

Python Cheat Sheet 3	DataFrames: estructuras en dos dimensiones	DataFrames: carga de datos	Metodos de DataFrames
<b>Pandas</b>	<b>Crear DataFrames</b> <code>df = pd.DataFrame(data, index, columns)</code> <b>data</b> : NumPy Array, diccionario, lista de diccionarios <b>index</b> : indice que por defecto se asigna como 0-(n-1), n siendo el número de filas; <b>index</b> = [lista] para asignar “etiquetas” (nombres de filas) <b>column</b> : nombre de las columnas; por defecto 0-(n-1); <b>columns</b> =[lista] para poner mas nombres	<b>Carga de datos</b> <code>df = pd.read_csv(“ruta/nombre_archivo.csv”)</code> crear un dataframe de un archivo de Comma Separated Values <code>df = pd.read_csv(“ruta/nombre_archivo”, sep= “;”)</code> crear un dataframe de un csv si el separador es ; <code>df = pd.read_csv(“ruta/nombre_archivo”, index_col= 0)</code> crear un dataframe de un csv si el archivo ya tiene una columna indice  <code>df = pd.read_excel(“ruta/nombre_archivo.xlsx”)</code> crear un dataframe de un archivo de Excel - si sale “ <b>ImportError:... openpyxl...</b> ”, en el terminal: <code>pip3 install openpyxl</code> o <code>pip install openpyxl</code>  <code>df = pd.read_json(“ruta/nombre_archivo.json”)</code> crear un dataframe de un archivo de JavaScript Object Notation (formato crudo) <code>df = df[‘data’].apply(pd.Series)</code> convertir el dataframe de json en un formato legible  <code>df = pd.read_clipboard(sep=‘\t’)</code> crear un dataframe de datos en forma de dataframe en el clipboard; el separador podria ser \n ; , etc.  <b>Pickle</b> : modulo que serializa objetos (convertir objetos complejos en una serie de bytes, en este caso en formato binario) para guardarlos en un archivo <code>with open(‘ruta/nombre_archivo.pkl’, ‘wb’) as f:</code> <code>pickle.dump(df,f)</code> pone los datos de un dataframe en el archivo.pkl  <code>pd.read_pickle(‘ruta/nombre_archivo.csv’).head(n)</code> leer n filas y 5 columnas del archivo pickle  <code>pd.read_parquet(‘ruta/nombre_archivo.parquet’)</code> leer un archivo parquet  <code>pd.read_sas(‘ruta/nombre_archivo.sas7bdat’, format = ‘sas7bdat’)</code> leer un archivo SAS de formato SAS7BDAT  <code>pd.read_spss(‘ruta/nombre_archivo.sav’</code> leer un archivo SAS de formato SAS7BDAT  <b>Guardado de datos</b> <code>df.to_csv(‘ruta/nombre_archivo.csv’)</code> guardar dataframe como archivo csv <code>df.to_excel(‘ruta/nombre_archivo.xlsx’)</code> guardar dataframe como archivo de Excel <code>df.to_json(‘ruta/nombre_archivo.json’)</code> guardar dataframe como archivo de JSON <code>df.to_parquet(‘ruta/nombre_archivo.parquet’)</code> guardar dataframe como archivo de parquet <code>df.to_pickle(‘ruta/nombre_archivo.pkl’)</code> guardar dataframe como archivo de pickle	<b>Metodos para explorar un dataframe</b> <code>df.shape</code> devuelve el número de filas y columnas <code>df.dtypes</code> devuelve el tipo de datos que hay en cada columna <code>df.columns</code> devuelve los nombres de las columnas <code>df.describe</code> devuelve un dataframe con un resumen de los principales estadísticos (media, mediana, desviación estándar etc.) de las columnas numéricas <code>df.describe(include = object)</code> devuelve un dataframe con un resumen de los principales estadísticos, incluyendo columnas con variables tipo string <code>df.info</code> devuelve un resumen sobre el no. de columnas, nombres de columnas, numero de valores no nulos y los tipos de datos de las columnas <code>df[“nombre_columna”].unique()</code> o <code>df.nombre_columna.unique()</code> devuelve un array con los valores únicos de la columna <code>df[“nombre_columna”].value_counts()</code> o <code>df.nombre_columna.value_counts()</code> devuelve una serie con el recuento de valores únicos en orden descendente <code>df.isnull()</code> o <code>df.isna()</code> devuelve True o False según si cada valor es nulo o no <code>df.isnull().sum()</code> o <code>df.isna().sum()</code> devuelve el número de valores nulos por columnas <code>df.corr()</code> devuelve la correlación por pares de columnas, excluyendo valores NA/nulos <code>df.set_index([“nombre_columna”], inplace = True)</code> establece el índice utilizando uno o mas columnas; puede sustituir o ampliar un índice existente <code>inplace = True</code> los cambios sobrescriben sobre el df * cuando una columna se cambia a índice ya no es columna *  <code>df.reset_index(inplace = True)</code> quitar una columna