

Python Cheat Sheet 2
Funciones
Definir una funcion: <code>def nombre_funcion(parametro1, parametro2, ...): return valor_del_return</code>
Lllamar una funcion: <code>nombre_funcion(argumento1, argumento2, ...)</code>
return: es opcional, pero sin return devuelve None parametros por defecto: – siempre deben ser lo ultimo
*args: una tupla de argumentos sin limite **kwargs: diccionarios cuyas keys se convierten en parámetros y sus valores en los argumentos de los parámetros
<code>def nombre_funcion(parametros, *args, **kwargs, parametro_por_defecto = valor) arg/kwarg: sin */** dentro de la funcion arg[0]</code>
Lllamar una funcion con *args: <code>nombre_funcion(argumento, argumento, argumento, ...) o nombre_funcion(*[lista_o_tupla_de_args])</code>
Lllamar una funcion con **kwargs: <code>nombre_funcion(**diccionario)</code>
Clases
Definir una clase: <code>class NombreClase: def __init__(self, atributo1, atributo2): self.atributo1 = atributo1 self.atributo2 = atributo2 self.atributo_por_defecto = ‘valor’ def nombre_funcion1(self, parametros) self.atributo += 1 return f“el nuevo valor es {self.atributo}”</code>
Definir una clase hija: <code>class NombreClaseHija(NombreClaseMadre): def __init__(self, atributo1, atributo2): super().__init__(atributo_hereditado1, ...) def nombre_funcion_hija (self, parametros):</code>
Crear un objeto de la clase: <code>variable_objeto = NombreClase(valor_atributo1, valor_atributo2)</code> instanciar (crear) un objeto <code>variable_objeto.atributo</code> devuelve el valor del atributo guardado para ese objeto <code>variable_objeto.atributo = nuevo_valor</code> para cambiar el valor del atributo <code>variable_objeto.nombre_funcion()</code> llamar una funcion
<code>print(help(NombreClase))</code> imprime informacion sobre la clase

Regex
- una abreviatura de `expresión regular`, `regex` es una cadena de texto que permite crear patrones que ayudana a emparejar, localizar y gestionar strings <code>import re</code> para poder trabajar con regex
Operadores communes de regex + coincide con el carácter precedente una o más veces * coincide con el carácter precedente cero o más veces u opcional ? indica cero o una ocurrencia del elemento precedente . coincide con cualquier carácter individual ^ coincide con la posición inicial de cualquier string \$ coincide con la posición final de cualquier string
Sintaxis básica de regex \\w cualquier caracter de tipo alfabético \\d cualquier caracter de tipo númeroico \\s espacios \\n saltos de línea \\w cualquier caracter que no sea una letra \\D cualquier caracter que no sea un dígitos \\S cualquier elemento que no sea un espacio () aísla sólo una parte de nuestro patrón de búsqueda que queremos devolver [] incluye todos los caracteres que queremos que coincidan e incluso incluye rangos como este: a-z y 0-9 es como el operador ‘or’ \\ señala una secuencia especial (escapar caracteres especiales) { } Exactamente el número especificado de ocurrencias {n} Exactamente n veces {n,} Al menos n veces {n,m} Entre n y m veces
Métodos Regex <code>re.findall(“patron”, string)</code> busca en todo el string y devuelve una lista con todas las coincidencias en nuestro string <code>re.search(“patron”, string_original)</code> busca en todo el string y devuelve un objeto con la primera coincidencia en nuestro string <code>re.match(“patron”, “string_original)</code> busca en la primera linea del string y devuelve un objeto con la primera coincidencia en nuestro string <code>resultado_match.span()</code> devuelve la referencia de las posiciones donde hizo el “match” <code>resultado_match.group()</code> devuelve el element resultando de la coincidencia del “match” <code>re.split(“patron”, “string_original”) </code> busca en todo el string y devuelve una lista con los elementos separados por el patron <code>re.sub(“patron”, “string_nuevo”, “string_original”) </code> busca en todo el string y devuelve un string con el element que coincide

Modulos/Librerias (paquetes de funciones)
Importar y usar modulos y sus funciones <code>import modulo</code> para importar un modulo <code>from modulo import funcion</code> importar solo una funcion <code>modulo.funcion()</code> usar una funcion de un modulo <code>modulo.clase.funcion()</code> para usar una funcion de una clase <code>import modulo as md</code> asignar un alias a un modulo
Libreria os <code>os.getcwd()</code> devuelve la ruta de donde estamos trabajando; se puede guardar en un variable e.g. ruta = os.getcwd() <code>os.listdir()</code> devuelve una lista de los archivos y carpetas donde estamos trabajando <code>os.listdir(‘carpeta’)</code> devuelve los contenidos de otra carpeta <code>os.chdir(‘ruta’)</code> cambia la carpeta en la que estes <code>os.mkdir(‘nueva_carpeta’)</code> crear una nueva carpeta <code>os.rename(‘nombre_carpeta’, ‘nueva_nombre’)</code> cambia el nombre de una carpeta <code>os.rmdir(‘carpeta’)</code> borra la carpeta
Libreria shutil from shutil import rmtree <code>rmtree(‘carpeta’)</code> borra la carpeta y subcarpetas
Abrir y cerrar ficheros Primero hay que guardar la ruta del archivo: ubicacion_carpeta = os.getcwd() nombre_archivo = “text.txt” ubicacion_archivo = ubicacion_carpeta + “/” + nombre_archivo <code>f = open(ubicacion_archivo)</code> abrir un archivo en variable f <code>f.close()</code> cerrar un archivo * IMPORTANTE * with open(ubicacion_archivo) as f: codigo e.g. variable = f.read() abre el archivo solo para ejecutar el codigo indicado (y despues lo deja)
Encoding from locale import getpreferredencoding getpreferredencoding() para saber que sistema de encoding estamos usando <code>f = open(ubicacion_archivo, encoding="utf-8")</code> abrir un archivo y leerlo con el encoding usado; guardar con .read()
mode: argumento opcional al abrir un archivo r – read w – write - sobreescribe x – exclusive creation, sólo crearlo si no existe todavía a – appending, añadir texto al archivo sin manipular el texto que ya habia hay que anadir otra letra: t – texto – leer en texto b – bytes – leer en bytes (no se puede usar con encoding)
<code>f = open(ubicacion_archivo, mode = “rt”)</code>
Leer ficheros <code>f.read()</code> leer el contenido de un archivo <code>f.read(n)</code> leer los primeros n caracteres de un archivo <code>variable = f.read()</code> guardar el contenido del archivo (o n caracteres de un archivo) en un variable <code>f.readline(n)</code> por defecto devuelve la primera linea o n lineas <code>f.readlines()</code> devuelve una lista de todas las lineas del archivo (cada linea es un elemento); se usa vacio sin n y list_name[x:] para seleccionar lineas especificas
Escribir en ficheros with open(ubicacion_archivo, “w”) as f: f.write(“Texto que va en el fichero.”) para escribir with open(ubicacion_archivo, “a”) as f: f.write(“Texto que va en el fichero.”) para anadir texto f.writelines(‘lista’) para anadir lineas de texto de una lista

