
TP N° 6 : Analyse de données en pandas

Objectifs du TP : Manipuler des bases des données avec Pandas, affichage graphique avec Matplotlib.

Commencer par nommer votre fichier en suivant la même procédure que dans les autres TPs, en utilisant `filename` pour votre nom de TP :

```
# Changer ici par votre Prenom Nom:
prenom = "Joseph" # à remplacer
nom = "Salmon" # à remplacer
extension = ".ipynb"
tp = "TP4_HMLA310"
filename = "_".join([tp, prenom, nom]) + extension
filename = filename.lower()
```

- DONNÉES DE POLLUTION SUR PARIS (SOURCE : AIRPARIF) -

On va repartir ici des données vue en cours.

```
# Download
url = "http://josephsalmon.eu/enseignement/datasets/20080421_20160927-PA13_auto.csv"
path_target = "./20080421_20160927-PA13_auto.csv"
download(url, path_target, replace=False)

pollution_df = pd.read_csv('20080421_20160927-PA13_auto.csv', sep=';',
                           comment='#', na_values="n/d",
                           converters={'heure': str})

# 24:00 issues, voir https://www.tutorialspoint.com/python/time_strptime.htm
pollution_df['heure'] = pollution_df['heure'].replace('24', '0')

time_improved = pd.to_datetime(pollution_df['date'] +
                                ' ' + pollution_df['heure'] + ':00',
                                format='%d/%m/%Y %H:%M')

pollution_df['DateTime'] = time_improved
del pollution_df['heure']
del pollution_df['date']

pollution_ts = pollution_df.set_index(['DateTime'])
pollution_ts = pollution_ts.sort_index()

# Seulement les 4 années pleines
day_ini = '01/01/2009'
day_end = '12/31/2015'
pollution_ts = pollution_ts.loc[day_ini:day_end]

pollution_ts.head()
```

- 1) Passer en revue la base de données avec la fonction `describe()` de `pandas`.
- 2) Transformer la base pour n'en extraire que les années qui sont complètes. On affichera alors l'évolution de la concentration moyenne par jour sur toute la durée d'étude (par exemple en utilisant `resample`).
- 3) Afficher l'évolution de la pollution journalière pour les deux polluants (NO2 et O3), avec un subplot tout au long de la période d'étude. On utilisera la commande `resample` pour afficher la moyenne journalière sur toute la période de l'étude, en distinguant les 7 jours de la semaine.
- 4) La pollution atmosphérique montre-t-elle une tendance à la baisse au fil des ans? Pour cela on pourra faire une visualisation simple des moyennes annuelles sur la période d'étude.
- 5) Afficher le profils par mois sous forme de graphique en barres : on affichera pour les douze mois de l'année autant de barres qu'il y a d'années complètes.
- 6) Au vue des seuils légaux <https://www.airparif.asso.fr/reglementation/normes-europeennes> trouver combien de fois les valeurs limites et les seuils d'alertes ont-ils été franchi sur la période d'étude.
- 7) Trouver les 10 pics de pollutions les plus importants pour les deux polluants (jour et heure).

- DONNÉES DE CONSOMMATION ÉLECTRIQUE (SOURCE : EDF) -

On utilise la base de données¹ **Individual household electric power consumption Data Set**. Pour cela utiliser les commandes ci-dessous :

```
# download part if needed.
from download import download
import os
import zipfile

url = 'https://archive.ics.uci.edu/ml/machine-learning-databases/00235/'
# Lien alternatif:
# url=http://josephsalmon.eu/enseignement/datasets/household_power_consumption.zip

filename = 'household_power_consumption.zip'
cwd = os.getcwd()
path_target = os.path.join(cwd, filename)
download(url, path_target, replace=False)

# # unzip part
zip = zipfile.ZipFile(path_target)
zip.extractall()

# Visualisation succincte du fichier décompressé:
!head -10 'household_power_consumption.txt'
```

- 8) Quelle est la nature des colonnes dans cette base données? Dans la suite et sauf contre-indication on n'utilisera que la variable `Global_active_power`.
- 9) Charger la base de données avec `read_csv` :

```
na_values = ['?', '']
fields = ['Date', 'Time', 'Global_active_power']
df_conso = pd.read_csv('household_power_consumption' + '.txt', sep=';',
                      na_values=na_values, usecols=fields)
```

- 10) Selon vous, quel est le symbole qui encode les données manquantes ici?

1. Voir <https://archive.ics.uci.edu/ml/datasets/Individual+household+electric+power+consumption> pour un descriptif des données

- 11) En utilisant `describe` donner le nombre de ligne de la base de données et la moyenne de la variable `Global_active_power` sur toutes la durée de l'étude.
- 12) Utiliser `df_conso.tail()` et `df_conso.head()` pour trouver les dates de début et de fin d'étude.
- 13) En utilisant `df_conso.count()` calculer le pourcentage de lignes manquantes dans la base de données `df_conso`.
- 14) Lancer et décrivez ce que fait l'instruction suivante :

```
df_conso = df_conso.dropna(axis=0)
print('Taille da la base sans valeurs manquantes: {} lignes'.format(
    df_conso.shape[0]))
```

- 15) Utiliser `to_datetime` et `set_index` pour créer un objet `Serie` indexé par le temps (on prendra garde au format des dates internationales qui diffère du format français, et on utilisera l'option `infer_datetime_format=True` pour accélérer). On supprimera ensuite les colonnes `Date` et `Time` qui ne serviront plus, et on ne gardera que les années complètes.
- 16) Afficher le graphique des moyennes journalières entre le 1er janvier et le 30 avril 2007. Proposer une cause expliquant l'évolution de la consommation fin février et début avril.
- 17) Reprendre le question précédentes mais cette fois en produisant quatre sous-graphiques (avec `plt.subplot`) pour représenter les unes en dessous des autres les quatre années complètes de la base de données.
- 18) Proposer une visualisation pour analyser le comportement par jour de la semaine et mettre en évidence les différences de comportement entre le weekend et le reste de la semaine.
- 19) Faire de même pour analyser le comportement par mois au sein du foyer.

- AJOUTER DES DONNÉES CLIMATIQUE -

On ajoute des informations de température pour cette étude : les données utiles étant disponibles ici http://josephsalmon.eu/enseignement/TELECOM/MDI720/datasets/TG_STAID011249.txt². Ici les températures relevées sont celles d'Orly (noter cependant qu'on ne connaît pas le lieux de relevé de la précédente base de données).

- 20) Charger les données avec `pandas`, et ne garder que les colonnes `DATE` et `TG`. Diviser par 10 la colonne `TG` pour obtenir des températures en degrés Celsius. Traiter les éléments de température aberrantes comme des `NaN`.
- 21) Créer une `Serie pandas` des températures journalières entre le 1er janvier et le 30 avril 2007. Afficher sur un même graphique ces températures et la série `Global_active_power`.

2. on peut aussi trouver d'autres informations sur le site <http://eca.knmi.nl/dailydata/predefinedseries.php>