

Python: Pandas	DataFrames	Metodos de exploracion	Tipos de datos	Valores nulos
<b>Series: estructuras en una dimension</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>df.head(n)</b> devuelve las primeras n lineas del dataframe <b>df.tail(n)</b> devuelve las últimas n lineas del dataframe <b>df.sample(n)</b> devuelve n filas aleatorias de nuestro dataframe, o uno por defecto <b>df.shape</b> devuelve el número de filas y columnas <b>df.dtypes</b> devuelve el tipo de datos que hay en cada columna <b>df.columns</b> devuelve los nombres de las columnas <b>df.describe</b> devuelve un dataframe con un resumen de los principales estadísticos de las columnas numéricas <b>df.info()</b> devuelve un resumen sobre el no. de columnas, nombres de columnas, numero de valores no nulos y los tipos de datos <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 <b>df.duplicated().sum()</b> devuelve el numero de filas duplicadas	Tipos de datos en Pandas: - object - int64 - float64 - datetime, timedelta[ns] - category - bool <b>df.dtypes</b> devuelve el tipo de datos que hay en cada columna <b>df.tipo = df.select_dtypes(include = "tipo")</b> crea un dataframe de las columnas del tipo de datos especificado <code>df['columna'] = df['columna'].astype('tipo', copy = True, errors = 'ignore')</code> convierte una columna en el tipo de dato especificado <b>copy = True</b> devuelve una copia <b>copy = False</b> *cuidado: los cambios en los valores pueden propagarse a otros objetos pandas* <b>errors = ignore</b> omite excepciones; en caso de error devuelve el objeto original <b>errors = raise</b> permite que se generen excepciones	<b>Identificar nulos</b> <b>df.isnull()</b> o <b>df.isna()</b> devuelve True o False según si cada valor es nulo o no <b>df.isnull().sum()</b> o <b>df.isna().sum()</b> devuelve una serie con el número de valores nulos por columnas <b>df_% nulos = ((df.isnull().sum() / df.shape[0] * 100).reset_index())</b> <b>df_% nulos.columns</b> = ['columna', '% nulos'] crea un dataframe de los porcentajes de los valores nulos
<b>Crear series</b> <b>serie = pd.Series()</b> crear serie vacía <b>serie = pd.Series(array)</b> crear serie a partir de un array con el indice por defecto <b>serie = pd.Series(array, index = ['a', 'b', 'c'...])</b> crear una serie con indice definida; debe ser lista de la misma longitude del array <b>serie = pd.Series(lista)</b> crear una seria a partir de una lista <b>serie = pd.Series(número, indice)</b> crear una serie a partir de un escalar con la longitude igual al número de indices <b>serie = pd.Series(diccionario)</b> crear una serie a partir de un diccionario	<b>df = pd.DataFrame(array)</b> crear un dataframe a partir de un array con indices y columnas por defecto <b>df = pd.DataFrame(diccionario)</b> crear un dataframe a partir de un diccionario - los keys son los nombres de las columnas	<b>Eliminar filas duplicadas</b> <b>df.drop_duplicates(inplace = True, ignore_index=True)</b> elimina filas duplicadas; ignore_index para no tener el indice en cuenta	<b>pd.options.display.max_columns = None</b> ejecutar antes del df.head() para poder ver todas las columnas	<b>Eliminar nulos</b> <b>df.dropna(inplace = True, axis=b, subset=[lista_de_columnas], how=)</b> quitar nulos <b>how = 'any'   'all'</b> por defecto 'any': si hay algun valor NA, se elimina la fila o columna; all: si todos los valores son NA, se elimina la fila o columna <b>subset</b> una columna o lista de columnas
<b>Acceder a informacion de una serie</b> <b>serie.index</b> devuelve los indices <b>serie.values</b> devuelve los valores <b>serie.shape</b> devuelve la forma (no. filas) <b>serie.size</b> devuelve el tamaño <b>serie.dtypes</b> devuelve el tipo de dato	<b>DataFrames: carga de datos</b> <b>Carga de datos</b> <b>df = pd.read_csv("ruta/nombre_archivo.csv")</b> crear un dataframe de un archivo de Comma Separated Values <b>df = pd.read_csv("ruta/nombre_archivo", sep= ";")</b> crear un dataframe de un csv si el separador es ; <b>df = pd.read_csv("ruta/nombre_archivo", index_col= 0)</b> crear un dataframe de un csv si el archivo ya tiene una columna indice  <b>df = pd.read_excel("ruta/nombre_archivo.xlsx")</b> crear un dataframe de un archivo de Excel - si sale " <b>ImportError:... openpyxl...