

# Projet de Data Science

## Analyse et Prevision de la Consommation Electrique

### Réalisé par :

- Othmane Ouzzine
- Zaynab EL moudni
- Hajar Bouzid
- Hakkou mouataz

**Encadré par :** M. Yasser El Alami

**Classe :** 1<sup>ère</sup> Année Master Ingénierie des Systèmes Informatiques (ISI) –Hybride

**Établissement :** SUPMTI RABAT

**Année universitaire :** 2025/2026

**Le :** 01/02/2026

## Objectif du Projet

Ce projet de data science a pour but d'étudier et analyser les tendances de consommation énergétique à partir de données réelles. Les objectifs principaux sont :

- Comprendre la structure des données de consommation électrique •  
Identifier les tendances et relations importantes entre les variables
- Construire des modèles prédictifs pour prévoir la consommation totale (`totalkW_mean`) •  
Communiquer clairement les résultats et recommandations

**Type de problème** : Regression (variable cible numérique continue)

**Variable cible** : `totalkW_mean` (consommation totale moyenne en kW)

## Dataset

- **Source** : Données de consommation électrique quotidienne
- **Taille** : 600 lignes, 34 colonnes
- **Période** : Juin 2018 à 2020
- **Variables** : Mesures de consommation de plusieurs compteurs (53, 71, 71A, 83), données météorologiques (température, humidité, insolation), et variables temporelles

## Structure du Projet

```

.
├── data/
│   ├── raw/                # Données brutes
│   ├── processed/          # Données nettoyées et préparées
│   └── external/           # Données externes
├── notebooks/
│   ├── 01_collecte_donnees.ipynb
│   ├── 02_analyse_exploratoire_EDA.ipynb
│   ├── 03_preparation_donnees.ipynb
│   └── 04_modelisation.ipynb
├── src/                    # Modules Python réutilisables
│   ├── data_processing.py
│   ├── evaluation.py
│   ├── models.py
│   └── visualization.py
├── models/                 # Modèles entraînés sauvegardés (.pkl)
├── visualizations/         # Graphiques générés
├── reports/               # Rapports et résultats
├── presentation/          # Support de présentation
└── requirements.txt        # Dépendances Python
  
```

## 1. Collecte des Données ([01\\_collecte\\_donnees.ipynb](#))

Ce notebook est dédié au chargement initial et à l'exploration préliminaire du dataset. Il couvre :

- **Chargement du dataset** : Lecture du fichier CSV contenant 600 observations et 34 variables
- **Aperçu des données** : Affichage des premières lignes pour comprendre la structure
- **Informations générales** : Types de variables (22 float, 5 int, 7 string), utilisation mémoire (0.34 MB)
- **Statistiques descriptives** : Résumé statistique de toutes les variables numériques
- **Analyse des valeurs manquantes** : Détection des colonnes avec des valeurs manquantes  
 (`insolation`: 44.5%, `totalkW_w-1` et `totalkW_w/1`: 1.17%, `totalkW_d-1` et `totalkW_d/1`: 0.17%)
- **Vérification des doublons** : Aucun doublon détecté
- **Cardinalité des variables catégorielles** : Analyse du nombre de valeurs uniques par colonne
- **Sauvegarde** : Export du dataset brut pour les étapes suivantes

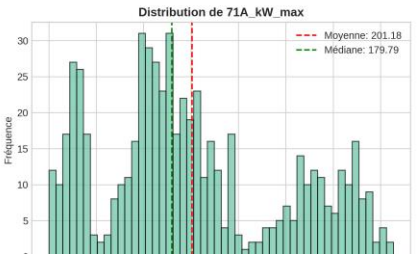
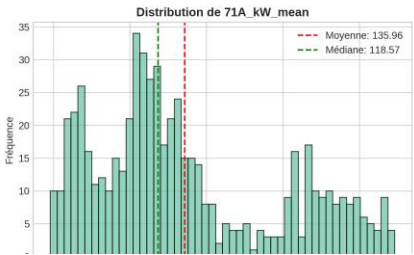
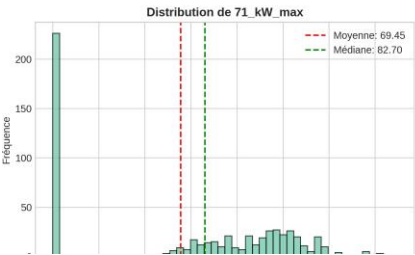
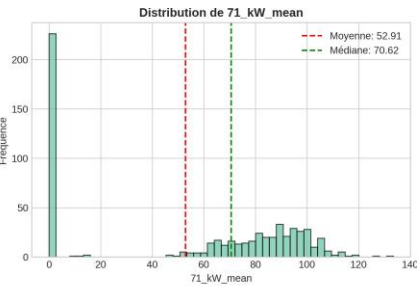
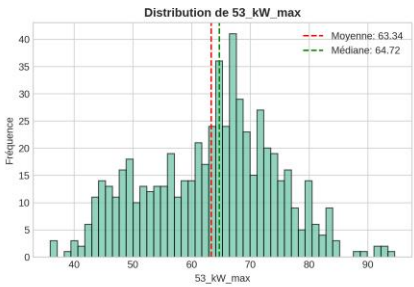
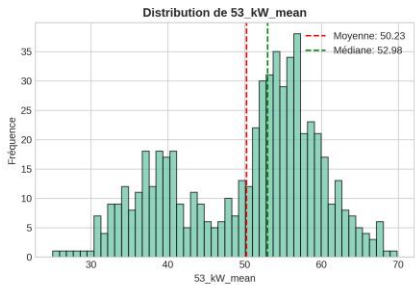
## 2. Analyse Exploratoire des Données - EDA ([02\\_analyse\\_exploratoire\\_EDA.ipynb](#))

