## **Pandas**



### Cours M2 CCN

Modélisation des problèmes scientifiques

2019/2020

# Pourquoi utiliser pandas?



 Pour charger, traiter et analyser des bases de données

 Pourquoi pandas plutôt qu'excel, ou libreoffice calc ? Parce qu'excel gère mal les grosses bases de données!

Comparaison excel-pandas

	<b>▼</b> Prénom	Age	Taille	
	Adele	45	168	
	Alex	30	182	
Individus	►Alice	28	153	
	Bob	52	161	

Variable

## Dataframe avec pandas



### **Import**

```
Entrée [1]: import numpy as np import pandas as pd
```

#### Création d'une base de données

#### Out[2]:

	id	bonjour	langue
0	1	Salut	Français
1	2	Hello	Anglais
2	3	Buongiorno	Italien
3	4	Guten Tag	Allemand
4	5	Buenos dias	Espagnol

### Import d'une base de données

```
Entrée [3]: pd.read_table("ronfle.csv", delimiter =',', index_col = 0).head()
Out[3]:
```

		AGE	POIDS	TAILLE	ALCOOL	SEXE	RONFLE	TABA
	P0001	47	71	158	0	0	0	1
	P0002	56	58	164	7	0	1	0
	P0003	46	116	208	3	0	0	1
	P0005	70	96	186	3	0	0	1
	P0006	51	91	195	2	0	1	1

Marche en local ou avec une page web

## Attributs des dataframes

Index



```
Entrée [4]: coucou.index # appel à l'index
            # L'index est la colonne grise qui identifie chacun des individus devant les données.
            # On peut fixer la variable id comme index par la commande suivante:
            coucou.set index("id")
   Out[4]:
                   bonjour
                           langue
                     Salut
                          Français
                     Hello
                            Anglais
              3 Buongiorno
                  Guten Tag Allemand
              5 Buenos dias Espagnol
                            Breton
                    Demat
            Types des variables
Entrée [3]:
            coucou.dtypes
   Out[3]: id
                        int32
                       object
            bonjour
            langue
                       object
            dtype: object
            Nombre total de données
Entrée [5]:
           coucou.size
   Out[5]: 18
            Nombre de colonnes
Entrée [6]: coucou.columns # libellés des colonnes
   Out[6]: Index(['id', 'bonjour', 'langue'], dtype='object')
            Dimensions du DataFrame
Entrée [7]:
            coucou.shape
   Out[7]: (6, 3)
```

L'index permet d'identifier les individus

Types de variables

# Accès et filtrage des données



### Appel aux données

```
Entrée [8]: # Variable(s)
             coucou["bonjour"] # pour appeler bonjour seulement
             coucou[["bonjour","langue"]] # pour appeler deux variables ou plus
             # Individus
             coucou.loc[2] # le troisième individu
     Out[8]: id
             bonjour
                         Buongiorno
             langue
                            Italien
             Name: 2, dtype: object
             Reauête
Entrée [9]: coucou.query("id < 4")</pre>
     Out[9]:
                      bonjour
                               langue
                               Anglais
              2 3 Buongiorno
                                Italien
Entrée [10]: coucou.query("langue == 'Français'")
   Out[10]:
                 id bonjour langue
                      Salut Français
```

coucou[coucou.columns[i]] donne la i+1-ème colonne

Requêtes directes sur la base de données :

9 → renvoie le dataframe avec les individus ayant un id inférieur à 4

10 → renvoie le dataframe avec les individus pour lesquels la variable langue vaut 'Français'

## Aggrégation de données



### Agrégation sur les dataframes



Variables quantitatives : tout ce qui est dénombrable, une mesure numérique lci la quantité de glace

Variables qualitatives : une catégorie, un renseignement, une qualification lci un parfum ou une couleur de glace

Group by, comme en SQL

# Pour aller plus loin



- Jointures de tables
- Ajouter ou supprimer des lignes/colonnes
- Comprendre les series
- Gestion des manquants (fillna)
- Export de données au format csv (to\_csv)
- S'intéresser au type de pandas pkl pour un import plus rapide

## Quelques liens sur pandas



- La doc : https://pandas.pydata.org/
- Encore lui!
   http://eric.univ-lyon2.fr/~ricco/tanagra/fichiers/fr\_Tanagra\_Data\_Manipulation\_Pandas.pdf
- Un script sur le github

 Pareil, il faut tester les fonctions et se les approprier. N'hésitez pas à googler 'pandas python' et à lire un ou deux article(s) pour vous faire une idée sur la librairie, ce qui est utile, les avantages, les faiblesses, etc.

Ça évite d'apprendre « parce que c'est au programme » ;)