

Python Cheat Sheet 3	DataFrames	DataFrames: carga de datos	Metodos de DataFrames	Filtrados de datos
Pandas	Crear DataFrames <code>df = pd.DataFrame(data, index, columns)</code> data : NumPy Array, diccionario, lista de diccionarios index : indice que por defecto se asigna como 0-(n-1), n siendo el número de filas; index = [lista] para asignar “etiquetas” (nombres de filas) column : nombre de las columnas; por defecto 0-(n-1); columns = [lista] para poner mas nombres	Carga de datos <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 “ ImportError:... openpyxl... ”, 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. Pickle : 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 Guardado de datos <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 Librería PyDataset <code>pip install pydataset</code> o <code>pip3 install pydataset</code> from pydataset import data data() para ver los datasets listados en un dataframe por su id y titulo <code>df = data(‘nombre_dataset’)</code> guardar un dataset en un dataframe	Metodos para explorar un dataframe <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 estadísticosde las 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 indice utilizando uno o mas columnas; puede sustituir o ampliar un índice existente inplace = True los cambios sobreescriben sobre el df * cuando una columna se cambia a indice 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: diccionario = {col : col.upper() for col in df.columns} <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 axis = 1 columnas axis = 0 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)	<code>pd.options.display.max_columns = None</code> ejecutar antes del df.head() para poder ver todas las columnas Filtrado por una columna con operadores de comparación <code>variable_filtro = df[df[“nombre_columna”] == valor]</code> extrae las filas donde el valor de la columna igual al valor dado * se puede usar con cualquier operador de comparación * Filtrado por multiples columnas con operadores logicos <code>&</code> and <code> </code> or <code>~</code> not <code>variable_filtro = df[(df[“columna1”] == valor) & (df[“columna2”] == valor) & (df[“columna3”] > n valor)]</code> extrae las filas donde los valores de las columnas cumplan las condiciones en parentesis <code>variable_filtro = df[(df[“columna1”] == valor) (df[“columna1”] == valor)</code> extrae las filas donde los valores de las columnas cumplan con una condición u otra <code>variable_filtro = ~(df[df[“columna1”] == valor])</code> extrae las filas donde los valores de las columnas NO cumplan con la condición Metodos de pandas de filtrar <code>variable_filtro = df[df[“nombre_columna”.isin(iterable)]</code> extrae las filas cuyas valores de la columna nombrada están en el iterable (una lista, serie, dataframe o diccionario) <code>variable_filtro = df[df[“nombre_columna”].str.contains (patron, regex = True)]</code> extrae las filas cuyas valores de la columna nombrada contienen el patron de regex <code>variable_filtro = df[df[“nombre_columna”.str.contains (“substring”, case = False, regex = False)]</code> extrae las filas cuyas valores de la columna nombrada contienen el substring, no siendo case sensitive <code>variable_filtro = df[df[“nombre_columna”.str.contains (“substring”, case = False, regex = False)]</code> extrae las filas cuyas valores de la columna nombrada contienen el substring, no siendo case sensitive <code>df[pd.notnull(df[“nombre_columna”])]</code> devuelve las filas que no tiene valores nulos en la columna especificada Reemplazar valores basados en indices y condiciones: <code>indices_filtrados = df.index[df[“columna”] == “valor”]</code> for indice in indices_filtrados: <code>df[“nombre_columna”.iloc[indice] = “valor_nuevo”</code> Reemplazar valores basados en metodos NumPy: <code>df[“nueva_columna”] = np.where(df[“nombre_columna”] > n, “categoria_if_true”, “categoria_if_false”)</code> crea una nueva columna con los valores basados en una condición <code>df[“nueva_columna”] = np.select(lista_de_condiciones, lista_de_opciones)</code> crea una nueva columna con los valores basados en multiples condiciones
				