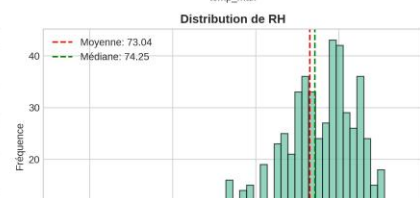
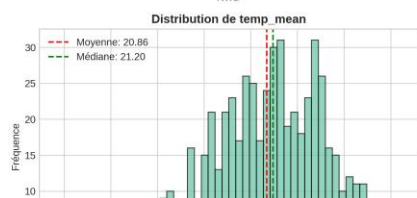
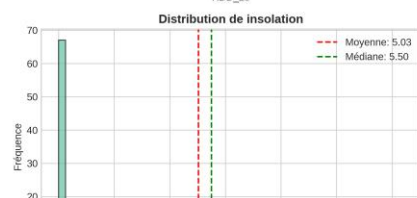
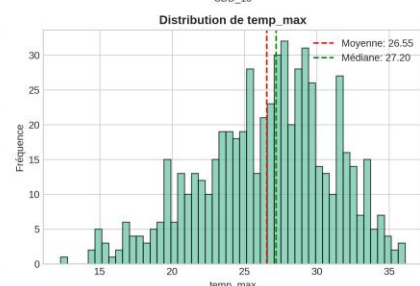
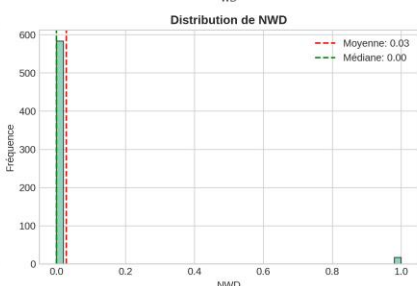
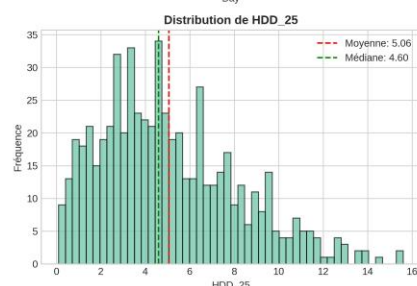
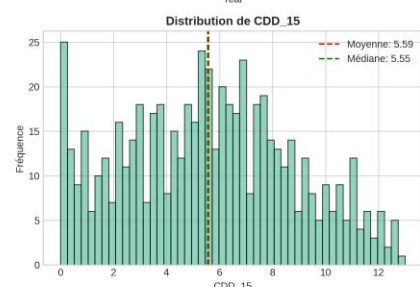
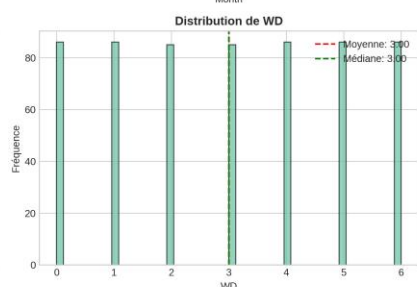
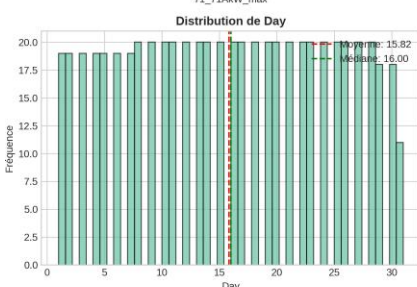
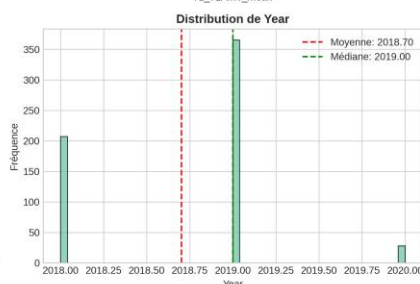
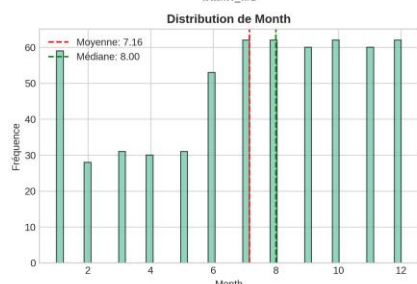
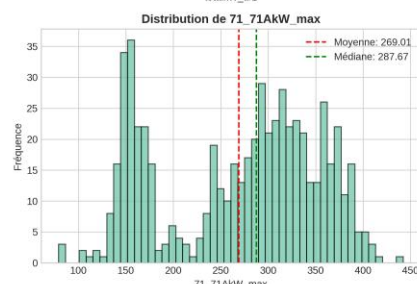
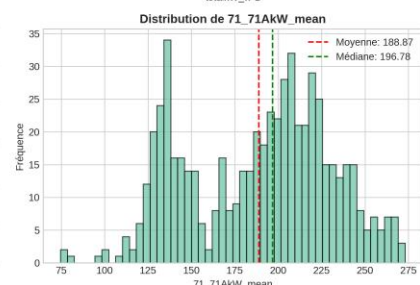
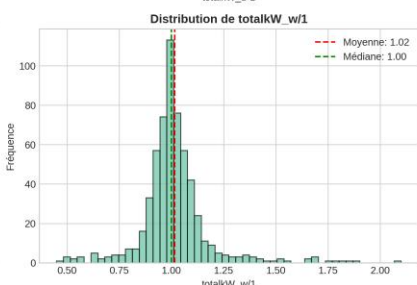
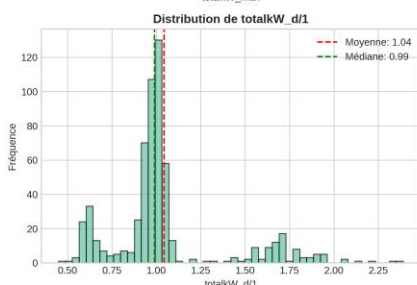
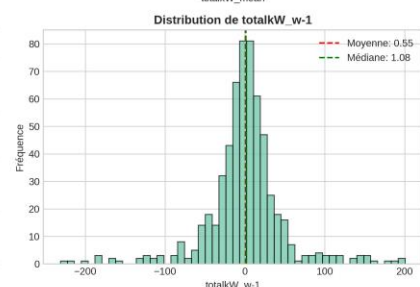
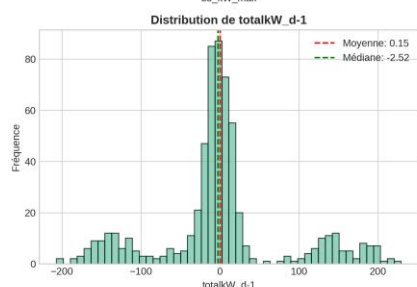
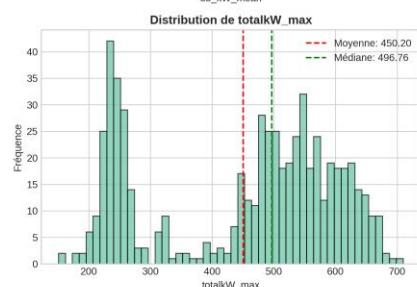
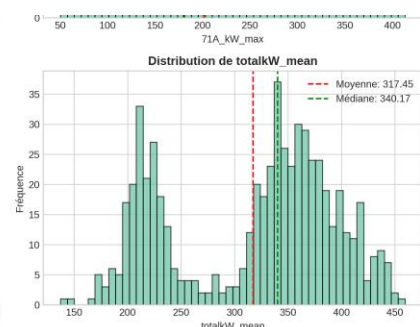
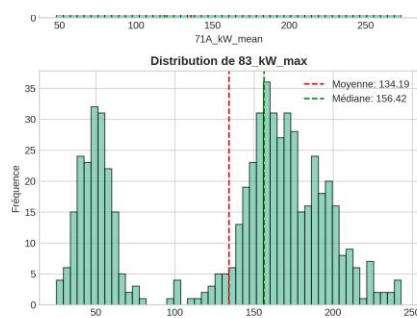
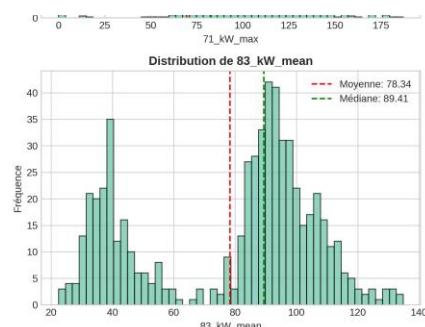
Ce notebook réalise une analyse approfondie des données pour identifier des tendances et formuler des hypothèses. Il comprend :