Ficheros xml
<code>import xml.etree.ElementTree as ET</code> importa la librería xml <code>variable_tree = ET.parse(‘ruta/archivo.xml’)</code> abre el archivo <code>variable_root = variable_tree.getroot()</code> saca el elemento que envuelve todo (el elemento raíz) en una lista <root> <child_tag atributo1=“valor” atributo2=valor> <subchild_tag> elemento </subchild_tag> </child_tag> </root> <code>variable_root.tag</code> devuelve el nombre del tag del raiz <code>variable_root.attrib</code> devuelve los atributos del fichero <code>variable_root.find(“tag”).find(“childtag”).text</code> devuelve la primera ocasión en que el tag de un elemento coincida con el string <code>variable_root.findall(“tag”).findall(“childtag”).text</code> devuelve todos los elementos cuyos tag coincide
MySQL Connector/Python
Conectar a una base de datos <code>import mysql.connector</code> para importar MySQL Connector pip install mysql-connector pip install mysql-connector-Python <code>connect()</code> para conectar a una base de datos: <code>variable_cnx = mysql.connector.connect(user='root', password='AlumnaAdalab', host='127.0.0.1', database='nombre_BBDD')</code> <code>from mysql.connector import errorcode</code> importar errores <code>mysql.connector.Error</code> se puede usar en un try/except <code>cnx.close()</code> desconectar de la base de datos
Realizar queries <code>variable_cursor = cnx.cursor()</code> crear el objeto cursor que nos permite comunicar con la base de datos <code>variable_cursor.close()</code> desconectar el cursor <code>variable_query = (“SQL Query”)</code> guardar un query en un variable <code>variable_cursor.execute(variable_query)</code> ejecutar el query; devuelve una lista de tuplas <code>import datetime</code> sacar fechas en el formato AAAA-MM-DD <code>datetime.date(AAAA, M, D)</code> devuelve el formato de fecha <code>variable_query = “SQL Query... %s AND %s”</code> query dinamica <code>variable_cursor.execute(query, (variable1, variable2))</code> valores que van en lugar de los %s <code>variable_cursor.execute(“SHOW DATABASES”)</code> mostrar las BBDD <code>variable_cursor.execute(“SHOW TABLES”)</code> mostrar las tablas de la BBDD indicado en la conexión <code>variable_cursor.execute(“SHOW TABLES”)</code> <code>variable_cursor.execute(“SHOW COLUMNS FROM bbdd.table”)</code> mostrar las columnas de la tabla especificada; hay que conectarse a la bbdd information_schema
Argumentos cursor: <code>variable_cursor = cnx.cursor([arg=value[, arg=value]...])</code> <code>buffered=True</code> devuelve todas las filas de la bbdd <code>raw=True</code> el cursor no realizará las conversiones automáticas entre tipos de datos <code>dictionary=True</code> devuelve las filas como diccionarios <code>named_tuple=True</code> devuelve las filas como named tuples <code>cursor_class</code> un argumento que se puede usar para indicar que subclase queremos usar para instanciar el nuevo cursor

MySQL Connector/Python
Obtener resultados de una query <code>variable_cursor.fetchone()</code> devuelve el primer resultado <code>variable_cursor.fetchall()</code> devuelve todos los resultados como iterable – cada fila es una tupla
Pandas dataframe with SQL <code>import pandas as pd</code> <code>variable_df = pd.DataFrame(variable_resultado_fetchall, columns = [‘columna1’, ‘columna2’, ...])</code> crear un dataframe con los resultados de una query en una variable <code>variable_df.head(n)</code> devuelve las n primeras filas del df, o 5 por defecto <code>variable_df = pd.read_sql_query(variable_query, variable_cnx)</code> convertir los resultados de la query en df <code>pd.read_sql(variable_query, variable_cnx)</code> <code>variable_df.to_csv(“nombre_archivo.csv”)</code> guardar en csv <code>variable_df.to_string()</code> formatear el dato en string <code>variable_df.to_latex()</code> formatear el dato en un string que facilite la inserción en un documento latex
Crear y alterar una base de datos <code>variable_cursor.execute(“CREATE DATABASE nombre_BBDD”)</code> <code>variable_cursor.execute(“CREATE TABLE nombre_tabla (nombre_columna TIPO, nombre_columna2 TIPO2)”)</code> <code>variable_cursor.execute(“ALTER TABLE nombre_tabla ALTERACIONES”)</code>
Insertar datos <code>variable_query = “INSERT INTO nombre_tabla (columna1, columna2) VALUES (%s, %s)”</code> <code>variable_valores = (valor1, valor2)</code> <code>variable_cursor.execute(variable_query, variable_valores)</code> otro método: <code>variable_query = “UPDATE nombre_tabla SET nombre_columna = “nuevo_valor” WHERE nombre_columna = “valor”</code>
Insertar múltiples filas a una tabla <code>variable_valores_en_tuplas = ((valor1columna1, valor1columna2), (valor2columna1, valor2columna2))</code> <code>variable_cursor.executemany(variable_query, variable_valores_en_tuplas)</code>
<code>variable_conexion.commit()</code> después de ejecutar la inserción, para que los cambios efectúen en la BBDD <code>variable_conexion.rollback()</code> se puede usar después de execute y antes de commit para deshacer los cambios <code>print(variable_cursor.rowcount, “mensaje”)</code> imprimir el número de filas en las cuales se han tomado la accion
Eliminar registros <code>variable_query = “DROP TABLE nombre_tabla”</code>
Añadir errores importar errorcode y usar try/except: <code>try: accion except mysql.connector.Error as err: print(err) print("Error Code:", err.errno) print("SQLSTATE", err.sqlstate) print("Message", err.msg)</code>

