

Reconstitution d'un historique de consommation électrique sur 1 an, pour deux sites

Mahamadou Klanan Diarra

Sorbonne Université

10 octobre 2019

Sommaire

- 1 Introduction
- 2 Présentation du module
- 3 Données
- 4 Pré-traitement
 - Labellisation
 - Remplacement des nan
- 5 Construction du modèle
 - Modèle
 - Évaluation
- 6 Conclusion

Introduction

- **Challenge data** : plateforme en ligne qui donne accès à des données d'apprentissage automatique supervisé à des fins éducatives.
- **Planet OUI** : Fournisseur d'électricité verte, créé en 2007 et passe filiale de **BCM energy** en 2016.
- Projet d'apprentissage automatique : Analysé les données et prédire.
- Données (Planet OUI) : météorologiques et de consommations électriques à télécharger sur la plateforme challenge data
- L'objectif de ce challenge est de construire un modèle de prédiction de la consommation d'électricité sur deux sites (Lille, Aix-en-province)

Présentation

- Big data et ses applications : Présenté par **Georges Uzbelger**
- Les grandes branches du Big Data :
 - ① IA (1956) : Modélisation/simulation de comportements cognitifs
 - ② Machine learning : ensemble d'outils **Statistique + Optimisation**
 - ③ Data science (2008) : analyse statistique multivariée sur des données de grande taille
- Big Data dans l'entreprise
- Gestion d'un projet Data :
 - ① Comprendre les données
 - ② Préparation des données
 - ③ Modélisation
 - ④ Évaluation
 - ⑤ Déploiement

Données

Pour ce projet nous disposons de 3 tableaux de 8760 observations : 2 fichiers de training (input_train, output_train) et 1 fichier de test (input_test).

- **input (15 colonnes)** : température et humidité pour les 2 sites concernés, consommation de 3 sites secondaires, coordonnées géographiques
- **output (2 colonnes)** : consommations pour les 2 sites d'intérêt.

Données

| ID | timestamp | temp_1 | temp_2 | mean_rational_temp | humidity_1 | humidity_2 | loc_1 | loc_2 | loc_secondary_1 | loc_secondary_2 | loc_secondary_3 | consumption_secondary_1 | consumption_secondary_2 | consumption_secondary_3 |
|----|---------------------|--------|--------|--------------------|------------|------------|----------------|-----------------|------------------|------------------|-----------------|-------------------------|-------------------------|-------------------------|
| 0 | 2016-11-01T00:00:00 | 8.5 | 11.1 | 95.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 143 | 74 | 168 |
| 1 | 2016-11-01T01:00:00 | 8.0 | 11.1 | 96.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 141 | 80 | 162 |
| 2 | 2016-11-01T02:00:00 | 6.8 | 11.0 | 97.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 142 | 80 | 164 |
| 3 | 2016-11-01T03:00:00 | 7.5 | 10.9 | 99.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 139 | 80 | 162 |
| 4 | 2016-11-01T04:00:00 | 6.1 | 10.8 | 96.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 154 | 80 | 164 |
| 5 | 2016-11-01T05:00:00 | 5.4 | 10.7 | 98.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 184 | 82 | 167 |
| 6 | 2016-11-01T06:00:00 | 5.0 | 10.5 | 99.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 182 | 78 | 173 |
| 7 | 2016-11-01T07:00:00 | 5.1 | 10.2 | 99.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 185 | 79 | 175 |
| 8 | 2016-11-01T08:00:00 | 4.7 | 10.2 | 99.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 157 | 80 | 178 |
| 9 | 2016-11-01T09:00:00 | 5.5 | 10.3 | 100.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 153 | 86 | 175 |
| 10 | 2016-11-01T10:00:00 | 6.6 | 10.7 | 100.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 149 | 83 | 177 |
| 11 | 2016-11-01T11:00:00 | 18.4 | 11.4 | 98.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 153 | 82 | 176 |
| 12 | 2016-11-01T12:00:00 | 12.4 | 12.3 | 86.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 137 | 76 | 176 |
| 13 | 2016-11-01T13:00:00 | 13.8 | 13.2 | 80.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 139 | 77 | 178 |
| 14 | 2016-11-01T14:00:00 | 15.7 | 13.4 | 70.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 138 | 79 | 178 |
| 15 | 2016-11-01T15:00:00 | 15.0 | 13.7 | 78.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 139 | 77 | 175 |
| 16 | 2016-11-01T16:00:00 | 14.7 | 13.8 | 82.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 150 | 84 | 175 |
| 17 | 2016-11-01T17:00:00 | 14.1 | 13.2 | 81.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 150 | 84 | 174 |
| 18 | 2016-11-01T18:00:00 | 13.5 | 12.4 | 87.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 162 | 89 | 177 |
| 19 | 2016-11-01T19:00:00 | 12.8 | 11.6 | 87.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 163 | 80 | 178 |
| 20 | 2016-11-01T20:00:00 | 11.7 | 11.4 | 95.0 | | | [0.633, 3.067] | [43.530, 5.447] | [44.838, -0.579] | [47.476, -0.563] | [48.067, 2.333] | 154 | 76 | 171 |