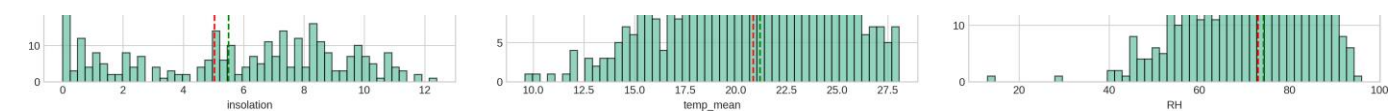
- **Statistiques descriptives avancées** : Calcul de la variance, asymétrie (skewness) et aplatissement (kurtosis) pour chaque variable
- **Analyse de l'asymétrie** : Identification des distributions asymétriques à droite (`71A_kW_mean`, `HDD_25`, `NWD`) et à gauche (`RH`)
- **Distributions des variables** : Histogrammes de toutes les variables numériques avec moyenne et médiane
- **Détection des outliers** : Boxplots avec comptage des outliers par la méthode IQR
- **Matrice de corrélation** : Heatmap des corrélations entre toutes les variables numériques
- **Corrélations fortes** : Identification de 48 paires de variables avec  $|r| > 0.7$
- **Scatterplots** : Visualisation des paires les plus fortement corrélées
- **Distribution des variables catégorielles** : Analyse des variables temporelles (heures de pic)

### Visualisations

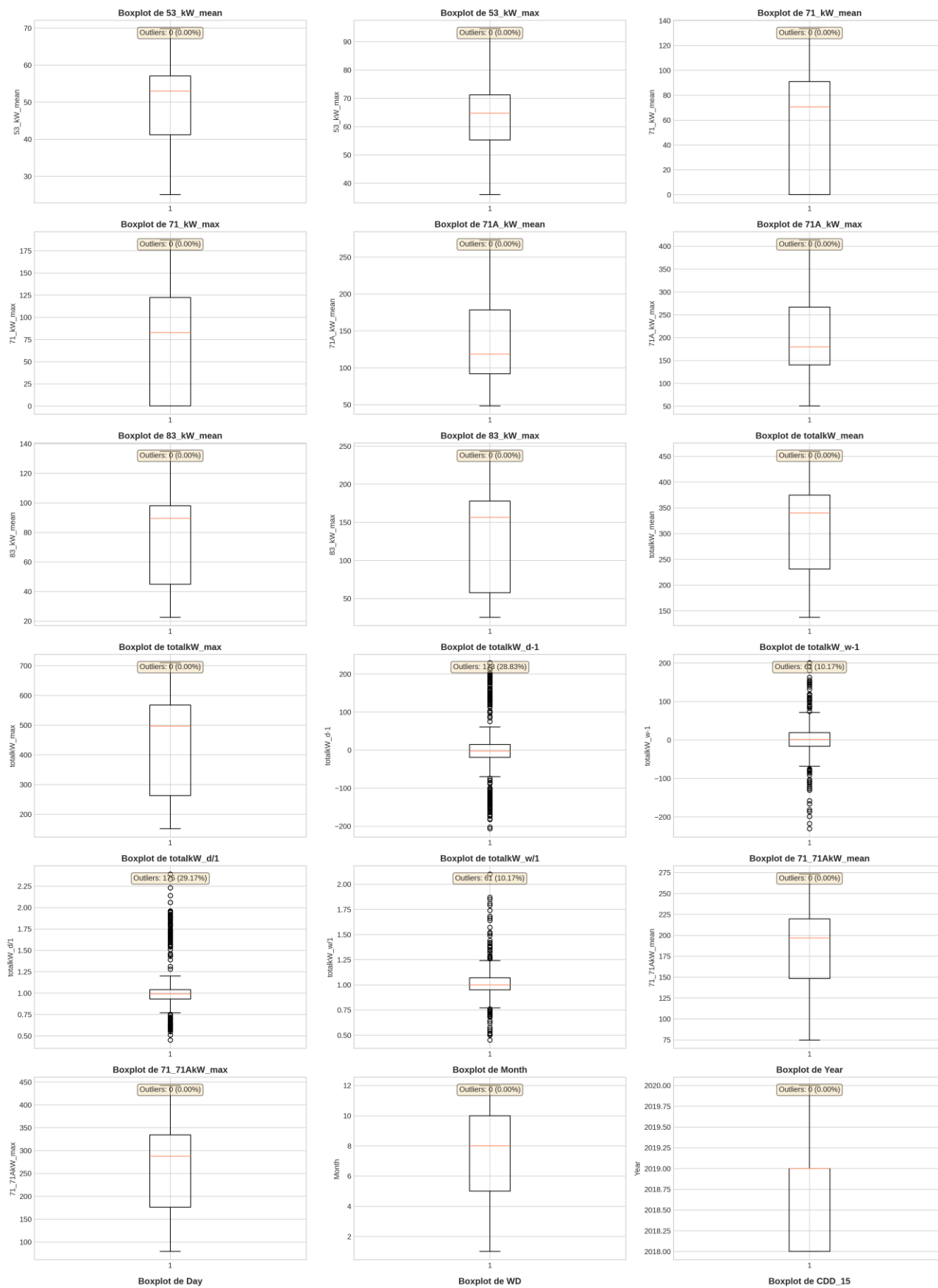
#### Distribution des variables numériques :

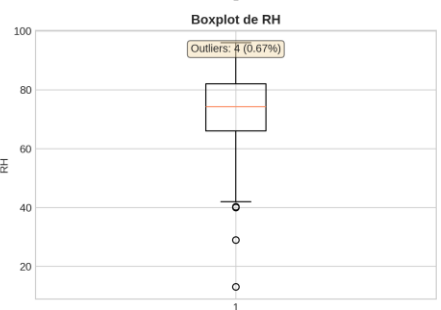
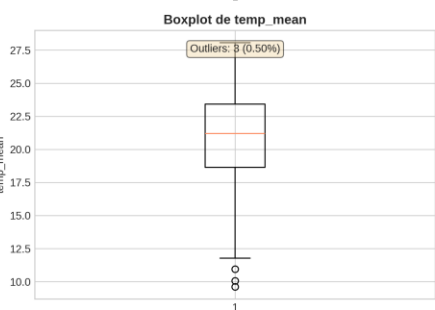
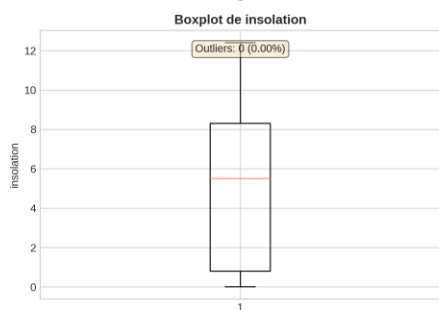
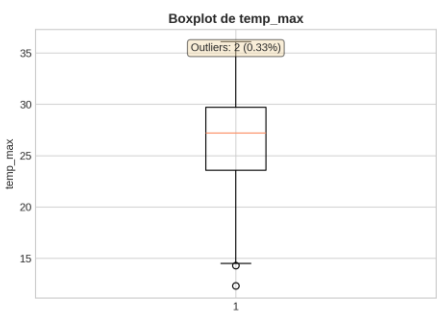
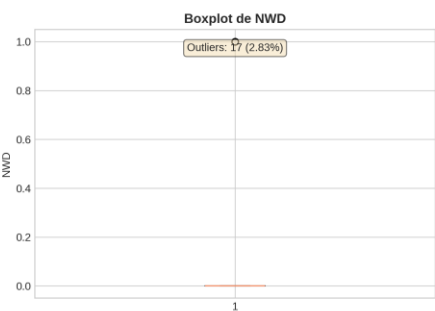
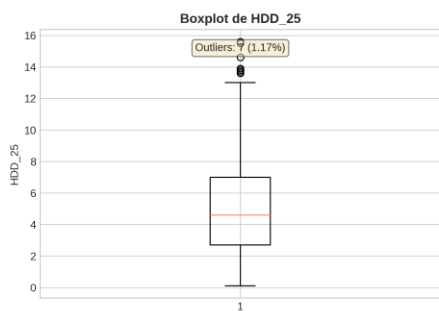
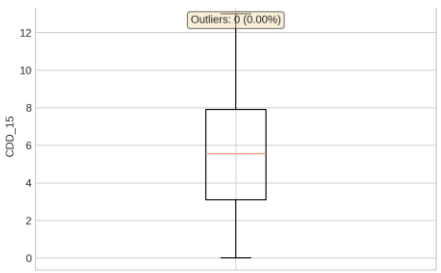
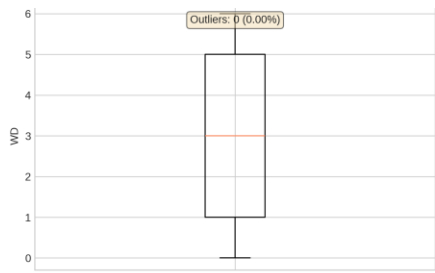
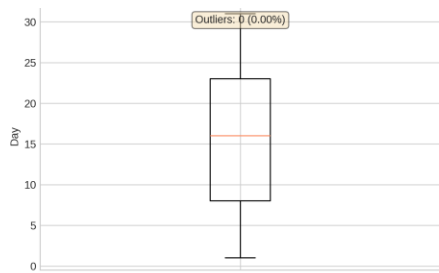






