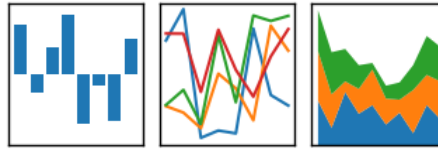


Mémento Pandas

pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



Lecture d'un fichier csv	
Importer le module pandas	<code>import pandas</code>
Lire le fichier <i>nom.csv</i> figurant dans le même dossier que le programme Python et affecter les données à une table de données, nommée T, dont chaque colonne correspond à une variable statistique.	<code>T = pandas.read_csv('nom.csv', sep = ';')</code> <code>sep = ','</code> si le séparateur du fichier csv est une virgule. <code>sep = '\t'</code> si le séparateur du fichier csv est une tabulation. Ajouter l'argument <code>encoding = 'latin-1'</code> si le fichier csv comporte des accents. Ajouter l'argument <code>decimal = ','</code> si les nombres décimaux du fichier csv sont séparés par une virgule.
Opérations de base sur une table pandas et ses colonnes	
Dimension de T : nombre de lignes (individus) et de colonnes (variables).	<code>T.shape</code>
(Nom) et type des variables (ou colonnes).	<code>T.dtypes</code>
Afficher les 5 premières lignes de T.	<code>T.head(5)</code>
Afficher les 5 dernières lignes de T.	<code>T.tail(5)</code>
Afficher des indicateurs statistiques pour chaque colonne de nombres de T.	<code>T.describe()</code>
Afficher des indicateurs statistiques pour la colonne de nombres nommée C.	<code>T['C'].describe()</code>
Calculer la moyenne de la variable C.	<code>T['C'].mean()</code>
Calculer l'écart type de la variable C. <u>Remarque</u> : la fonction <code>std()</code> fournit, par défaut, l'estimation s_{n-1} de l'écart type en considérant C comme un échantillon.	<code>T['C'].std(ddof = 0)</code>
Calculer la médiane de la variable C.	<code>T['C'].median()</code>
Calculer la somme de la variable C.	<code>T['C'].sum()</code>
Compter le nombre de valeurs de la colonne C.	<code>T['C'].count()</code>
Déterminer les effectifs de chacune des valeurs prises par la variable C.	<code>T['C'].value_counts()</code>
Trier T par ordre croissant selon la variable C.	<code>T.sort_values(by = 'C')</code>
Trier T par ordre décroissant selon la variable C.	<code>T.sort_values(by = 'C', ascending = False)</code>

Extraction d'un sous-tableau sous condition, tableau croisé avec pandas	
Créer une table T1 correspondant aux lignes de T où la variable C est non nulle.	<pre>T1 = T.query('C != 0') ou T1 = T[T['C'] != 0]</pre>
Créer une table T2 correspondant aux lignes de T où la variable A est supérieure à 2 et la variable B inférieure à 6.	<pre>T2 = T.query('A >= 2 and B <= 6') ou T2 = T[(T['A'] >= 2) & (T['B'] <= 6)]</pre>
Créer une table T3 correspondant aux lignes de T où la variable A est supérieure à 2 ou la variable B inférieure à 6.	<pre>T3 = T.query('A >= 2 or B <= 6') ou T3 = T[(T['A'] >= 2) (T['B'] <= 6)]</pre>
Créer une table T4 correspondant aux lignes de T où la variable A contient la chaîne de caractères "kiwi".	<pre>T4 = T.query('A == "kiwi"') ou T4 = T[(T['A'] == "kiwi")]</pre>
Créer une table T5 correspondant aux lignes de T où la colonne A contient une variable nommée var.	<pre>T5 = T.query('A == @var') ou T5 = T[(T['A'] == var)]</pre>
Créer une table T6 échantillon aléatoire sans remise de 100 lignes de T.	<pre>import random T6 = T.sample(n = 100)</pre>
Créer une table des effectifs des colonnes de T en regroupant les données selon les valeurs de la variable C.	<pre>T.groupby('C').count()</pre>
Créer une table des moyennes des colonnes de T en regroupant les données selon les valeurs de la variable C.	<pre>T.groupby('C').mean()</pre>
Créer une table des sommes des colonnes de T en regroupant les données selon les valeurs de la variable C.	<pre>T.groupby('C').sum()</pre>
Créer un tableau croisé en effectifs des variables C1 et C2 de la table T.	<pre>pandas.crosstab(T['C1'], T['C2'], margins = True)</pre>