</b> ", en el terminal: <b>pip3 install openpyxl</b> o <b>pip install openpyxl</b>  <b>df = pd.read_json("ruta/nombre_archivo.json")</b> crear un dataframe de un archivo de JavaScript Object Notation (formato crudo) <b>df = df['data'].apply(pd.Series)</b> convertir el dataframe de json en un formato legible  <b>df = pd.read_clipboard(sep='\\t')</b> crear un dataframe de datos en forma de dataframe en el clipboard; el separador podria ser \\n ; , etc.	<b>Metodos de estadistica</b>  <b>df['columna'].mean()   mode()   median()   var()   std()</b> calcula la media/moda/mediana/variación/desviación estándar de los valores de una columna <b>df['columna1'].corr(df['columna2'])</b> calcula la correlacion entre dos variables <b>matriz_correlacion = df.corr()</b> crea una matriz mostrando las correlaciones entre todos los variables <b>df.crosstab = pd.crosstab(df['columna1'], df['columna2'], normalize = True, margins = True)</b> <b>normalize</b> muestra los valores en porcentajes (por uno) <b>margins</b> muestra los totales y subtotales <b>media_ponderada = np.average(df['columna'], weights = w)</b> calcula la media ponderada según los pesos <b>percentil_n = np.percentile(df['columna'], n)</b> saca el valor en el percentil n <b>q3, q1 = np.percentile(df["columna"], [75, 25])</b> saca los tercer y primer cuartiles	<b>pd.set_option("display.precision", 2)</b>	<b>Tipos de nulos</b> <b>np.nan</b> significa "not a number"; es un tipo numérico <b>None</b> valores nulos en columnas tipo string <b>NaT</b> valores nulos tipo datetime <b>valores texto:</b> "n/a", "NaN", "nan", "null" strings que normalmente se convierten automaticamente a np.nan <b>99999</b> o <b>00000</b> integers que se pueden convertir a nulos
<b>Operaciones con series</b> <b>serie1 +/-*/ serie2</b> suma/resta/multiplica/divide las filas con indices comunes entre las dos series <b>serie1.add(serie2, fill_value = número)</b> suma las filas con indices comunes, y suma el fill value a los valores sin indice comun <b>serie1.sub(serie2, fill_value = número)</b> restan las filas de la seria2 de la serie1 cuando tienen indices comunes, y resta el fill value de las otras indices de serie1 <b>serie1.mul(serie2, fill_value = número)</b> multiplica las filas con indices comunes y multiplica el fill value con las otras *usar 1 para conservar el valor* <b>serie1.mul(serie2, fill_value = número)</b> divida las filas de la serie1 entre las de la serie2 cuando tienen indices comunes, y divide las otras por el fill value <b>serie1.mod(serie2, fill_value = número)</b> devuelve el modulo (division sin resta) <b>serie1.pow(serie2, fill_value = número)</b> calcula el exponencial <b>serie1.ge(serie2)</b> compara si serie1 es <u>mayor</u> que serie2 y devuelve True o False <b>serie1.le(serie2)</b> compara si serie1 es <u>menor</u> que serie2 y devuelve True o False	<b>Guardado de datos</b> <b>df.to_csv('ruta/nombre_archivo.csv')</b> guardar dataframe como archivo csv <b>df.to_excel('ruta/nombre_archivo.xlsx')</b> guardar dataframe como archivo de Excel <b>df.to_json('ruta/nombre_archivo.json')</b> guardar dataframe como archivo de JSON <b>df.to_parquet('ruta/nombre_archivo.parquet')</b> guardar dataframe como archivo de parquet <b>df.to_pickle('ruta/nombre_archivo.pkl')</b> guardar dataframe como archivo de pickle  <b>ExcelWriter</b> <b>with pd.ExcelWriter("ruta/archivo.ext") as writer:</b> <b>df.to_Excel(writer, nombre_hoja = 'nombre')</b> guardar un dataframe en una hoja de Excel	<b>Sidetable: frecuencias de datos</b>  <b>df.stb.freq(['columna'])</b> devuelve un dataframe con informacion sobre la frecuencia de ocurrencia de cada categoría de un variable categorica parametros: <b>thresh = n</b> limita los valores mostrados a los más frecuentes hasta un umbral de n% cumulative y agrupando los restantes bajo la etiqueta "other" <b>other_label = 'etiqueta'</b> cambia la etiqueta 'other' <b>value = 'columna'</b> ordena los resultados por la columna especificada <b>df.stb.freq(['columna1', 'columna2'])</b> combina dos columnas y devuelve las frecuencias de las subcategories	<b>Outliers</b>  <b>Calcular tres desviaciones estandares:</b> <b>media = df.column.mean()</b> <b>desviacion = df.column.std()</b>  <b>lcb = media - desviacion * 3</b> <b>ucb = media + desviacion * 3</b>  <b>Eliminar Outliers</b> <b>outlier_step = 1.5 * IQR</b> calcular outlier step <b>outliers_data = df[(df['columna'] &lt; Q1 - outlier_step)   (df['columna'] &gt; Q3 + outlier_step)]</b> identificar datos fuera del rango del maximo hasta el minimo <b>lista_outliers_index = list(outliers_data.index)</b> crear una lista de los indices de las filas con outliers  <b>if outliers_data.shape[0] &gt; 0:</b> <b>dicc_indices[key] = (list(outliers_data.index))</b> crear un diccionario de los indices de las filas con nulos; se puede hacer iterando por columnas	<b>Reemplazar nulos</b> <b>df = pd.read_csv('archivo.csv', na_values = ['n/a'])</b> <b>.fillna(np.nan)</b> reemplaza los strings 'n/a' con np.nan al cargar el dataframe
<b>Filtrado booleanos</b> <b>serie &lt; &gt; &gt;= &lt;= == valor</b> devuelve True o False segun si cada condición cumple la condición <b>serie1[serie1 &lt; &gt; &gt;= &lt;= == valor]</b> devuelve solo los valores que cumplen la condición <b>np.nan</b> crear valor nulo (NaN) <b>serie.isnull()</b> devuelve True o False segun si los valores existen o son nulos ("" no cuenta como nulo) <b>serie.notnull()</b> devuelve True o False segun si los valores existen o son nulos ("" no cuenta como nulo)			<b>Reemplazar Outliers</b> <b>for k, v in dicc_indices.items():</b> <b>media = df[k].mean()</b> <b>for i in v:</b> <b>df.loc[i,k] = media</b> reemplazar <b>outliers por la media</b>	<b>Reemplazar nulos</b> <b>df.fillna(df[value=n, axis=b, inplace=True])</b> reemplazar todos los NaN del dataframe con el valor que especifiquemos <b>df['columna'].fillna(df['columna'].median, axis=b, inplace=True)</b> reemplazar los nulos de una columna por la mediana de esa columna <b>value=n</b> por defecto NaN; es el valor por el que queremos reemplazar los valores nulos que puede ser un escalar, diccionario, serie o dataframe <b>axis</b> por defecto 0 (filas) <b>df.replace(valor_nulo, valor_nuevo, inplace=True, regex=False)</b> reemplazar los nulos por el valor nuevo
				<b>Imputacion de nulos</b> <b>from sklearn.impute import SimpleImputer</b> <b>imputer = SimpleImputer(strategy='mean', missing_values = np.nan)</b> inicia la instancia del metodo, especificando que queremos reemplazar los nulos por la media <b>imputer = imputer.fit(df['columna1'])</b> aplicamos el imputer <b>df['media_columna1'] = imputer.transform(df[['price']])</b> rellena los valores nulos segun como hemos especificado <b>from sklearn.experimental import enable_iterative_imputer</b> <b>from sklearn.impute import IterativeImputer</b> <b>imputer = IterativeImputer(n_nearest_features=n, imputation_order='ascending')</b> crea la instancia <b>n_nearest_features</b> por defecto None; numero de columnas a utilizar para estimar los valores nulos <b>imputation_order</b> por defecto ascendente; el orden de imputacion <b>imputer.fit(df_numericas)</b> aplicamos el imputer <b>df_datos_trans = pd.DataFrame(imputer.transform(df_numericas), columns = numericas.columns)</b> crea un dataframe de los datos transformados; metemos estas columnas en el dataframe original <b>from sklearn.impute import KNNImputer</b> <b>imputerKNN = KNNImputer(n_neighbors=5)</b> crea la instancia <b>imputerKNN.fit(df_numericas)</b> <b>df_knn_imp = pd.DataFrame(imputerKNN.transform(df_numericas), columns = numericas.columns)</b> crea un dataframe de los datos transformados; metemos estas columnas en el dataframe original

Pandas
Union de datos
<div><div><div><div><div><div><span></span></div></div></div><div><div><div><span></span></div></div><div><div><span></span></div></div></div><div><div><div><span></span></div></div><div><div><span></span></div></div></div><div><div><div><span></span></div></div><div><div><span></span></div></div></div></div></div><div><div><div><div><span></span></div></div><div><div><span></span></div></div></div><div><div><div><span></span></div></div><div><div><span></span></div></div></div><div><div><div><span></span></div></div><div><div><span></span></div></div></div></div></div> <div><div><div><div><span></span></div></div><div><div><span></span></div></div></div><div><div><div><span></span></div></div><div><div><span></span></div></div></div><div><div><div><span></span></div></div><div><div><span></span></div></div></div></div>