| ID | consumption_1 | consumption_2 |
|----|---------------|---------------|
| 0 | 100 | 93 |
| 1 | 101 | 94 |
| 2 | 100 | 96 |
| 3 | 101 | 95 |
| 4 | 100 | 100 |
| 5 | 100 | 110 |
| 6 | 115 | 110 |
| 7 | 145 | 108 |
| 8 | 141 | 110 |
| 9 | 143 | 109 |
| 10 | 145 | 106 |
| 11 | 144 | 105 |
| 12 | 143 | 107 |
| 13 | 138 | 107 |
| 14 | 137 | 102 |
| 15 | 136 | 100 |
| 16 | 136 | 100 |
| 17 | 147 | 103 |
| 18 | 142 | 115 |
| 19 | 140 | 116 |
| 20 | 137 | 118 |

(a)

(b)

Figure – input vs output

Labellisation

But : Découper les observations par saison, période de la semaine.

Avantages : Permet une meilleure exploration et traitement des donnée.

- **Calendrier des saisons** : <https://www.kalendrier.com/dates-changements-saisons-2016.html>
- `days_label = ["{: %A} ".format(x) for x in input_train['timestamp']]`
- **Calendrier de vacances** :
<https://vacances-scolaires.education/annee-2016-2017.php>

Remplacement des nan

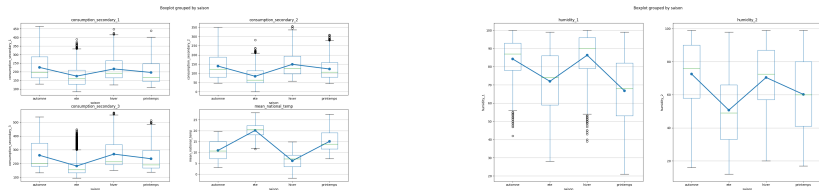
- 3.894 % de nan
- Supprimés = perte d'informations
- Remplacés par des valeurs adéquates

- 1 **K plus proche voisins**
- 2 **Stratégie : remplacé par la moyenne des K plus proche voisins**

[illegible]

Figure – Tableau complet

Modèle



- Utilisation de l'algorithme du **RandomForests** dans la bibliothèque **sklearn**.
- Entraînement par saison et par tranche de valeur :

Entraînement spécialisé

Combinaison de prédicteurs spécialisés par saison et périodes de faibles consommations.

Random Forest

- Algorithme proposé par **Léo Briemen** en 2001
- La méthode s'appuie sur deux techniques principales :
 - ① **Bagging : Bootstrap Aggregating**
 - ② **CART : Classification and Regression Trees**
- Algorithme très efficace : non paramétrique (s'adapte aux données), prédicteur plus stable (Bagging)

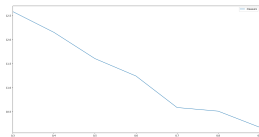


Figure – moyenne des écarts absolues vs
taile de l'échantillon de training

- Benchmark **24.33**
- Taille en pourcentage
- **MAE** sur les données
d'entraînement est à 12.48
(échantillon d'entraînement :
30%)

Taille vs mae

| taille | 0.30000000000000004 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 |
|---------|---------------------|--------------------|-------------------|--------------------|-------------------|--------------------|-------------------|
| measure | 12.484250520407485 | 11.976345718731892 | 11.52014927180635 | 11.029604356595687 | 10.43060365714876 | 10.234570138013474 | 9.968446641171903 |

Conclusion

Plusieurs points de conclusion importants :

- Tenue des étapes de suivie d'un projet,
- Une analyse de données cohérente,
- Meilleure utilisation des algorithmes de machine learning (KNN),
- Une meilleure compréhension du **Random Forest**,
- Un modèle avec de meilleurs résultats

Merci pour votre attention.
Des questions ?