

Dataframe : Tableau à deux dimensions composé de lignes et de colonnes. Il peut être comparé aux dictionnaires de Python : les clés sont les noms des colonnes et les valeurs sont les séries.

Librairies incontournables (à ajouter au début de son script Python) :

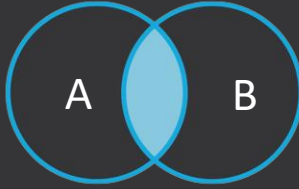
- **Pandas** (Outil d'analyse, de manipulation de données open source rapide, puissant, flexible et facile à utiliser). => `import pandas as pd`
- **Numpy** (Outil permettant d'effectuer des calculs scientifiques) => `import numpy as np`

Objectif	Remarque	Code Python
<p>Copier une table</p> <p>Remarque : En Python l'opérateur d'assignation crée aussi un pointeur vers la variable initiale</p>	<p><code>table_cpy = table</code> → Si on modifie les valeurs de "table", les valeurs de "table_cpy" seront également changées.</p>	<p><code>table_cpy = table.copy()</code></p>
<p>Regarder les premières lignes d'un dataframe</p>	<p>Observation des 5 premières lignes du dataframe "table".</p>	<p><code>table.head(5)</code></p>
<p>Créer une variable</p>	<p>Dans le dataframe "table" : <code>Total = var1 + var2</code></p>	<p><code>table["total"] = table["var1"] + table["var2"]</code></p>
<p>Filtrer un dataframe</p>	<p>Condition du filtre sur le dataframe "table" : <code>Total > 10 ET sexe = "Femme"</code></p>	<p><code>table[(table["total"] > 10) & (table["sex"] == "Female")]</code></p>
<p>SI ... ALORS ... SINON ...</p>	<p>Dans le dataframe "table" : Si <code>total > 10</code> alors la classe du total est "High". Sinon on attribue la valeur "Low" à la classe.</p>	<p><code>table["class_tot"] = np.where(table["class_tot"] > 10, "High", "Low")</code></p>
<p>Trier par valeur</p>	<p>Tri des valeurs en fonction des variables "total" et "sex" du dataframe "table". Remarque : Par default tri ascendant.</p>	<p><code>table.sort_values(["total", "sex"])</code></p>
<p>Sélectionner les données manquantes Ou non manquantes</p>	<p>Sélection des données égales "NA". Sélection des données non égales à "NA".</p>	<p><code>table[table["sex"].isna()]</code> <code>table[table["sex"].notna()]</code></p>
<p>Group by</p>	<p>Compte des données par sexe dans le dataframe "table". Remarque : Peut-être utilisé avec <code>.max</code>, <code>.min</code>, etc.</p>	<p><code>table.groupby("sex").count()</code></p>
<p>Transformation</p>	<p>Récupération des valeurs moyennes de la variable "total" par sexe.</p>	<p><code>moyenne = table.groupby("sex")["total"]</code> <code>table["transf_total"] = table["total"] - moyenne.transform("mean")</code></p>
<p>Processing by group</p>	<p>Renvoi la première ligne par groupe "sex" et "smoker".</p>	<p><code>table.groupby(["sex", "smoker"]).first()</code></p>
<p>Sélectionner Supprimer Renommer</p>	<p>Sélection des variables "sex", "total" et "total_class".</p>	<p><code>table[["sex", "total", "total_class"]]</code></p>
	<p>Suppression de la variable "sex".</p>	<p><code>table.drop("sex", axis)</code></p>
	<p>Attribution du nom "total_final" à la variable "total".</p>	<p><code>table.rename(columns={"total": "total_final"})</code></p>



Le Merge

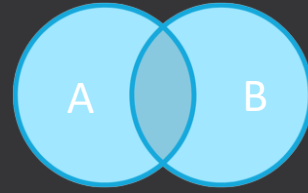
INNER JOIN



Joint les individus communs à A et B.

```
A.merge(B, on=["key"], how="inner")
```

FULL OUTER JOIN

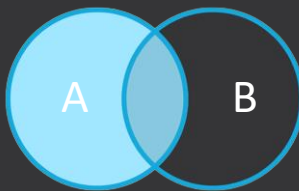


Joint tous les individus contenus dans A et B.

Remarque : La valeur « null » est attribuée quand l'enregistrement d'une table ne correspond pas à un enregistrement de l'autre.

```
A.merge(B, on=["key"], how="outer")
```

LEFT JOIN

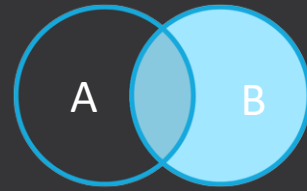


Joint les individus contenus dans A avec ceux ayant une correspondance dans B.

Remarque : La valeur « null » est attribuée pour les enregistrements de A qui n'ont pas de valeurs correspondantes dans B.

```
A.merge(B, on=["key"], how="left")
```

RIGHT JOIN



Joint les individus contenus dans B avec ceux ayant une correspondance dans A.

Remarque : La valeur « null » est attribuée pour les enregistrements de B qui n'ont pas de valeurs correspondantes dans A.

```
A.merge(B, on=["key"], how="right")
```

Les Formats

1

Créer un dictionnaire contenant en clé la valeur initiale et en valeur le format à appliquer.

```
dico_format = {0: "Non", 1: "Oui"}
```

2

Utiliser la fonction .map() pour appliquer les formats.

```
table["diabete"] = table["diabete"].map(p, (dico_format))
```

diabete

0

1

0

diabete

Non

Oui

Non