## Boxplots pour la detection des outliers :

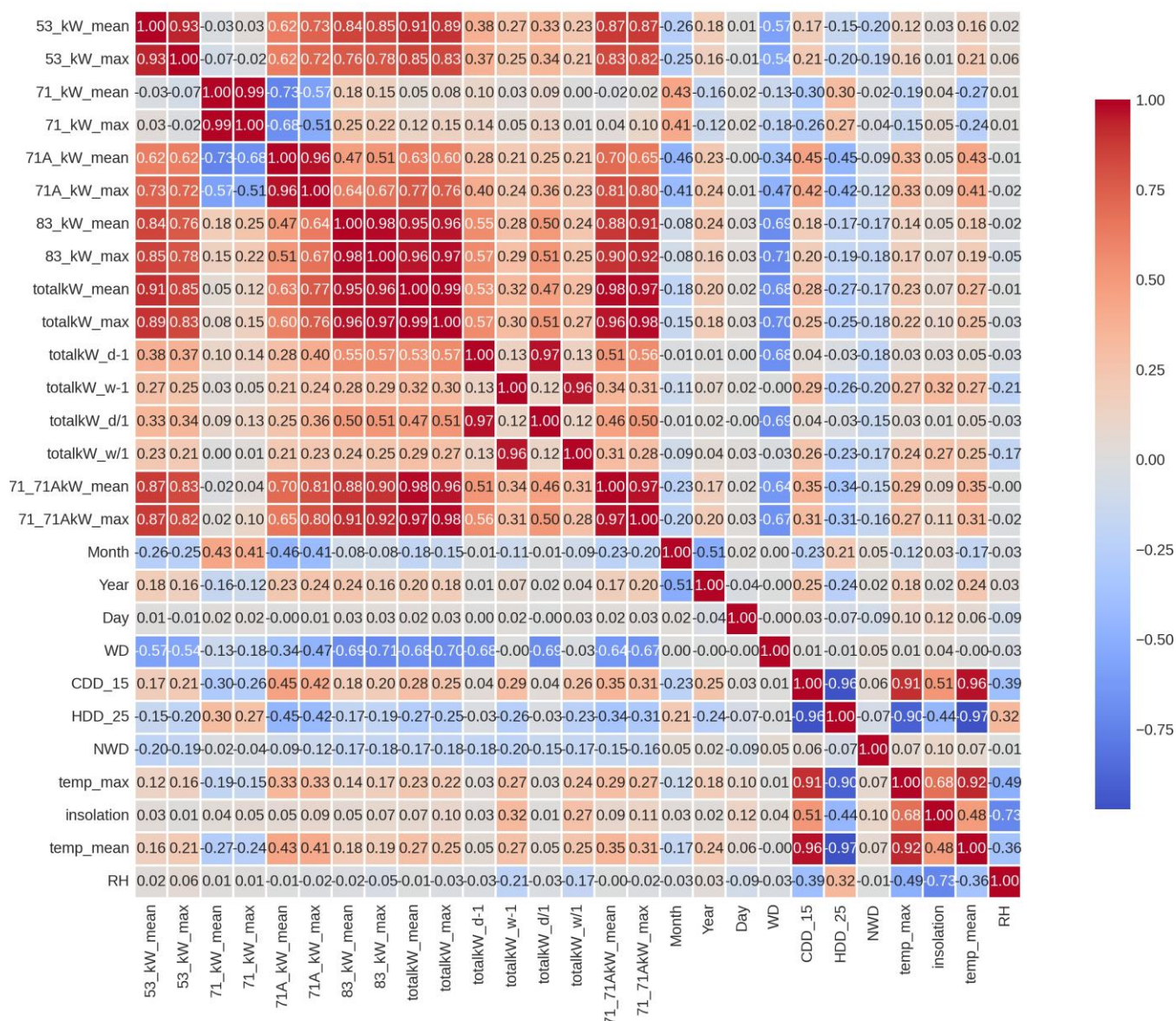




Matrice de corrélation :