Python Cheat Sheet 3
Pandas
Series: estructuras en una dimension
Crear series <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 escalar con la longitud igual al número de indices <code>serie = pd.Series(diccionario)</code> crear una serie a partir de un diccionario
Acceder a informacion de una serie <code>serie.index</code> devuelve los indices <code>serie.values</code> devuelve los valores <code>serie.shape</code> devuelve la forma (no. filas) <code>serie.size</code> devuelve el tamaño <code>serie.dtypes</code> devuelve el tipo de dato
<code>serie[i]</code> devuelve el valor del elemento en indice i <code>serie[[i,j]]</code> devuelve el valor de los dos elementos <code>serie[i:m]</code> devuelve el valor de un rango
<code>serie["etiqueta"]</code> devuelve el valor de los elementos en indices i y j
Operaciones con series <code>serie1 +-*/ serie2</code> suma/resta/multiplica/divide las filas con indices comunes entre las dos series <code>serie1.add(serie2, fill_value = número)</code> suma las filas con indices comunes, y suma el fill value a los valores sin indice comun <code>serie1.sub(serie2, fill_value = número)</code> restan las filas de la seria2 de la serie1 cuando tienen indices comunes, y resta el fill value de las otras indices de serie1 <code>serie1.mul(serie2, fill_value = número)</code> multiplica las filas con indices comunes y multiplica el fill value con las otras *usar 1 para conservar el valor* <code>serie1.mul(serie2, fill_value = número)</code> divida las filas de la serie1 entre las de la serie2 cuando tienen indices comunes, y divide las otras por el fill value <code>serie1.mod(serie2, fill_value = número)</code> devuelve el modulo (division sin resta) <code>serie1.pow(serie2, fill_value = número)</code> calcula el exponencial <code>serie1.ge(serie2)</code> compara si serie1 es <u>mayor</u> que serie2 y devuelve True o False <code>serie1.le(serie2)</code> compara si serie1 es <u>menor</u> que serie2 y devuelve True o False
Filtrado booleanos <code>serie < > >= <= == valor</code> devuelve True o False segun si cada condición cumple la condición <code>serie1[serie1 < > >= <= == valor]</code> devuelve solo los valores que cumplen la condición <code>np.nan</code> crear valor nulo (NaN) <code>serie.isnull()</code> devuelve True o False segun si los valores existen o son nulos ("" no cuenta como nulo) <code>serie.notnull()</code> devuelve True o False segun si los valores existen o son nulos ("" no cuenta como nulo)

DataFrames
Crear DataFrames <code>df = pd.DataFrame(data, index, columns)</code> <code>data</code> : NumPy Array, diccionario, lista de diccionarios <code>index</code> : indice que por defecto se asigna como 0-(n-1), n siendo el número de filas; <code>index = [lista]</code> para asignar “etiquetas” (nombres de filas) <code>column</code> : nombre de las columnas; por defecto 0-(n-1); <code>columns = [lista]</code> para poner mas nombres
<code>df = pd.DataFrame(array)</code> crear un dataframe a partir de un array con indices y columnas por defecto <code>df = pd.DataFrame(diccionario)</code> crear un dataframe a partir de un diccionario - los keys son los nombres de las columnas
DataFrames: carga de datos
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 <code>"ImportError:... openpyxl..."</code> , 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.
with open('ruta/nombre_archivo.pkl', 'wb') as f: pickle.dump(df,f) 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
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
ExcelWriter with pd.ExcelWriter("ruta/archivo.ext") as writer: df.to_excel(writer, nombre_hoja = 'nombre') guardar un dataframe en una hoja de Excel

Metodos de exploracion
<code>df.head(n)</code> devuelve las primeras n lineas del dataframe <code>df.tail(n)</code> devuelve las últimas n lineas del dataframe <code>df.sample(n)</code> devuelve n filas aleatorias de nuestro dataframe, o uno por defecto <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 de las columnas numéricas <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 <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.duplicated().sum()</code> devuelve el numero de filas duplicadas
Eliminar filas duplicadas <code>df.drop_duplicates(inplace = True)</code> elimina filas duplicadas
Metodos de estadistica
<code>df['columna'].mean()</code> <code>mode()</code> <code>median()</code> <code>var()</code> <code>std()</code> calcula la media/moda/mediana/variación/desviación estándar de los valores de una columna <code>df['columna1'].corr(df['columna2'])</code> calcula la correlacion entre dos variables <code>matriz_correlacion = df.corr()</code> crea una matriz mostrando las correlaciones entre todos los variables <code>df_crosstab = pd.crosstab(df['columna1'], df['columna2'], normalize = True, margins = True)</code> <code>normalize</code> muestra los valores en porcentajes (por uno) <code>margins</code> muestra los totales y subtotales <code>media_ponderada = np.average(df['columna'], weights = w)</code> calcula la media ponderada según los pesos <code>percentil_n = np.percentile(df['columna'], n)</code> saca el valor en el percentil n <code>q3, q1 = np.percentile(df["columna"], [75, 25])</code> saca los tercer y primer cuartiles
Sidetable: frecuencias de datos
<code>df.stb.freq(['columna'])</code> devuelve un dataframe con informacion sobre la frecuencia de ocurrencia de cada categoría de un variable categorica parametros: <code>thresh = n</code> limita los valores mostrados a los más frecuentes hasta un umbral de n% cumulative y agrupando los restantes bajo la etiqueta “other” <code>other_label = 'etiqueta'</code> cambia la etiqueta ‘other’ <code>value = 'columna'</code> ordena los resultados por la columna especificada <code>df.stb.freq(['columna1', 'columna2'])</code> combina dos columnas y devuelve las frecuencias de las subcategorias
<code>df.stb.missing(['columna'])</code> devuelve informacion sobre la frecuencia de datos nulos

Tipos de datos
Tipos de datos en Pandas: - object - int64 - float64 - datetime, timedelta[ns] - category - bool <code>df.dtypes</code> devuelve el tipo de datos que hay en cada columna <code>df_tipo = df.select_dtypes(include = "tipo")</code> 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 <code>copy = True</code> devuelve una copia <code>copy = False</code> *cuidado: los cambios en los valores pueden propagarse a otros objetos pandas* <code>errors = ignore</code> omite excepciones; en caso de error devuelve el objeto original <code>errors = raise</code> permite que se generen excepciones
<code>pd.options.display.max_columns = None</code> ejecutar antes del <code>df.head()</code> para poder ver todas las columnas
<code>pd.set_option("display.precision", 2)</code>

Outliers
Calcular tres desviaciones estandares: <code>media = df.column.mean()</code> <code>desviacion = df.column.std()</code>
<code>lcb = media - desviacion * 3</code> <code>ucb = media + desviacion * 3</code>
Eliminar Outliers <code>outlier_step = 1.5 * IQR</code> calcular outlier step <code>outliers_data = df[(df['columna'] < Q1 - outlier_step) (df['columna'] > Q3 + outlier_step)]</code> identificar datos fuera del rango del maximo hasta el minimo <code>lista_outliers_index = list(outliers_data.index)</code> crear una lista de los indices de las filas con outliers <code>if outliers_data.shape[0] > 0:</code> <code>dicc_indices[key] = (list(outliers_data.index))</code> crear un diccionario de los indices de las filas con nulos; se puede hacer iterando por columnas
<code>valores = dicc_indices.values()</code> sacar todos los valores e.g. todos los indices <code>valores = {indice for sublista in valores for indice in sublista}</code> set comprehension para eliminar duplicados <code>df_sin_outliers = df.drop(df.index[list(valores)])</code> crear nuevo dataframe sin outliers
Reemplazar Outliers <code>for k, v in dicc_indices.items():</code> <code>media = df[k].mean()</code> for i in v: <code>df.loc[i,k] = media</code> reemplazar outliers por la media

