

Mémo Pandas

ACTION	COMMANDE
Lire un fichier CSV en créant le Dataframe T	<code>T = pd.read_csv('nom.csv')</code>
Importer le module pandas	<code>import pandas as pd</code>
Donner le nombre de lignes de T	<code>len(T)</code>
Donner la taille de T	<code>T.shape</code>
Donner le nom et type de chaque colonne	<code>T.dtypes</code>
Donner les 5 premières lignes de T	<code>T.head()</code>
Donner les 10 premières lignes de T	<code>T.head(10)</code>
Donner les 5 dernières lignes de T	<code>T.tail()</code>
Calculer les indicateurs statistiques de T	<code>T.describe()</code>
Calculer les indicateurs statistiques de la colonne C	<code>T['C'].describe()</code>
Calculer la moyenne de la colonne C	<code>T['C'].mean()</code>
Calculer l'écart-type de la colonne C	<code>T['C'].std()</code>
Calculer la médiane de la colonne C	<code>T['C'].median()</code>
Calculer la somme de la colonne C	<code>T['C'].sum()</code>
Calculer l'effectif de la colonne C	<code>T['C'].count()</code>
Donner les effectifs de chaque valeur de la colonne C	<code>T['C'].value_counts()</code>
Trier T par ordre croissant selon la colonne C	<code>T.sort_values(by='C')</code>
Trier T par ordre décroissant selon la colonne C	<code>T.sort_values(by='C', ascending=False)</code>
Extraire la colonne C	<code>T['C']</code>
Extraire les colonnes C1 et C2	<code>T[['C1', 'C2']]</code>
Filtrer les lignes avec les valeurs de C non nulle	<code>T.query('C != 0')</code>
Filtrer les lignes avec A > 2 et B < 6	<code>T.query('A > 2 & B < 6')</code>
Filtrer les lignes avec A < 2 ou B > 6	<code>T.query('A < 2 B > 6')</code>
Filtrer les lignes de la colonne C correspondant à 'foo'	<code>T.query('C == "foo"')</code>
Filtrer les lignes de la colonne C contenant 'foo'	<code>T.query('C.str.contains("pu")')</code>
Filtrer les lignes de la colonne C avec les éléments d'une liste L	<code>T.query('C in @L')</code>
Échantillon aléatoire de 100 lignes	<code>T.sample(n=100)</code>
Filtrer les lignes d'index 7 et 8	<code>T.query('index in [7,8]')</code>
Supprimer les lignes en double selon une liste de noms de colonne	<code>T.drop_duplicates(subset=[c1, c2, ...])</code>

Grouper et compter selon la colonne C	T.groupby('C').count()
Grouper et calculer la moyenne selon la colonne C	T.groupby('C').mean()
Grouper et calculer la somme selon la colonne C	T.groupby('C').sum()
Grouper et trouver le max selon la colonne C	T.groupby('C').max()
Créer la colonne C comme somme terme de colonne A et de la colonne B	T['C'] = T['A'] + T['B']
Créer la colonne C comme produit de chaque terme de la colonne A par 3	T['C'] = 3 * T['A']
Créer la colonne C comme la somme de chaque time de la colonne A par 4	T['C'] = T['A'] + 4
Appliquer une fonction f à la colonne A	T['C'] = T['A'].apply(f)