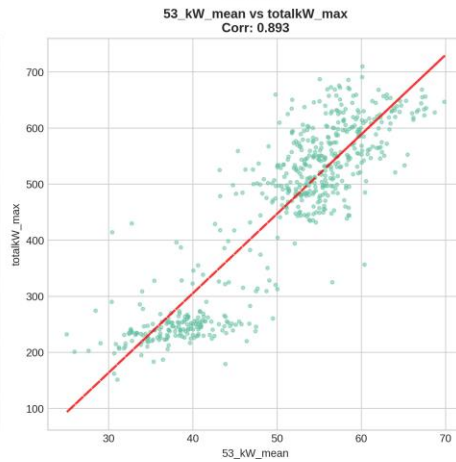
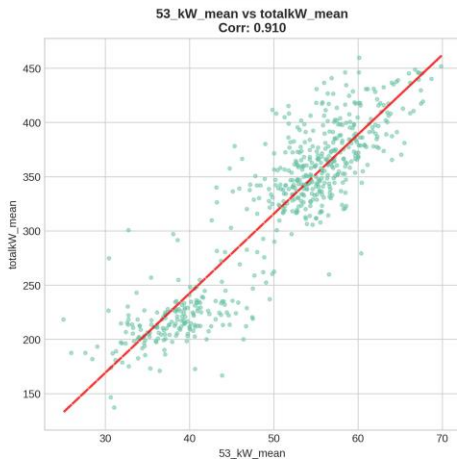
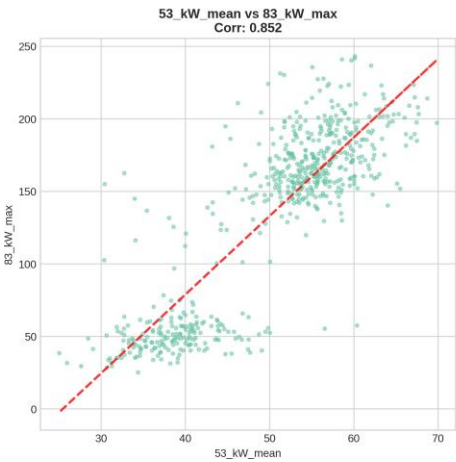
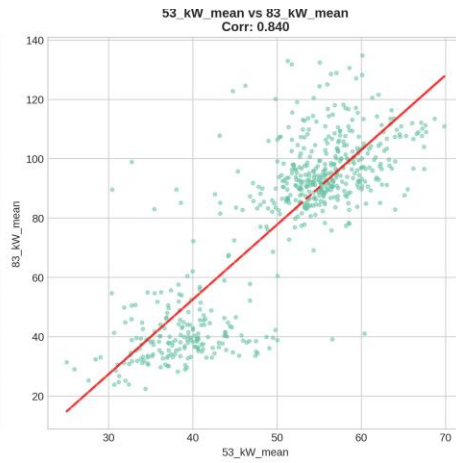
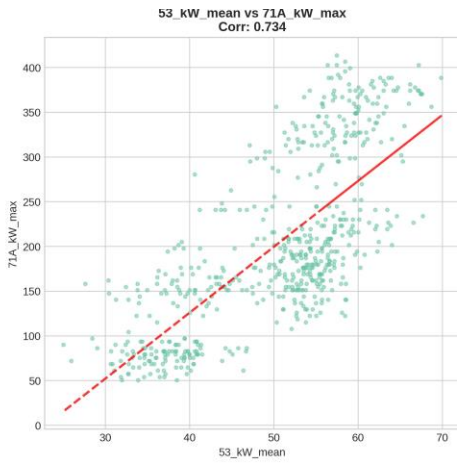
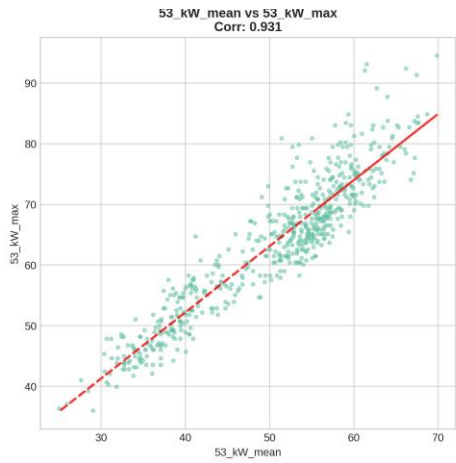


### Matrice de Corrélation des Variables Numériques

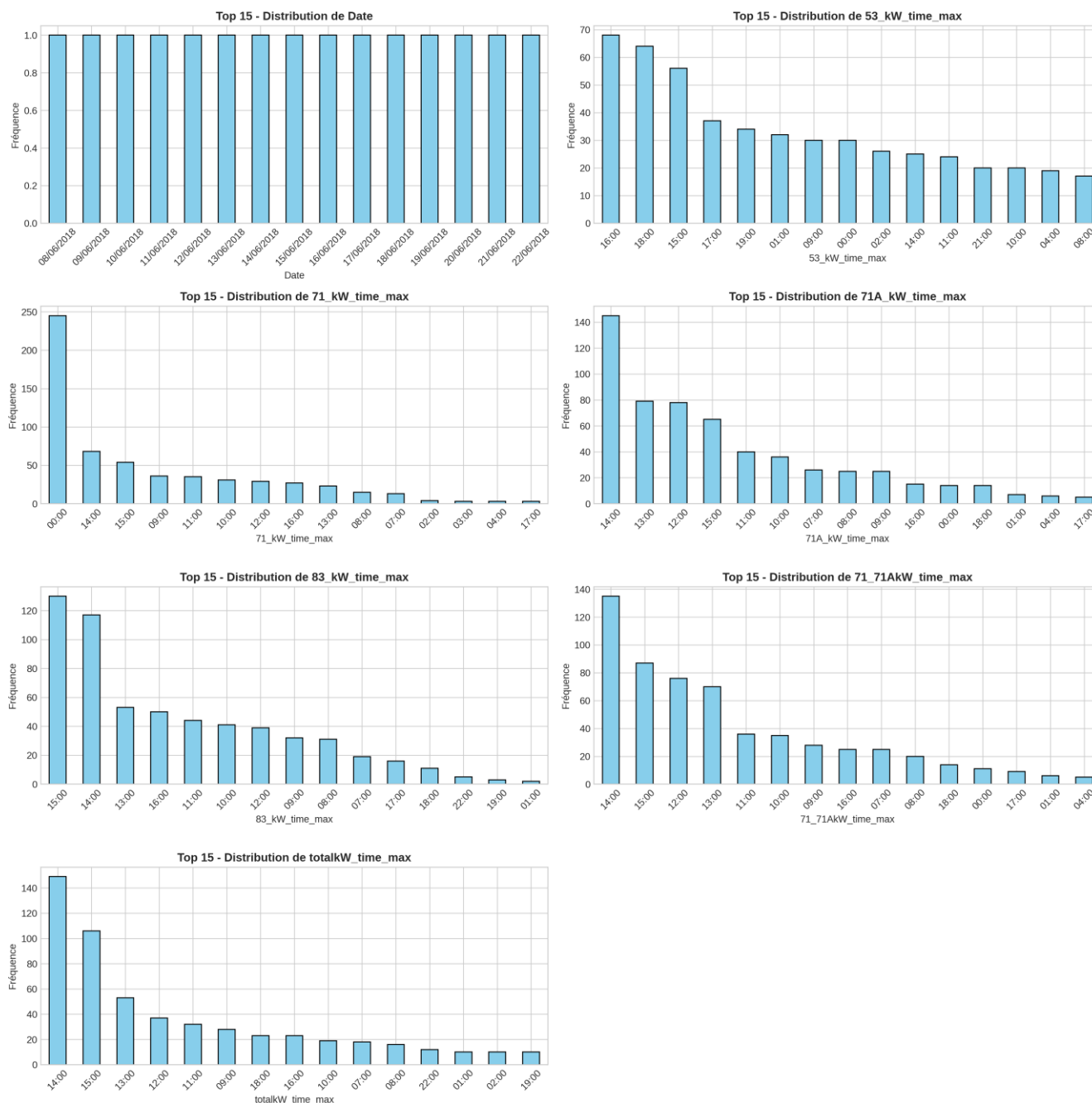


Scatterplots des correlations fortes :





Distribution des variables categorielles :



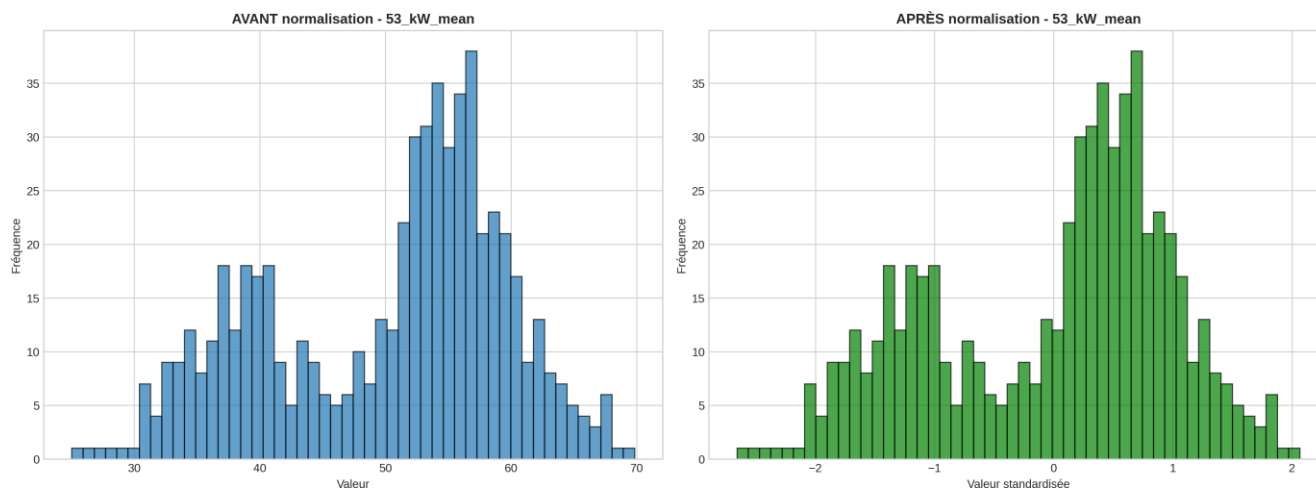
## 1. Préparation des Données (03\_preparation\_donnees.ipynb)