Valores nulos
Identificar nulos <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 una serie con el número de valores nulos por columnas <code>df_%_nulos = ((df.isnull().sum() / df.shape[0] * 100).reset_index())</code> <code>df_%_nulos.columns = ['columna', '% nulos']</code> crea un dataframe de los porcentajes de los valores nulos
Eliminar nulos <code>df.dropna(inplace = True, axis=b, subset=[lista_de_columnas], how=)</code> quitar nulos <code>how = 'any' 'all'</code> 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 <code>subset</code> una columna o lista de columnas
Tipos de nulos <code>np.nan</code> significa “not a number”; es un tipo numérico <code>None</code> valores nulos en columnas tipo string <code>NaT</code> valores nulos tipo datetime <code>valores texto: “n/a”, “NaN”, “nan”, “null”</code> strings que normalmente se convierten automaticamente a np.nan <code>99999</code> o <code>00000</code> integers que se pueden convertir a nulos
Reemplazar nulos <code>df = pd.read_csv('archivo.csv', na_values = ['n/a'])</code> <code>.fillna(np.nan)</code> reemplaza los strings ‘n/a’ con np.nan al cargar el dataframe
<code>df.fillna(df[value=n, axis=b, inplace=True])</code> reemplazar todos los NaN del dataframe con el valor que especifiquemos <code>df['columna'].fillna(df['columna'].median, axis=b, inplace=True)</code> reemplazar los nulos de una columna por la mediana de esa columna <code>value=n</code> por defecto NaN; es el valor por el que queremos reemplazar los valores nulos que puede ser un escalar, diccionario, serie o dataframe <code>axis</code> por defecto 0 (filas) <code>df.replace(valor_nulo, valor_nuevo, inplace=True, regex=False)</code> reemplazar los nulos por el valor nuevo
Imputacion de nulos from sklearn.impute import SimpleImputer <code>imputer = SimpleImputer(strategy='mean', missing_values = np.nan)</code> inicia la instancia del metodo, especificando que queremos reemplazar los nulos por la media <code>imputer = imputer.fit(df['columna1'])</code> aplicamos el imputer <code>df['media_columna1'] = imputer.transform(df[['price']])</code> rellena los valores nulos segun como hemos especificado from sklearn.experimental import enable_iterative_imputer from sklearn.impute import IterativeImputer <code>imputer = IterativeImputer(n_nearest_features=n, imputation_order='ascending')</code> crea la instancia <code>n_nearest_features</code> por defecto None; numero de columnas a utilizar para estimar los valores nulos <code>imputation_order</code> por defecto ascendente; el orden de imputacion <code>imputer.fit(df_numericas)</code> aplicamos el imputer <code>df_datos_trans = pd.DataFrame(imputer.transform(df_numericas), columns = df_numericas.columns)</code> crea un dataframe de los datos transformados; metemos estas columnas en el dataframe original from sklearn.impute import KNNImputer <code>imputerKNN = KNNImputer(n_neighbors=5)</code> crea la instancia <code>imputerKNN.fit(df_numericas)</code> <code>df_knn_imp = pd.DataFrame(imputerKNN.transform(df_numericas), columns = numericas.columns)</code> crea un dataframe de los datos transformados; metemos estas columnas en el dataframe original

Python Cheat Sheet 4
Pandas
Union de datos
<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>.concat() unir dataframes con columnas en comun df_union = pd.concat([df1, df2, df3], axis=b, join = ‘inner/outer’, ignore_index = True/False) parametros: axis = 0 une por columnas – los dataframes van uno encima del otro; las columnas tienen que ser de formatos compatible axis = 1 une por filas – los dataframes van uno al lado del otro; los datos deben ser relacionados para que tenga sentido join = ‘inner’ solo se quedan elementos que aparecen en todos los dataframes join = ‘outer’ se queda todo los datos de todos los dataframes ignore_index = True/False por defecto es False; si es True no usa las índices para la union (por ejemplo para union por el axis 0) .merge() unir las columnas de un dataframe a otro df_nuevo = df1.merge(df2, on = ‘columna’) inner merge df_nuevo = pd.merge(left = df1, right = df2, how=‘left’, left_on = ‘columna_df1’, right_on = ‘columna_df2’) left merge parametros: how = ‘left’ ‘right’ ‘outer’ ‘inner’ ‘cross’ on = columna [columna1, columna2, etc] si las columnas se llaman igual en los dos dataframes left_on = columna_df1 right_on = columna_df2 para especificar por donde hacer el merge suffixes = [‘left’, ‘right’] por defecto nada, el sufijo que apareciera en columnas duplicadas .join() unir dataframes por los indices df_nuevo = df1.join(df2, on = ‘columna’, how = ‘left’) inner merge parametros: how = ‘left’ ‘right’ ‘outer’ ‘inner’ por defecto left on = columna la columna o indice por el que queremos hacer el union; tienen que tener el mismo nombre en los dos dataframes lsuffix = ‘string’ rsuffix = ‘string’ por defecto nada, el sufijo que apareciera en columnas duplicadas</div>
Group By
<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>df_groupby = df.groupby(“columna_categoria”) crea un objeto DataFrameGroupBy; agrupa los valores segun las categorías de los valores de la columna indicada (o múltiples columnas en una lista) df_groupby.ngroups devuelve el numero de grupos df_groupby.groups devuelve un diccionario donde los keys son las categorías y los valores son listas de los índices de cada elemento en la categoría df_grupo1 = df_groupby.get_group(“grupo1”) devuelve un dataframe con los resultados de un grupo (la categoría indicada como grupo1) Cálculos con groupby: df_nuevo = df.groupby(“columna_categoria”).mean() devuelve un dataframe con la media de todas las columnas de valores numéricos, por categoría df_nuevo = df.groupby(“columna_categoria”)[“columna1”].mean() devuelve un dataframe con la media de la columna especificada count() número de observaciones no nulas describe() resumen de los principales estadísticos sum() suma de todos los valores mean() media de los valores df_nuevo = df.groupby(“columna_categoria”, dropna = False)[“columna_valores”].agg([nombre_columna = ‘estadistico1’, nombre_columna2 = ‘estadistico2’]) añade columnas con los cálculos de los estadísticos especificados dropna = False para tener en cuenta los Nan en los cálculos (por defecto es True)</div>

