

## TP9

## STATISTIQUES

**But du TP**

- Manipuler un tableur (excel ou libreoffice) et calculer des indicateurs statistiques.
- Ouvrir un tableur avec python et la bibliothèque pandas.
- Utiliser une bibliothèque nouvelle à l'aide d'une liste de commandes.
- Comparer des séries statistiques.

Dans ce TP, nous allons étudier les données de deux tableurs :

- « tous-les-pays-du-monde-INED.xls » : estimations 2021 de la population mondiale (par continent) et des données de natalités/mortalité.  
*Source : Institut national d'études démographiques - [www.ined.fr](http://www.ined.fr)  
World Population Prospects. Nations Unies. 2019*
- « Deces-France.xls » : nombre de décès par semaine, par sexe et par tranche d'âge en France de 2000 à 2021.  
*Source : Human Mortality Database, Short-Term Mortality Fluctuation data series (STMF), INSERM*

**I Manipuler un tableur simple**

**Exercice 1** 1. Ouvrir avec **Excel** ou **LibreOffice Calc** le tableur « tous-les-pays-du-monde-INED.xls ».

2. Remplir les cases vides en utilisant des formules (voir ci-dessous)

- Population totale ;
- Population moyenne par continent ;
- Taux de natalité moyen ;
- Taux de mortalité moyen ;
- Espérance de vie moyenne ;
- Population 65+ totale puis moyenne.

|                            |   |
|----------------------------|---|
| MOYENNE(...)               | Moyenne   |
| SOMME(...)                 | Somme   |
| SOMMEPROD(Liste1 ; Liste2) | Renvoie la somme des produits des éléments des deux listes. |

## II Python – Bibliothèque pandas

**Exercice 2** Pour la suite, nous allons passer par **Jupyter Notebook**. Téléverser les fichiers suivants :

- Python-ECG1-09-statistiques.ipynb
- Deces-France.xls

puis suivez les instructions du notebook.

Nous allons utiliser plusieurs nouvelles bibliothèques dans cet exercices :

```
import pandas as pd          # pandas : gestion de données
import matplotlib.pyplot as plt # matplotlib : graphiques
```

Voici une liste de commandes utiles pour ce TP. Le tableau des données importées sera ici noté `df` (dataframe) et `Colonne` est le *nom* d'une colonne de ce tableau.

**Attention : adapter le `df` en fonction du nom donné à votre tableau au moment de l'import.**

|  |   |
|--|---|
| <code>df = pd.read_csv(fichier)</code>                 | import d'un fichier csv.  |
| <code>df = pd.read_excel(fichier)</code>               | import d'un fichier excel   |
| <code>df</code>  | aperçu du tableau   |
| <code>df.head()</code> , <code>df.head(n)</code>       | 5 premières lignes, <i>n</i> premières lignes   |
| <code>df.tail()</code> , <code>df.tail(n)</code>       | 5 dernières lignes, <i>n</i> dernières lignes   |
| <code>df.shape</code>                                  | taille du tableau (lignes, colonnes)  |
| <code>df.info()</code>                                 | affiche les colonnes et leur type   |
| <code>df[Colonne]</code>                               | données de la colonne <code>Colonne</code> . Attention, si le nom de la colonne est une chaîne de caractères, on écrira son nom entre guillemets. |
| <code>df[[Colonne1,Colonne2,...]]</code>               | sélection des colonnes indiquées.   |
| <code>df[ df[Colonne] == 5 ]</code>                    | applique un filtre, ici la valeur de <code>Colonne</code> doit être égale à 5.  |
| <code>df[ (filtre1) &amp; (filtre2) ]</code>           | applique plusieurs filtres (booléens), le <code>&amp;</code> signifie 'et'  |
| <code>df.sort_values(Colonne)</code>                   | trie la <code>Colonne</code> par ordre croissant.   |
| <code>df.describe()</code>                             | statistiques de base  |
| <code>df.mean()</code>                                 | moyenne   |
| <code>df.std()</code>                                  | écart-type  |
| <code>df.count()</code>                                | nombre de valeurs   |
| <code>df.median()</code>                               | médiane   |
| <code>df.max()</code>                                  | maximum   |
| <code>df.min()</code>                                  | minimum   |
| <code>df.sum()</code>                                  | somme des données   |
| <code>plt.bar(abscisses, ordonnees)</code>             | diagramme en barres   |
| <code>plt.boxplot(df[Colonne])</code>                  | boîte à moustaches  |
| <code>plt.hist(df[Colonne],range=(a,b),bins=n)</code>  | histogramme de <code>Colonne</code> , sur l'intervalle $[a, b]$ découpé en <i>n</i> sous-intervalles  |
| <code>df.hist(column = Colonne1, by = Colonne2)</code> | histogrammes de <code>Colonne1</code> pour chaque valeur de <code>Colonne2</code> .   |

Pour la suite de ce TP, ouvrez le Jupyter Notebook du TP9 et suivez les instructions.