

Importare pandas e numpy
<pre>import pandas as pd</pre> importa la libreria pandas
<pre>import numpy as np</pre> importa la libreria numpy
<pre>pd.set_option('max_columns', 300)</pre> limitare il numero di colonne stampate in output:*
<pre>pd.set_option('max_rows', 1000)</pre> limitare il numero di righe stampate in output:*

Caricare un file in un dataframe
<pre>df=pd.read_csv("file.csv")</pre> Legge un file csv
parametri aggiuntivi:
header specifica la riga da usare come header (int type), nel caso non sia presente: header=None
sep=';' specifica il separatore usato nel file csv
names=['col1','col2'] esplicita il nome delle colonne del dataframe, utile soprattutto in caso non ci sia header nel file
parse_dates = ['col1','col2'] lista di interi nomi di colonne da tentare di parsare come data.
<pre>df=pd.read_excel("nomefile.xlsx")</pre> legge un file excel

Esplorare un dataframe
<pre>df.shape</pre> restituisce (numero righe, numero colonne)
<pre>df[df.isnull().any(axis=1)]</pre> controlla se ci sono valori nulli
<pre>df.dtypes</pre> restituisce il tipo di variabile di ogni colonna
<pre>df.head()</pre> restituisce le prime (n) righe di un dataframe (default = 5)
<pre>df.tail()</pre> restituisce le ultime (n) righe di un dataframe (default = 5)
<pre>df.columns</pre> restituisce l'elenco delle colonne del dataframe
<pre>df.describe()</pre> restituisce una descrizione del dataframe con le statistiche di base

Gestire i dati mancanti
<pre>df.dropna()</pre> Elimina le righe con valori null presenti in qualsiasi colonna
<pre>df.fillna(value)</pre> Riempi le righe con valori null usando il valore <i>value</i> indicato tra parentesi

Salvare un dataframe in csv
<pre>df.to_csv('path/file_name.csv', sep=';', index=False)</pre> salva il dataframe in un file csv (rimuovendo l'indice e usando come separatore il ;)

Reshape dei dati
<pre>df.sort_values(by=['col1','col2'])</pre> ordina il dataframe per colonna o array di colonne
<pre>df.sort_values(by='col1', ascending = False)</pre> ordina il dataframe per colonna in ordine discendente
<pre>df.sort_index()</pre> ordina l'indice del dataframe
<pre>df.reset_index()</pre> resetta l'indice del dataframe usando l'ordine delle righe, l'indice attuale viene salvato in una nuova colonna
<pre>df.reset_index(name='new_name')</pre> assegna il nome alla colonna indice durante la fase di reset
<pre>df.drop(columns=['col1','col2'])</pre> elimina le colonne 'col1' e 'col2'
<pre>df.rename(columns={'old_column': 'new_column', 'old_column2': 'new_column2'})</pre> rinomina le colonne (usando un dizionario)
<pre>df.concat([df1,df2])</pre> aggiunge le righe del df2 al df1
<pre>df.concat([df1,df2], axis=1)</pre> aggiunge le colonne del df2 al df1

ricavare un sottoinsieme di un dataframe
<b>per righe</b>
<code>df[df['colonna1']&gt;n]</code> seleziona tutte le righe dove nella <i>colonna1</i> i valori sono maggiori di <i>n</i>
<code>df.drop_duplicates()</code> rimuove le righe duplicate
<code>df.iloc[10:20]</code> seleziona tutte le righe dalla 10 alla 20
<code>df[df['column'].isin(['value1','value2'])]</code> seleziona tutte le righe dove il valore nella colonna <i>column</i> è uguale a <i>value1</i> o <i>value2</i>
<b>per colonne</b>
<code>df['colonna1']</code> o <code>df.colonna1</code> seleziona la colonna 1
<code>df[['colonna1','colonna2','colonna4']]</code> Seleziona le colonne 1,2,4
<code>df.iloc[10:20, [1,2,5]]</code> seleziona tutte le righe dalla 10 alla 20 e le colonne in posizione 1,2,5
<code>df.loc[:, 'x2':'x4']</code> seleziona tutte le colonne dalla posizione x2 alla posizione x4 (compresa)
<code>df.loc[df['a']&gt;10, ['a','c']]</code> seleziona le colonne 'a' e 'c' dove il valore di a è maggiore di 10

Riassumere i dati
<code>len(df)</code> number of rows in a dataframe
<code>df['column1'].nunique()</code> numero di valori unici in una colonna
<code>df['column1'].sum()</code> somma dei valori di una colonna
<code>df['column1'].min()</code> minimo dei valori di una colonna
<code>df['column1'].max()</code> massimo dei valori di una colonna
<code>df['column1'].mean()</code> media dei valori di una colonna
<code>df['column1'].var()</code> varianza dei valori di una colonna
<code>df['column1'].std()</code> deviazione standard dei valori di una colonna
<code>df['column1'].count()</code> count dei valori non nulli di una colonna
<code>df['column1'].median()</code> mediana dei valori di una colonna
<code>df['column1'].quantile([.25, .75])</code> quantile dei valori di una colonna

Gestione oggetti di tipo datetime
<code>df['Data']=pd.to_datetime(df['Data'])</code> converte stringhe in oggetti datetime
<code>df['Data'].dt.round('1s')</code> arrotonda la data al secondo
<code>df['Duration']=(df['Data2']-df['Data1']).dt.seconds</code> Calcola la differenza tra due date restituendo un risultato in seocndi
<code>df['Data'].dt.date</code> estrae solo la data da data ora
<code>df['Data'].dt.month</code> estrae il mese
<code>df['Data'].dt.week</code> estrae la settimana dell'anno
<code>df['Data'].dt.weekday</code> estrae il giorno della settimana
<code>df['Data'].dt.weekday_name</code> estrae il giorno della settimana (come nome)
<code>df['Data'].dt.time</code> estrae l'ora completa
<code>df['Data'].dt.hour</code> estrae l'ora
<code>df['Stagione']=pd.cut((df['Data'].dt.dayofyear + 11) % 366, [0, 91, 183, 275, 366], labels=['inverno', 'primavera', 'estate', 'autunno'])</code> crea una colonna con la corrispondente stagione