Subsets: loc e iloc
<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>df.loc[“etiqueta_fila”, “etiqueta_columna”] devuelve el contenido de un campo en una columna de una fila df.loc[“etiqueta_fila”,:] devuelve los valores de todas las columnas de una fila df.loc[:,“etiqueta_columna”] devuelve los valores de todas las filas de una columna df.iloc[indice_fila, indice_columna] devuelve el contenido de un campo en una columna de una fila df.iloc[indice_fila, :] devuelve los valores de todas las columnas de una fila df.iloc[:,indice_columna] devuelve el contenido de un campo en una columna de una fila df.loc[[lista_etiquetas_filas], [lista_etiquetas_columnas]] devuelve el contenido de varias filas / varias columnas df.loc[[lista_indices_filas], [lista_indices_columnas]] 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 df.loc[df.etiqueta > x] seleccionar datos basado en una condición usando operadores comparativos df.loc[(df.etiqueta > x) & (df.etiqueta == y)] seleccionar datos que tienen que cumplir las dos condiciones (and) df.loc[(df.etiqueta > x) (df.etiqueta == y)] seleccionar datos que tienen que deben cumplir una de las dos condiciones (or) df.iloc[list(df.etiqueta > x), :] iloc no acepta una Serie booleana; hay que convertirla en lista variable_df.head(n) devuelve las n primeras filas del df, o 5 por defecto</div>

Filtrados de datos

Filtrado por una columna con operadores de comparación
<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>df_filtrado = df[df[“nombre_columna”] == valor] extrae las filas donde el valor de la columna igual al valor dado</div>
Filtrado por multiples columnas con operadores logicos
<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>df_filtrado = df[(df[“columna1”] == valor) & (df[“columna2”] == valor) & (df[“columna3”] > n valor)] extrae las filas donde los valores de las columnas cumplan las condiciones en parentesis</div>
<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>df_filtrado = df[(df[“columna1”] == valor) (df[“columna1”] == valor) extrae las filas donde los valores de las columnas cumplan con una condición u otra</div>
<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>df_filtrado = ~(df[df[“columna1”] == valor]) extrae las filas donde los valores de las columnas NO cumplan con la condición</div>

Filtrados de datos
<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>Metodos de pandas de filtrar df_filtrado = df[df[“nombre_columna”].isin(iterable)] extrae las filas cuyas valores de la columna nombrada están en el iterable (una lista, serie, dataframe o diccionario) df_filtrado= df[df[“nombre_columna”].str.contains (patron, regex = True, na = False)] extrae las filas cuyas valores de la columna nombrada contienen el patron de regex df_filtrado = df[df[“nombre_columna”].str.contains (“substring”, case = False, regex = False)] extrae las filas cuyas valores de la columna nombrada contienen el substring, no siendo case sensitive df_filtrado = df[df[“nombre_columna”].str.contains (“substring”, case = False, regex = False)] extrae las filas cuyas valores de la columna nombrada contienen el substring, no siendo case sensitive df[pd.notnull(df[“nombre_columna”])] devuelve las filas que no tiene valores nulos en la columna especificada</div>

Cambiar columnas

<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>lista_columnas = df.columns.to_list() crea una lista de los nombres de las columnas del dataframe</div>
<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>df.set_index([“nombre_columna”, inplace = True) establece el índice utilizando uno o mas columnas; puede sustituir o ampliar un índice existente inplace = True los cambios sobrescriben sobre el df * cuando una columna se cambia a índice ya no es columna * df.reset_index(inplace = True) quitar una columna como indice para que vuelva a ser columna; crea un dataframe de una serie</div>

Renombrar columnas
<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>df.rename(columns = {“nombre_columna”: “nombre_nueva”}, inplace = True) 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} df.rename(columns = diccionario, inplace = True) cambia los nombres de las columnas según el diccionario</div>
Eliminar columnas
<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>df.drop(columns = [“columna1”, “columna2”], axis = b, inplace=True) eliminar una o mas columnas o filas segun lo que especificamos</div>

Reordenar columnas
<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>df = df.reindex(columns = lista_reordenada) cambia el orden de las columnas del dataframe segun el orden de la lista reordenada</div>

Crear columnas
<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>Creacion de ratios df[“columna?ratio”] = df.apply(lambda df: df[“columna1”] / df[“columna2”], axis = 1) Creacion de porcentajes def porcentaje(columna1, columna2): return (columna1 * 100) / columna2 df[“columna_%”] = df.apply(lambda df: porcentaje(df[“columna1”], datos[“columna2”]), axis = 1) df[“nueva_columna”] = np.where(df[“nombre_columna”] > n, “categoria_if_true”, “categoria_if_false”) crea una nueva columna basada en una condición df[“nueva_columna”] = np.select(lista_de_condiciones, lista_de_opciones) crea una nueva columna con los valores basados en multiples condiciones df[“columna_nueva”] = pd.cut(x = df[“nombre_columna”, bins = [n,m,l..], labels = [‘a’, ‘b’, ‘c’]) separa los elementos de un dataframe en diferentes intervalos (n-m, m-l, etc), creando una columna nueva que indica en cual intervalo cae el valor; con labels se puede asignar un string a cada intervalo Crear columnas df[“nueva_columna”] = (df[“etiqueta_columna”] + x) crea una nueva columna basada en otra df = df.assign(nueva_columna= df[“etiqueta_columna” + x]) crea una nueva basada en otra df = df.assign(nueva_columna= [lista_valores]) crea una nueva columna de una lista de valores *tiene que ser de la misma longitud como el número de filas del dataframe* df.insert(indice_nueva_columna, “nombre_columna”, valores) crea una nueva columna en la indice indicada allow_duplicates = True parametro cuando queremos permitir columnas duplicadas (por defecto es False)</div>

Apply

<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>apply() toma una función como argumento y la aplica a lo largo de un eje del DataFrame df[‘columna_nueva’] = df[‘col_1’].apply(función) crea una columna nueva con los valores de otra columna transformados según la función indicada df[‘columna_nueva’] = df[‘col_1’].apply(lambda x: x.método() if x > 1) crea una columna nueva con los valores de otra columna transformados según la lambda indicada df[‘columna_nueva’] = df.apply(lambda nombre: función(nombre[‘columna1’], nombre[‘columna2’]), axis = b) crea una columna nueva usando una función que coge dos parámetros (columna 1 y columna2) pd.applymap: se puede aplicar solo a todo el DataFrame. df.applymap(func, na_action=None, **kwargs) map solo lo podremos aplicar a una columa en particular. MiColumna.map(arg, na_action=None)</div>
--

apply() con datetime
<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>df[‘columna_fecha’] = df[‘columna_fecha’].apply(pd.to_datetime) cambia una columna de datos tipo fecha en el formato datetime def sacar_año(x): return x.strftime(“%Y”) df[‘columna_año’] = (df[‘columna_fecha’].apply(sacar_año) crea una columna nueva del año solo usando un método de la libreria datetime; (“%B”) para meses</div>

