

# Analyse descriptive avec Pandas

[yzpt.github.io/module\\_pandas\\_b2](https://yzpt.github.io/module_pandas_b2)

Attendus du module :

- # Maîtriser la manipulation de données avec Pandas
  - lecture et écriture de fichiers CSV
  - sélection, filtrage, tri, groupement
- # Nettoyer un jeu de données
  - valeurs manquantes, suppression des doublons,
- # Calculer et interpréter des statistiques descriptives avec Pandas
  - tendance centrale, dispersion, distribution
- # Explorer et analyser les relations entre plusieurs variables
  - créer une matrice de corrélation
  - analyse bi-variée scatter plots
  - réaliser une régression linéaire simple
- # Comprendre comment identifier et traiter les valeurs aberrantes avec Python
  - méthode des z-scores
  - IQR
  - Box-plot

## Mercredi 11 décembre 2024

### 1. Cheatsheets:

- Pandas official cheatsheet
- EDA Cheatsheet
- Datacamp cheatsheet

### 2. Notebooks très bien rédigés :

Disponibles dans les dossiers 'workshop\_notebooks' et 'cookbook': \* Stefanie Molin Panda's wokrshop

### 3. Collection de ressources et datasets:

- ressources.md
- data/

### Premiers notebooks & divers:

- 1.0.introduction.ipynb
- 2.0.gapminder\_analysis.html

- 3.0.plotly\_express.ipynb
- 4.0.outliers.ipynb
- 5.0.structure\_pop.ipynb
- 6.0.join.ipynb

## Jeudi 12 décembre 2024

### 1. Analyse de données Gapminder

- 2.0.gapminder\_analysis.html
- 2.0.gapminder\_analysis.ipynb

### 2. Template d'EDA à conserver/compléter

- 7.0.EDA\_template.html
- 7.0.EDA\_template.ipynb

**3. Conversion de fichier HTML en notebook** Très utile avec la fonctionnalité github Pages pour afficher des notebooks en ligne. - Capsule vidéo Github Pages

- 8.0.html\_conversion.html
- 8.0.html\_conversion.ipynb

### 4. Fichier clients\_v2 à analyser : data/clients\_v2.csv

## Vendredi 13 décembre

### 1. Analyser les données des parkings de la MEL

- Etat des parkings de la MEL:
- Informations : <https://data.lillemetropole.fr/catalogue/dataset/disponibilite-parkings>
- Depuis le fichier excel : data/parking.xls

**Notebook : 12.0.parkings.ipynb** HTML : 12.0.parkings.html

**Consigne : Effectuer une analyse simple des données récupérées.** \* Donner un histogramme de l'occupation des parkings de la MEL. \* Effectuer un classement des 5 parkings les plus occupés. \* Effectuer une petite visualisation en groupant les parkings par ville.

**Récupérer les données en temps réel au format JSON depuis une API**  
: [https://data.lillemetropole.fr/geoserver/wfs?SERVICE=WFS&REQUEST=GetFeature&VERSION=2.0.0&TYPENAMES=mel\\_mobilite\\_et\\_transpor](https://data.lillemetropole.fr/geoserver/wfs?SERVICE=WFS&REQUEST=GetFeature&VERSION=2.0.0&TYPENAMES=mel_mobilite_et_transpor%3Aparking&OUTPUTFORMAT=application%2Fjson)  
t%3Aparking&OUTPUTFORMAT=application%2Fjson

## 2. Outliers

- Faire une EDA du dataset 10.0.dataset\_\_health.csv en détectant et éliminant les outliers.
- Pour ce, s'appuyer sur le notebook 10.1.health\_outliers.ipynb (html: 10.1.health\_outliers.html)
- Autre notebook d'exemples de traitement des outliers: divers/wine\_outliers.ipynb

## 3. S'exercer sur des datasets sélectionnés :

- amazon\_sales
- animals
- reddit\_data
- weather
- data/iris.csv
- data/titanic.csv
- data/winequality.csv

### Dataset du devoir surveillé :

- 11.0.dataset\_orders.csv : dataset de commandes de produits
- 11.0.dataset\_ogen\_orders.ipynb : notebook d'exemple d'analyse (+ script de génération du dataset)

Effectuer une EDA qui couvre l'ensemble des points et l'envoyer par mail à l'adresse **yohann.zapart@gmail.com**, en-tête de sujet : [GEMA\_B2]

## Travail à rendre pour vendredi 20 décembre

- Effectuer une analyse descriptive sur un dataset de données réelles que vous choisirez.
- Le dataset devra contenir au moins 1000 lignes et 10 colonnes (si ce n'est pas le cas, me demander par mail une homologation).
- L'analyse devra contenir l'ensemble des points explicités dans le fichier fiche.md