Ce notebook transforme les données brutes en un format exploitable pour la modélisation. Les étapes sont:

- **Suppression des doublons** : Vérification et suppression (0 doublons trouvés)
- **Gestion des valeurs manquantes** : Imputation par la médiane pour les variables numériques (totalkW\_d-1, totalkW\_w-1, totalkW\_d/1, totalkW\_w/1, insolation)
- **Traitement des outliers** : Winsorization (plafonnement) de 503 valeurs extrêmes sur 9 variables avec la méthode IQR
- **Encodage des variables catégorielles** : Suppression des colonnes temporelles (Date, heures de pic) car trop de catégories uniques
- **Normalisation** : Standardisation (Z-score) de toutes les features avec StandardScaler
- **Séparation train/test** : 80% entraînement (480 lignes) / 20% test (120 lignes)

- **Sauvegarde** : Export des ensembles X\_train, X\_test, y\_train, y\_test et du scaler
- Visualisation**

**Comparaison avant/apres normalisation :**



### 3. Modelisation (04\_modelisation.ipynb)

Ce notebook entraine et evalue trois modeles de regression pour predire la consommation electrique totale. Il comprend :

- **Modele 1 - Regression Lineaire** : Modele de base avec validation croisee 5-fold
- **Modele 2 - Random Forest Regressor** : Modele d'ensemble avec 100 arbres, profondeur max 15
- **Modele 3 - Gradient Boosting Regressor** : Modele de boosting avec taux d'apprentissage 0.1
- **Evaluation** : Metriques  $R^2$ , RMSE, MAE et MAPE pour chaque modele
- **Comparaison** : Tableau et graphiques comparatifs des performances
- **Importance des variables** : Identification des features les plus influentes pour chaque modele
- **Analyse des residus** : Distribution des erreurs de prediction
- **Sauvegarde** : Export des modeles entraines au format `.pkl`

### Resultats

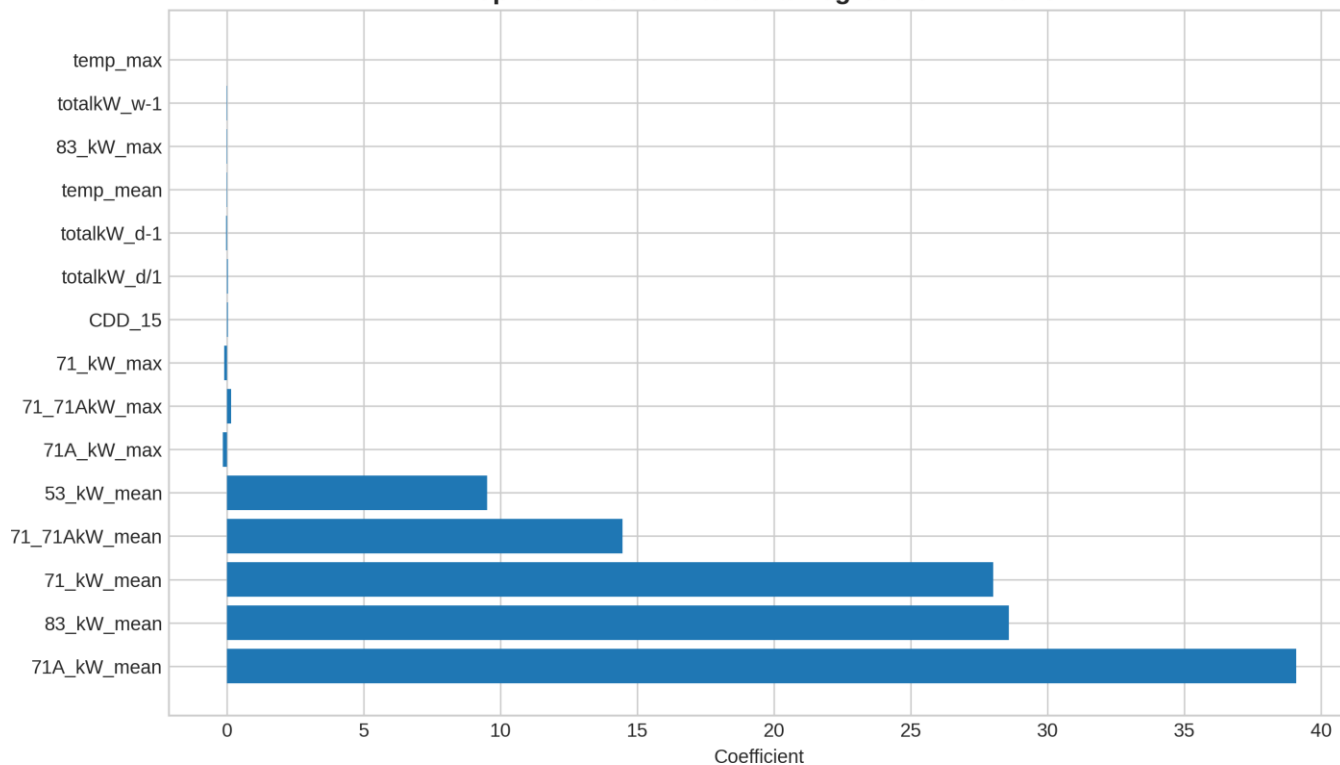
Modele	$R^2$	RMSE	MAE	MAPE (%)
Linear Regression	1.0000	0.0828	0.0662	0.02
Random Forest	0.9934	6.0606	4.1775	1.38
Gradient Boosting	0.9960	4.7214	3.2165	1.07