Cambiar valores
<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>Reemplazar valores basados en indices y condiciones: indices_filtrados = df.index[df[“columna”] == “valor”] for indice in indices_filtrados: df[“nombre_columna”].iloc[indice] = “valor_nuevo” Reemplazar valores basados en metodos NumPy: df.replace(to_replace = valor, value = valor_nuevo, inplace = True) reemplaza cierto valor por otro que especificamos df[“nombre_columna”.replace(to_replace = valor, value = valor_nuevo, inplace = True) reemplaza cierto valor en una columna por otro que especificamos df[[“columna1”, “columna2”]] = df[[“columna1”, “columna2”]].replace(r“string”, “string”, regex=True) cambiar un patron/string por otro en multiples columnas df[“nombre_columna”] = df[“nombre_columna”] + x reemplaza los valores de la columna por el valor + x (o otro valor que indicamos)</div>

datetime

<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>import datetime datetime.now() devuelve la fecha actual timedelta(n) representa una duración la diferencia entre dos instancias; n es un numero de días datetime.strptime(variable_fecha, ‘%Y-%m-%d’) formatea la fecha al formato indicado ayer = datetime.now() - timedelta(1) ayer = datetime.strptime(ayer, ‘%Y-%m-%d’) df[“fecha”] = ayer crea una columna con la fecha de ayer</div>

APIs

<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>import requests libreria para realizar petitions HTTP a una URL, para hacer web scraping url = ‘enlace’ el enlace de la que queremos extraer datos header = {} opcional; contiene informacion sobre las peticiones realizadas (tipo de ficheros, credenciales) response = requests.get(url=url, header = header) pedimos a la API que nos de los datos variables = {‘parametro1’:‘valor1’, ‘parametro2’:‘valor2’} response = request.get(url=url, params=variables) pedimos a la API que nos de los datos con los parametros segun el diccionario de parametros que le pasamos response.status_code devuelve el status de la peticion response.reason devuelve el motive de codigo de estado response.text devuelve los datos en formato string response.json() devuelve los datos en formato json df = pd.json_normalize(response.json) devuelve los datos en un dataframe</div>

Codigos de respuesta de HTTP
<div><div><div><div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div><div><div><div></div></div><div><div></div></div></div></div></div></div> <div>1XX informa de una respuesta correcta 2XX codigo de exito 200 OK 201 creado 202 aceptado 204 sin contenido 3XX redireccion</div> <div>4XX error durante peticion 401 peticion incorrecta 402 sin autorizacion 403 prohibido 404 no encontrado 5XX error del servidor 501 error interno del servidor 503 servicio no disponible</div>

