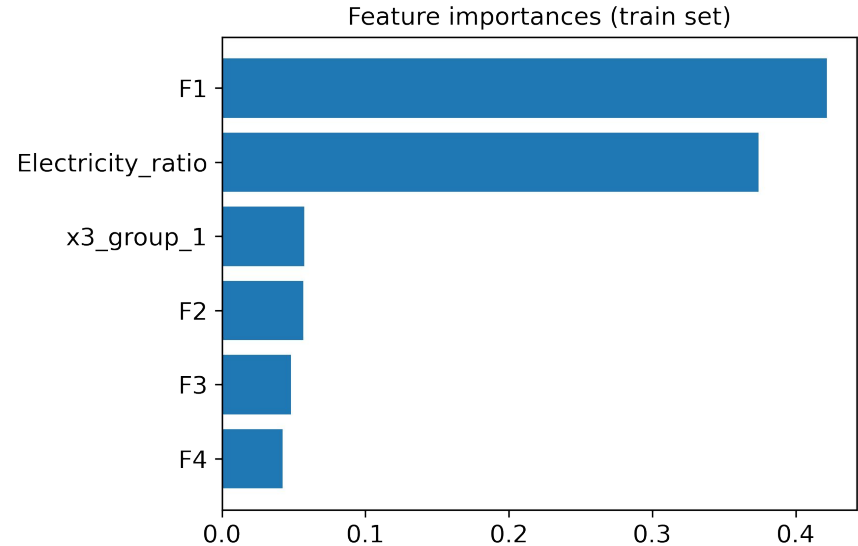
- Amélioration du temps d'entraînement
- Perte de performance maîtrisée



## Modèle final - Améliorations (1): Sélection des features

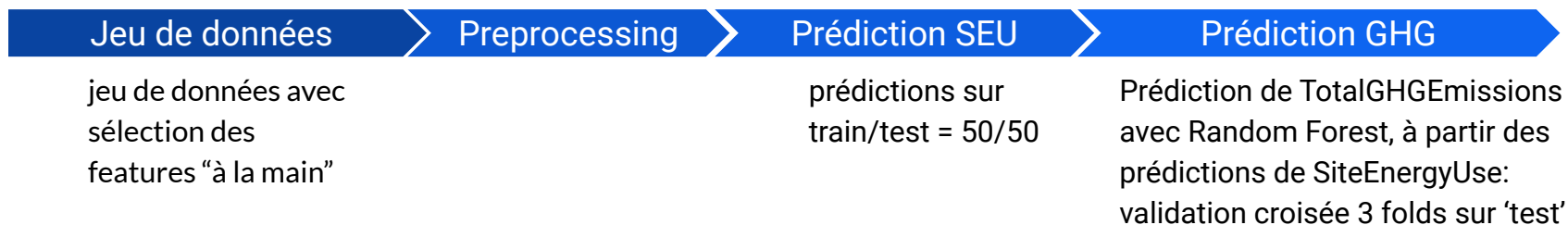
Sélection des features pour lesquelles feature importance > 2.5%

- Amélioration du temps d'entraînement
- Perte de performance maîtrisée
- Meilleure interprétabilité

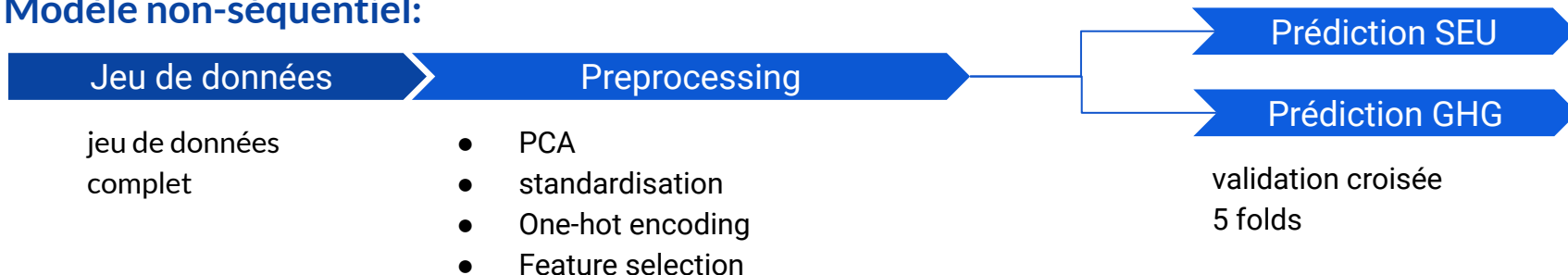


## Modèle final - Améliorations (2): modèle séquentiel vs non-séquentiel

### Modèle séquentiel:



### Modèle non-séquentiel:



## Résultats finaux

	EnergyStarScore	SiteEnergyUse	TotalGHGEmissions
Modèle séquentiel	sans ESS	0.82 +/- 0.09	0.82 +/- 0.02
	<b>avec ESS</b>	<b>0.90 +/- 0.04</b>	<b>0.92 +/- 0.01</b>
Modèle non-séquentiel	sans ESS	0.84 +/- 0.03	0.85 +/- 0.02
	<b>avec ESS</b>	<b>0.91 +/- 0.01</b>	<b>0.92 +/- 0.03</b>

- Prendre en compte ENERGYSTARScore permet d'améliorer les prédictions, à la fois par la mise à l'écart d'éléments introduisant du bruit dans le modèle et la prise en compte d'une variable corrélée avec la cible
- Les résultats sur chaque modèle sont équivalents