**Meilleur modele** : Regression Lineaire ( $R^2 = 1.0000$ )

### Visualisations

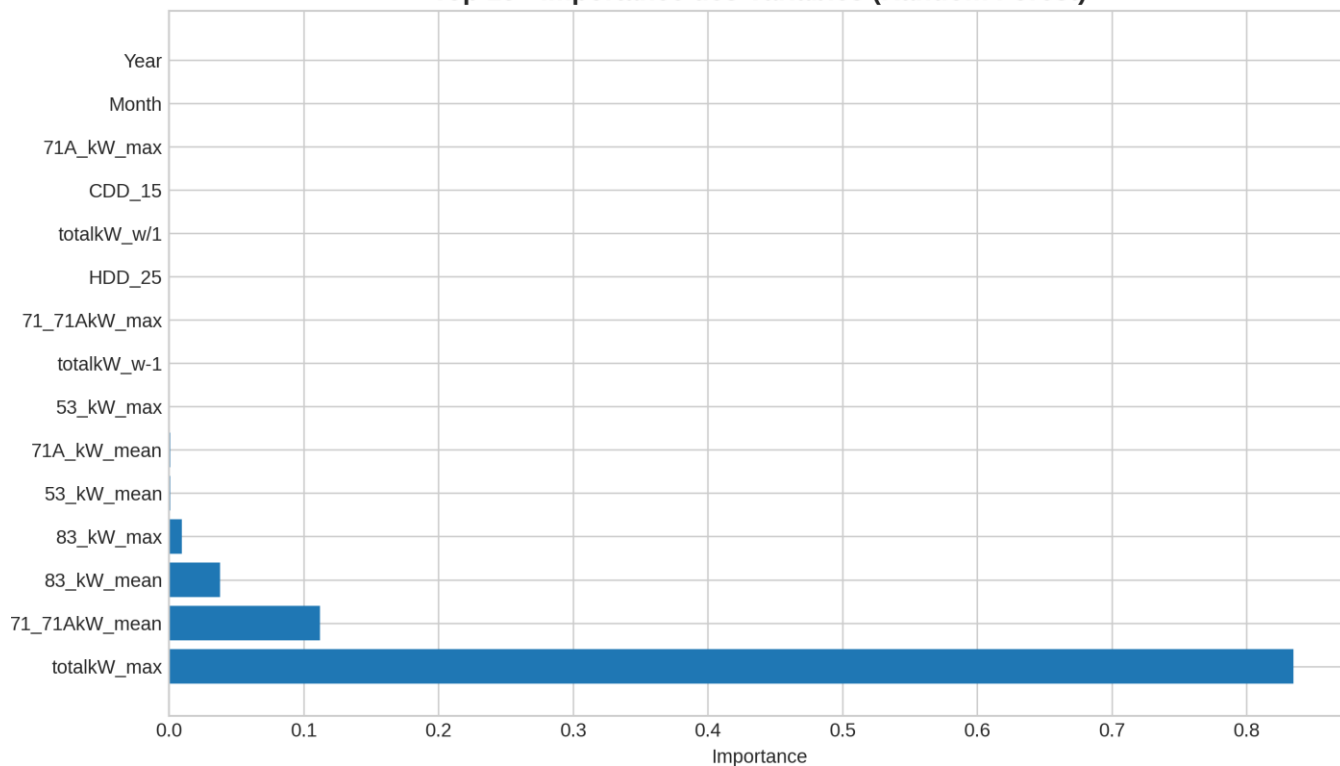
**Coefficients de la regression lineaire :**

Top 15 - Coefficients de la Régression Linéaire



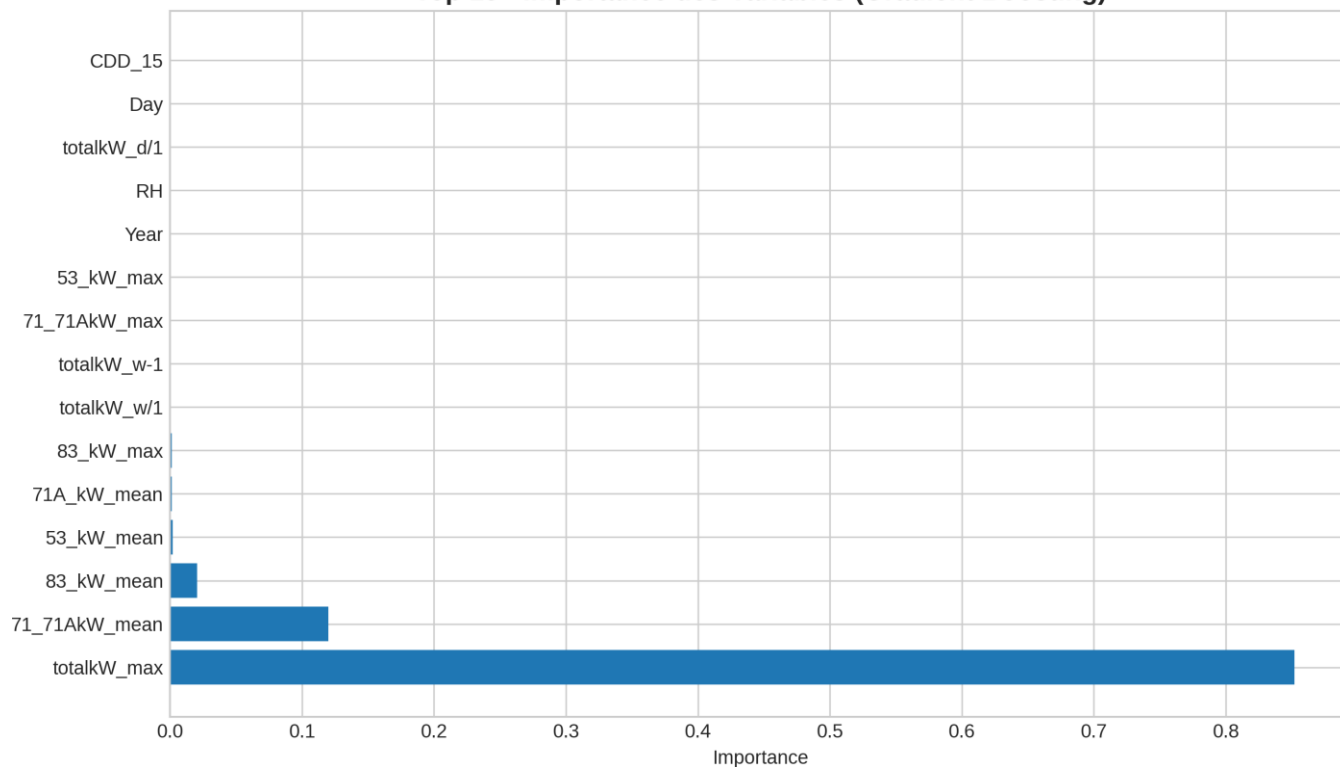
#### Importance des variables - Random Forest :

Top 15 - Importance des Variables (Random Forest)