NumPy (Numerical Python)	Indices, Subsets, Metodos de Arrays	Operaciones estadísticas y matemáticas	Funciones de conjuntos	Estadística
<h3>Crear arrays</h3> <h4>Crear arrays de listas</h4> <code>array = np.array(lista, dtype= tipo)</code> crea un array unidimensional de una lista <code>array = np.array([lista1, lista2])</code> crea un array bidimensional de dos listas <code>array = np.array([listadelistas1, listadelistas2])</code> crea un array bidimensional de dos listas	<h3>Indices de arrays</h3> <code>array[i]</code> devuelve la indice i; las indices de los arrays unidimensionales funcionan igual que las listas <code>array[i, j]</code> o <code>array[i][j]</code> devuelve el elemento de la columna j de la fila i <code>array[:,n]</code> seleccionar todas las filas y las columnas hasta n-1 <code>array[h, i, j]</code> o <code>array[h][i][j]</code> devuelve el elemento de la columna j de la fila i del array h <code>array[h][i][j] = n</code> cambiar el valor del elemento en esta posicion al valor n <h3>Subsets</h3> <code>array > n</code> devuelve la forma del array con True o False según si el elemento cumple con la condición o no <code>array[array > n]</code> devuelve un subset: todos los valores que cumplen la condición en una lista dentro de un array <code>array[(array > n) & (array < m)]</code> devuelve un subset: todos los valores que cumplen las condiciones en una lista dentro de un array; se puede usar para “or” <h3>Metodos de arrays</h3> <code>nuevo_array = array.copy()</code> crea un a copia del array <code>np.transpose(array_bidimensional)</code> cambia los filas del array a columnas y las columnas a filas <code>np.transpose(array_multidimensional)</code> cambia el número de columnas al número de arrays y viceversa; el número de filas no cambia <code>np.transpose(array_multidimensional, (z,y,x))</code> hace la transposicion segun lo que especifecemos usando las posiciones de la tupla (0,1,2) de la forma original <code>array = np.arange(n).reshape((y,x))</code> crea un array usando reshape para definir la forma <code>array = np.reshape(array, (z,y,x))</code> crea un array con los valores de otro array usando reshape para definir la forma <code>array = np.swapaxes(array, posicion, posicion)</code> intercambia dos ejes de una matriz usando las posiciones (z=0,y=1,x=2) de la forma original	<h3>Operaciones estadísticas y matemáticas</h3> El parametro axis en arrays bidimensionales: <code>axis = 0</code> columnas <code>axis = 1</code> filas - si especificamos el axis, la operación devuelve el resultado por cada fila o columna. Por ejemplo: <code>np.sum(array, axis = 0)</code> devuelve un array con la suma de cada fila El parametro axis en arrays multidimensionales: <code>axis = 0</code> dimensión <code>axis = 1</code> columnas <code>axis = 2</code> filas - si especificamos el axis, la operación devuelve el resultado por cada dimensión, fila o columna. Por ejemplo: <code>np.sum(array_3D, axis = 0)</code> devuelve un array de una matriz con la suma de todas las matrices <code>np.sum(array_3D, axis = 1)</code> devuelve un array donde las filas contienen las sumas de las columnas de cada matriz Operaciones con parámetro del axis: <code>np.sum(array_3D)</code> devuelve la suma de todos los elementos de los matrices <code>np.mean(array)</code> devuelve la media de todo el array <code>np.std(array)</code> devuelve la desviación estándar de todo <code>np.var(array)</code> devuelve la varianza de valores de todo <code>np.min(array)</code> devuelve el valor mínimo del array <code>np.max(array)</code> devuelve el valor máximo del array <code>np.sum(array)</code> devuelve la suma de los elementos del array <code>np.cumsum(array)</code> devuelve un array con la suma acumulada de los elementos a lo largo del array <code>np.cumprod(array)</code> devuelve un array con la multiplicación acumulada de los elementos a lo largo del array Operaciones sin parámetro del axis: <code>np.sqrt(array)</code> devuelve un array con la raíz cuadrada no negativa de cada elemento del array <code>np.exp(array)</code> devuelve un array con el exponencial de cada elemento del array <code>np.mod(array1, array2)</code> devuelve un array con el resto de la división entre dos arrays <code>np.mod(array1, n)</code> devuelve un array con el resto de la división entre el array y el valor de n <code>np.cos(array)</code> devuelve un array con el coseno de cada elemento del array <code>np.sin(array)</code> devuelve un array con el seno de cada elemento del array <code>np.sin(array)</code> devuelve un array con la tangente de cada elemento del array	<code>np.unique(array)</code> devuelve un array con los valores únicos del array ordenados <code>np.unique(array, return_index=True)</code> devuelve un array con los valores únicos del array ordenados y un array con la posición de la primera instancia de cada valor <code>np.unique(array, return_inverse=True)</code> devuelve un array con los valores únicos del array ordenados y un array con las posiciones de cada elemento de cada valor <code>np.unique(array, return_counts=True)</code> devuelve un array con los valores únicos del array ordenados y un array con el número de veces que aparece cada valor <code>np.unique(array, axis = b)</code> devuelve un array con los valores únicos ordenados de las filas o columnas Funciones para arrays unidimensionales <code>np.intersect1d(array1, array2)</code> devuelve un array con los valores únicos de los elementos en común de dos arrays <code>np.intersect1d(array1, array2, return_indices=True)</code> devuelve un array con los valores únicos de los elementos en común de dos arrays y arrays con los índices de cada valor, por array <code>np.union1d(array1, array2)</code> devuelve un array ordenado con los elementos resultantes de unir dos arrays (valores únicos) <code>np.in1d(array1, array2)</code> devuelve un array con True o False por cada elemento de array1 según si aparece el mismo valor en array2 <code>np.setdiff1d(array1, array2)</code> devuelve un array ordenado con los valores únicos que están en array1 pero no en array2 <code>np.setxor1d(array1, array2)</code> devuelve un array ordenado con los valores únicos que NO están en común de los dos arrays	<h3>Tablas de frecuencias</h3> Frecuencias absolutas el número de veces que se repite un número en un conjunto de datos <code>df = df.groupby('columna').count().reset_index()</code> Frecuencias relativas las veces que se repite un número o categoría en un conjunto de datos respecto al total, en porcentajes <code>df_group_sin_str = df_group.drop('columna_str', axis=1)</code> <code>frecuencia_relativa = df_group_sin_str / df.shape[0] * 100</code> <code>columnas = df_group_sin_strings.columns</code> <code>df_group[columnas] = frecuencia_relativa</code> Tablas de contingencia tabla de frecuencias que cuenta todas las combinaciones posibles de cada pareja de valores de las columnas que estamos intentando comparar <code>df_crosstab = pd.crosstab(df['columna1'], df['columna2'], normalize = True, margins = True)</code> <code>normalize</code> muestra los valores en porcentajes (por uno) <code>margins</code> muestra los totales y subtotales Coefficiente de correlación de Pearson - nos permite conocer la intensidad y dirección de la relación entre las dos variables - coeficiente > 0: correlación positiva - coeficiente < 0: correlación negativa - coeficiente = 1 o -1: correlación total - coeficiente = 0: no existe relación lineal <code>df['columna1'].corr(df['columna2'])</code> calcula la correlacion entre dos variables <code>matriz_correlacion = df.corr()</code> crea una matriz mostrando las correlaciones entre todos los variables <code>sns.heatmap(df.corr()[['column1', 'column2']], cmap = 'color_palette', annot = True, vmin = -1, vmax = 1)</code> crea una grafica heatmap de la matriz de correlaciones
<h3>NumPy Random</h3> <code>np.random.seed(x)</code> establece la semilla aleatoria del generador de números aleatorios, para que las funciones random que van después siempre cogerán los mismos valores “aleatorios” Crear arrays con valores aleatorios <code>array = np.random.randint(inicio, final, forma_matriz)</code> crea un array de números aleatorios entre dos valores; forma_matriz: (z,y,x) z: número de arrays y: número de filas x: número de columnas <code>array = np.random.randint(inicio, final)</code> devuelve un número aleatorio en el rango <code>array = np.random.rand(z,y,x)</code> crea un array de floats aleatorias con la forma que le especificemos; por defecto genera números aleatorios entre 0-1 <code>array = np.random.random_sample((z,y,x))</code> crea un array de floats aleatorias con la forma que le especificemos; por defecto genera números aleatorios entre 0-0.9999999... <code>array = np.random.z,y,x=None)</code> devuelve un número aleatorio en 0 y 0.9999999999999... <code>np.round(np.random.rand(z,y,x), n)</code> crear array con floats de n decimales <code>np.random.uniform(n,m, size = (z,y,x))</code> genera muestras aleatorias de una distribución uniforme en el intervalo entre n y m <code>np.random.binomial(n,m, size = (z,y,x))</code> genera muestras con una distribución binomial; n es el numero total de pruebas; m es la probabilidad de éxito <code>np.random.normal(loc = n, scale = m, size = (z,y,x))</code> genera números aleatorios de una distribución normal (curva de campana); loc es la media; scale es la desviación estándar <code>np.random.permutation(array)</code> devuelve un array con los mismos valores mezclados aleatoriamente	<h3>Otras operaciones</h3> <code>np.sort(array)</code> devuelve un array con los valores de cada fila ordenados en orden ascendente por defecto <code>np.sort(array, axis = 0)</code> devuelve un array con los valores de cada columna ordenados en orden ascendente <code>np.sort(-array)</code> devuelve un array con los valores de cada fila ordenados en orden descendente <code>np.round(array, decimals = x)</code> devuelve un array con los valores del array redondeados a x decimales <code>np.round(array, decimals = x)</code> devuelve un array con los valores del array redondeados a x decimales <code>np.where(array > x)</code> devuelve los indices de los valores que cumplan la condición, por fila y columna	<h3>Operaciones de comparación en arrays bidimensionales</h3> <code>np.any(array > n)</code> devuelve True o False segun si cualquier valor del array cumpla con la condicion <code>np.any(array > n, axis = b)</code> devuelve un array con True o False por cada columna o fila según si algún valor de la fila o columna cumpla con la condición <code>np.all(array > n)</code> devuelve True o False segun si todos los valores del array cumpla con la condicion <code>np.all(array > n, axis = b)</code> devuelve un array con True o False por cada columna o fila según si todos los valores de la fila o columna cumplan con la condición	Medidas de dispersión Desviación respecto a la media la diferencia en valor absoluto entre cada valor de los datos y su media aritmética <code>diferencias = df['columna'] - df['columna'].mean()</code> <code>desviación_media = np.abs(diferencias)</code> Varianza medida de dispersión; la variabilidad respecto a la media <code>df['columna'].var()</code> Desviación estándar o desviación típica la raíz cuadrada de la varianza; cuanto mayor sea, mayor será la dispersión o variabilidad en nuestros datos <code>df['columna'].std()</code> Robustez - cuanto más cantidad de datos, más robustos 1/n donde n es el numero de registros Coefficiente de variación el cociente entre la desviación típica y la media; cuanto mayor sea, mayor será la dispersión en nuestros datos <code>df['columna'].std() / df['columna'].mean()</code> Percentiles divide datos ordenados de menor a mayor en cien partes; muestra la proporción de datos por debajo de su valor <code>percentil_n = np.percentile(df['columna'], n)</code> saca el valor en el percentil n Rangos intercuartílicos medida de dispersión: diferencia entre cuartiles 75 y 25 <code>q3, q1 = np.percentile(df[“columna”], [75, 25])</code> saca los tercer y primer cuartiles <code>rango_intercuartílico = q3 - q1</code>	Intervalos de confianza describe la variabilidad entre la medida obtenida en un estudio y la medida real de la población (el valor real) <code>import scipy.stats as st</code> <code>st.t.interval(alpha = n, df = len(df['columna'])-1, loc = np.mean(df['columna']), scale = st.sem(df['columna']))</code> devuelve el rango de valores para lo cual hay un n% de probabilidad que un valor real cae en ese rango alpha: porcentaje de confianza (p.ej. 90%, 95%, o 99%) df: los datos loc: la media scale: la desviación estándar