como indice para que vuelva a ser columna <code>df.rename(columns = {“nombre_columna”: “nombre_nueva”}, inplace = True)</code> cambia los nombres de una o mas columnas ejemplo de dict comprehension para crear diccionario sobre las columnas existentes de un dataframe: <code>diccionario = {col : col.upper() for col in df.columns}</code> <code>df.rename(columns = diccionario, inplace = True)</code> cambia los nombres de las columnas según el diccionario <code>df.drop([“columna1”, “columna2”], axis = b)</code> eliminar una o mas columnas o filas segun lo que especificamos <code>axis = 1</code> columnas <code>axis = 0</code> filas <code>df.rename(columns = diccionario, inplace = True)</code> cambia los nombres de las columnas según el diccionario <code>df[“columna_nueva”] = pd.cut(x=df[“nombre_columna”], bins=[n,m,l..])</code> separa los elementos de un dataframe en diferentes intervalos (n-m, m-l, etc); con este sintaxis se crea una columna nueva que indica en cual intervalo cae el valor <code>df.replace(to_replace = valor, value = valor_nuevo, inplace = True)</code> reemplaza cierto valor por otro que especificamos <code>df[“nombre_columna”].replace(to_replace = valor, value = valor_nuevo, inplace = True)</code> reemplaza cierto valor en una columna por otro que especificamos <code>df[“nombre_columna”] = df[“nombre_columna”] + x</code> reemplaza los valores de la columna por el valor + x (o otro valor que indicamos)
<b>Series: estructuras en una dimension</b>			
<b>Crear series</b> <code>serie = pd.Series()</code> crear serie vacía <code>serie = pd.Series(array)</code> crear serie a partir de un array con el indice por defecto <code>serie = pd.Series(array, index = [‘a’, ‘b’, ‘c’...])</code> crear una serie con indice definida; debe ser lista de la misma longitud del array <code>serie = pd.Series(lista)</code> crear una seria a partir de una lista <code>serie = pd.Series(número, indice)</code> crear una serie a partir de un escalor con la longitud igual al número de indices <code>serie = pd.Series(diccionario)</code> crear una serie a partir de un diccionario	<b>Acceder a informacion de un DataFrame</b> <code>df.loc[“etiqueta_fila”, “etiqueta_columna”]</code> devuelve el contenido de un campo en una columna de una fila <code>df.loc[“etiqueta_fila”,:]</code> devuelve los valores de todas las columnas de una fila <code>df.loc[:,“etiqueta_columna”]</code> devuelve los valores de todas las filas de una columna <code>df.iloc[indice_fila, indice_columna]</code> devuelve el contenido de un campo en una columna de una fila <code>df.iloc[indice_fila, :]</code> devuelve los valores de todas las columnas de una fila <code>df.iloc[:,indice_columna]</code> devuelve el contenido de un campo en una columna de una fila <code>df.loc[[lista_etiquetas_filas], [lista_etiquetas_columnas]]</code> devuelve el contenido de varias filas / varias columnas <code>df.loc[[lista_indices_filas], [lista_indices_columnas]]</code> devuelve el contenido de varias filas / varias columnas - se puede usar los indices/rangos de las listas [start:stop:step] dentro de los loc/iloc <code>df.loc[df.etiqueta &gt; x]</code> seleccionar datos basado en una condición usando operadores comparativos <code>df.loc[(df.etiqueta &gt; x) &amp; (df.etiqueta == y)]</code> seleccionar datos que tienen que cumplir las dos condiciones (and) <code>df.loc[(df.etiqueta &gt; x)   (df.etiqueta == y)]</code> seleccionar datos que tienen que deben cumplir una de las dos condiciones (or) <code>df.iloc[list(df.etiqueta &gt; x), :]</code> iloc no acepta una Serie booleana; hay que convertirla en lista <code>variable_df.head(n)</code> devuelve las n primeras filas del df, o 5 por defecto		
	<b>Crear columnas</b> <code>df[“nueva_columna”] = (df[“etiqueta_columna”] + x)</code> crea una nueva columna basada en otra <code>df = df.assign(nueva_columna= df[“etiqueta_columna”] + x)</code> crea una nueva basada en otra <code>df = df.assign(nueva_columna= [lista_valores])</code> crea una nueva columna de una lista de valores *tiene que ser de la misma longitud como el número de filas del dataframe* <code>df.insert(indice_nueva_columna, “nombre_columna”, valores)</code> crea una nueva columna en la indice indicada <code>allow_duplicates = True</code> parametro cuando queremos permitir columnas duplicadas (por defecto es False)		
	<b>Eliminar columnas</b> <code>df = df.drop(columns = [“column1”, “column2”])</code> eliminar columnas		