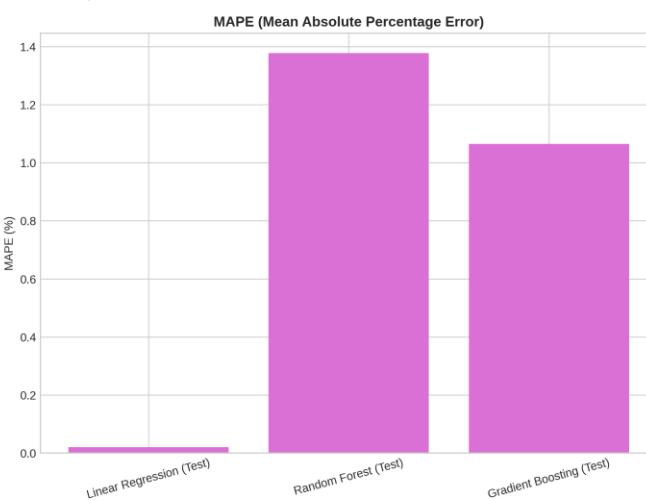
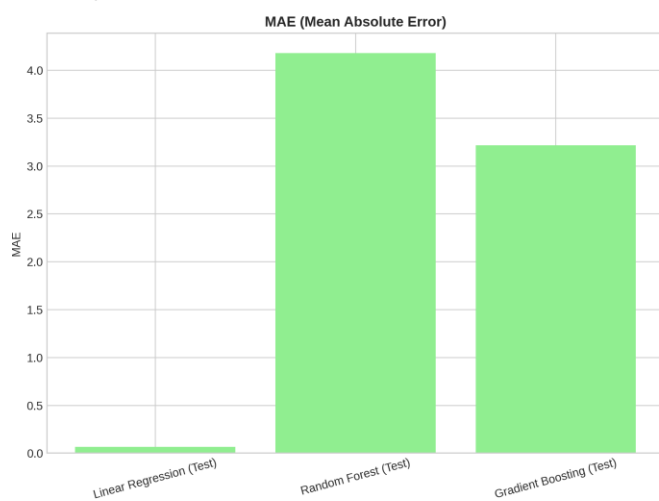
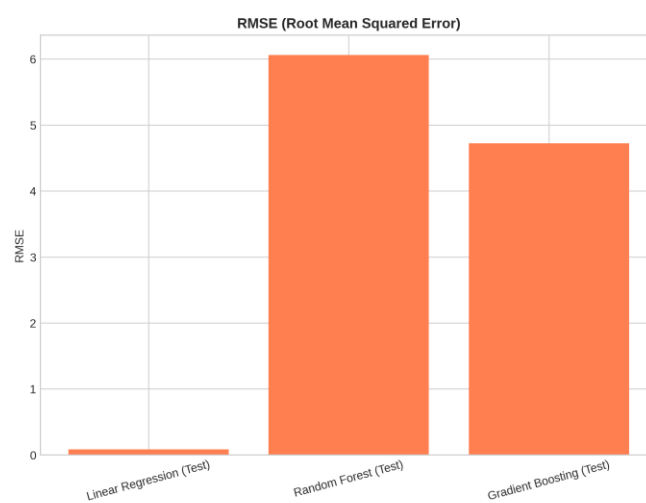
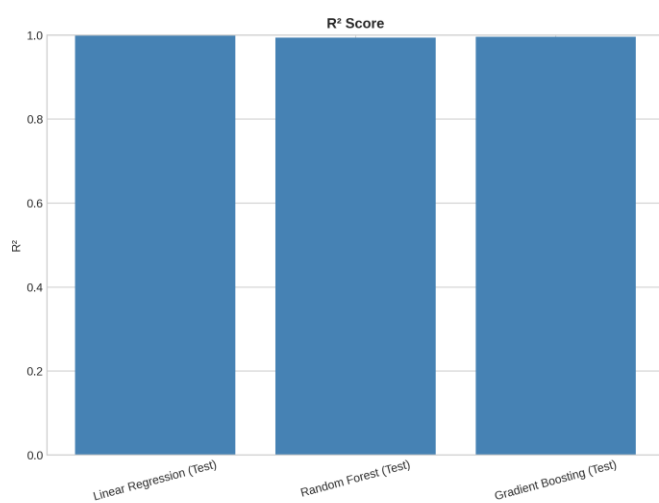


#### Importance des variables - Gradient Boosting :

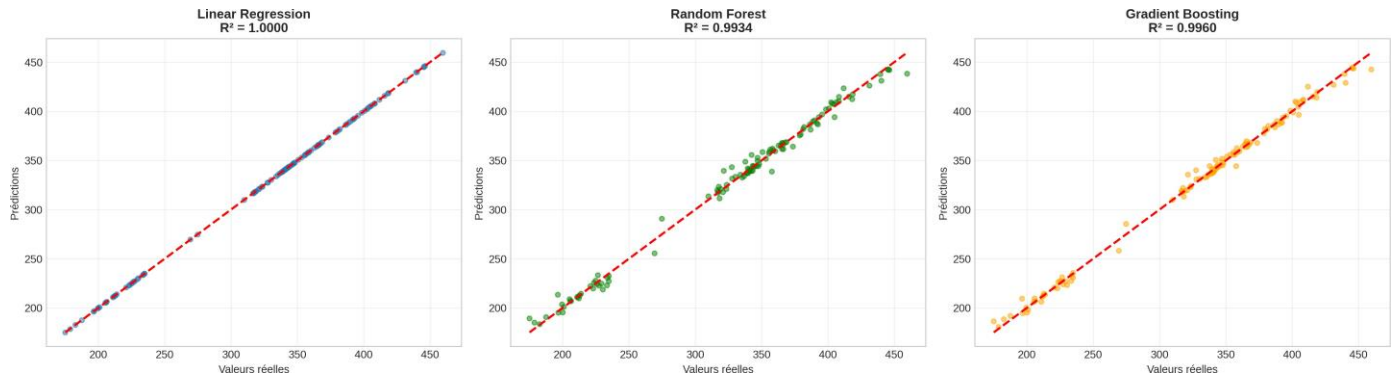
Top 15 - Importance des Variables (Gradient Boosting)



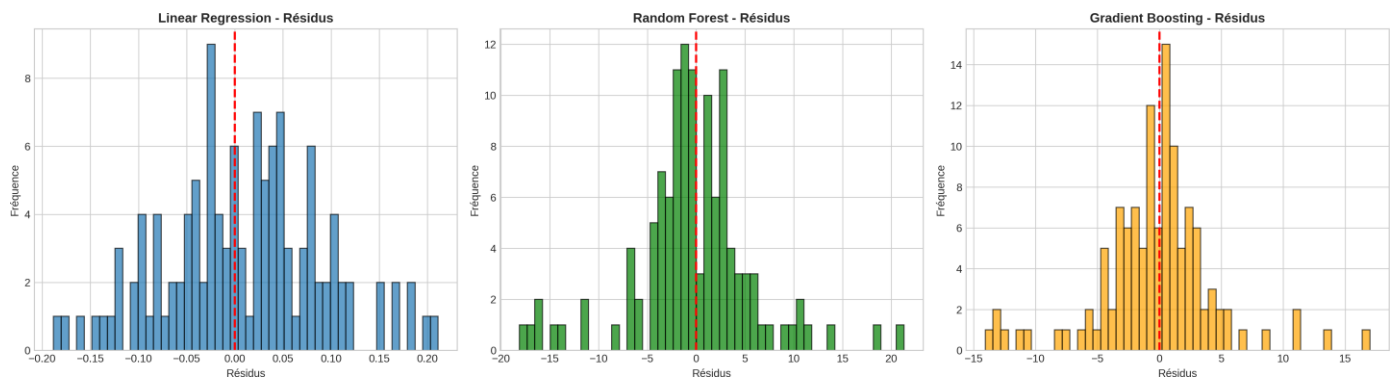
Comparaison des modeles :



## Predictions vs Valeurs réelles :



## Distribution des résidus :



## Technologies Utilisees

- **Python 3**
- **Pandas / NumPy** : Manipulation et traitement des donnees
- **Matplotlib / Seaborn** : Visualisations graphiques
- **Scikit-learn** : Modelisation (LinearRegression, RandomForestRegressor, GradientBoostingRegressor)
- **SciPy** : Analyses statistiques
- **Joblib** : Sauvegarde des modeles

## Installation

```
pip install -r requirements.txt
```

## Utilisation

Executer les notebooks dans l'ordre :

1. **01\_collecte\_donnees.ipynb** - Chargement et exploration initiale
2. **02\_analyse\_exploratoire\_EDA.ipynb** - Analyse exploratoire et visualisations
3. **03\_preparation\_donnees.ipynb** - Nettoyage et preparation des donnees
4. **04\_modelisation.ipynb** - Entrainement et evaluation des modeles