Estadística
Estadísticos generales
mínimo el valor mas bajo dentro de un conjunto de datos
máximo el valor mas alto dentro de un conjunto de datos
Medidas de tendencia central
Media el valor obtenido al sumar todos los datos y dividir el resultado entre el número total de elementos
df[‘columna’].mean()
Media ponderada el valor obtenido al multiplicar cada uno de los datos por un peso antes de sumarlos (suma ponderada); después se divide la suma ponderada entre la suma de los pesos
media_ponderada = np.average(df[‘columna’], weights = w)
Mediana el valor de la variable de posición central en un conjunto de datos
df[‘columna’].median()
Moda el valor que tiene mayor frecuencia absoluta de entre todos los datos
df[‘columna’].mode()

Estadística
Medidas de dispersión
Desviación respecto a la media la diferencia en valor absoluto entre cada valor de los datos y su media aritmética
diferencias = df[‘columna’] - df[‘columna’].mean()
desviación_media = np.abs(diferencias)
Varianza una medida de dispersión que representa la variabilidad de una serie de datos respecto a su media
df[‘columna’].var()
Desviación estándar o desviación típica la raíz cuadrada de la varianza; nos dice como de dispersos están los valores de los datos en un conjunto de datos; cuanto mayor sea, mayor será la dispersión o variabilidad en nuestros datos
df[‘columna’].std()
Robustez - normalmente, cuanto más cantidad de datos hayamos usado para calcular los estadísticos, más robustos serán estos, ya que la influencia de unos pocos datos inusuales sobre el total de los mismos será menor - punto de ruptura: la fracción de los datos a los que podríamos dar valores arbitrarios sin hacer que el estadístico se vea tan afectado como para no ser útil
1/n donde n es el numero de registros
Coeficiente de variación muestra la relación entre el tamaño de la media y la variabilidad de la variable; se calcula como el cociente entre la desviación típica y la media; cuanto mayor sea, mayor será la dispersión o variabilidad en nuestros datos
df[‘columna’].std() / df[‘columna’].mean()
Percentiles divide una serie de datos ordenados de menor a mayor en cien partes iguales; se trata de un indicador que busca mostrar la proporción de la serie de datos que queda por debajo de su valor
percentil_n = np.percentile(df[‘columna’], n) saca el valor en el percentil n
percentil 25 → Q1 percentil 50 → Q2 (la mediana)
percentil 75 → Q3
Rangos intercuartílicos una medida de dispersión estadística igual a la diferencia entre los cuartiles 75 y 25
q3, q1 = np.percentile(df[“columna”], [75, 25]) saca los tercer y primer cuartiles
rango_intercuartílico = q3 - q1
Tablas de frecuencias
Frecuencias absolutas son el número de veces que se repite un número en un conjunto de datos
df_group = df.groupby(‘columna’).count().reset_index()
Frecuencias relativas las veces que se repite un número o categoría en un conjunto de datos respecto al total, en porcentajes
df_group_sin_str = df_group.drop(‘columna_str’, axis=1)
frecuencia_relativa = df_group_sin_str / df.shape[0] * 100
columnas = df_group_sin_strings.columns
df_group[columnas] = frecuencia_relativa

Estadística
Tablas de contingencia - una tabla de frecuencias que cuenta todas las combinaciones posibles de cada pareja de valores de las columnas que estamos intentando comparar - deberá contener al menos dos filas y dos columnas para representar datos categóricos - permite medir la interacción entre las variables
df_crosstab = pd.crosstab(df[‘columna1’], df[‘columna2’], normalize = True, margins = True)
normalize muestra los valores en porcentajes (por uno)
margins muestra los totales y subtotales
Coeficiente de correlación de Pearson - nos permite conocer la intensidad y dirección de la relación entre las dos variables - coeficiente > 0: correlación positiva (según aumente el valor de una variable aumenta el valor de la otra) - coeficiente < 0: correlación negativa (según aumente el valor de una variable disminuye el valor de la otra) - coeficiente = 1 o -1: correlación total (positiva o negativa) - coeficiente = 0: no existe relación lineal entre las dos variables, pero puede haber algún otro tipo de relación que no se capture bien con este coeficiente
df[‘columna1’].corr(df[‘columna2’]) calcula la correlacion entre dos variables
matriz_correlacion = df.corr() crea una matriz mostrando las correlaciones entre todos los variables
Heatmap
sns.heatmap(df.corr(), cmap = ‘color_palette’, annot = True, vmin = -1, vmax = 1) crea un heatmap con una escala de colores que refleja los valores de correlacion
annot = True para que aparezcan los valores
vmin/vmax establecen la escala de color
sns.heatmap(df.corr()[[‘column1’, ‘column2’]], cmap = ‘color_palette’, annot = True, vmin = -1, vmax = 1)
Sesgos (skewness) - una medida de la asimetría de la distribución de los valores de una variable alrededor de su valor medio - valor de sesgo positivo: sesgado a la derecha - valor de sesgo negativo: sesgado a la izquierda - valor de sesgo igual a 0: valores simetricos
sns.displot(df[‘columna’], kde = True) crea un histograma que muestra la distribution de los valores
import scipy.stats import skew skew(df[‘columna’]) muestra el valor del sesgo de una variable
Intervalos de confianza describe la variabilidad entre la medida obtenida en un estudio y la medida real de la población (el valor real)
Confidence Interval = x +/- t*(s/√n)
where: x: sample mean t: t-value that corresponds to the confidence level s: sample standard deviation n: sample size
import scipy.stats as st
st.t.interval(alpha = n, df = len(df[‘columna’])-1, loc = np.mean(df[‘columna’]), scale = st.sem(df[‘columna’])) devuelve el rango de valores para lo cual hay un n% de probabilidad que un valor real cae en ese rango
alpha: porcentaje de confianza (p.ej. 90%, 95%, o 99%)
df: los datos
loc: la media
scale: la desviación estándar