

## Python Cheat Sheet 1

### Variables ampliadas por text (CONCATENATION)

#### Para encadenar texto

```
categoria1 = "verde"
color_detalle = categoria1 + ' ' + 'oscuro'

print(categoria1 + ' oscuro')
print(categoria1, 'oscuro')
```

### type() and isinstance()

**float/int/str(variable)** cambia el tipo de data/type

**type(variable)** devuelve: class 'float/int/str'

**isinstance(variable, float/int/str)** comprobar el tipo de dato (devuelve True/False)

### Operaciones Algebraicas

+ sumar	/ dividir
- restar	// divider y redondear (modulus)
* multiplicar	% resto de una division (floor division)
** elevar	round(x) redondear número x

### Operaciones Binarias

**==** comprobar si valores coinciden  
**is** comprobar si valores son exacamente igual  
**!=** comprobar si valores son diferentes  
**is not** comprobar si valores no son exactamente iguales  
**> (>=)** mayor que (mayor o igual que)  
**< (<=)** menor que (menor o igual que)  
**and** ambas verdaderas  
**or** ambas o solo una verdadera  
**in/not in** comprobar si hay un valor en una lista etc.

### Metodos String

**string.upper()** MAYUSCULAS  
**string.lower()** minusculas  
**string.capitalize()** Primera letra de la frase en may.  
**string.title()** Primera Letra De Cada Palabra En May.  
**string.swapcase()** mINUSCULAS A mAYUSCULAS O vICEVERSA  
**string.strip()** quita espacios del principio y final

**string.split()** divide string en lista - por espacios por defecto, o especifica otro divisor en ()  
**string.replace("frase", "frase")** reemplaza la primera frase del string por el otro  
**" ".join(string)** une los elementos de una lista en una string con el separador espificado en " "  
**list(string)** convierte un variable string en una lista  
**string.find("substring")** encuentra el indice en que empiece el substring/'-1' si no existe el substring

**string[i]** devuelve el elemento en la indice i  
**string[i:j]** devuelve un rango de caracteres

# metodos permanentes (cambia el variable, no devuelve nada)

### Listas [ ] Metodos no permanentes

**lista = []** crea una lista vacia  
**len(lista)** devuelve el no. de elementos  
**min(lista)/max(lista)** saca el valor minimo y maximo  
**lista.count()** devuelve el no. de elementos que hay en la lista de un valor determinado en los()  
**sorted(lista)** ordenar una lista de menor a mayor  
**lista.copy()** hacer una copia de la lista

#### Metodos con indices

**list.index(x)** devuelve la indice de x en la lista  
**lista[i]** devuelve el elemento en la indice i  
**[start:stop:step]**  
**lista[i:j:x]** devuelve los elementos por el rango de i a j (incluye i pero no j) saltando por x  
**lista[-i:-j]** devuelve los elementos por los indices negativos (incluye -j pero no -i)

### Listas – Acciones Permanentes

#### Ampliar una lista

[lista1, lista2] junta listas pero se mantienen como listas separadas  
lista1 + lista2 hace una lista mas larga

#### .append()

**lista.append(x)#** añade un solo elemento (lista, string, integer o tuple) a la lista

#### .extend()

**lista.extend(lista2)#** añade los elementos de una lista al final de la lista

#### .insert()

**.insert(i, x)#** mete un elemento (x) en un índice(i)

#### Ordenar una lista

#### .sort()

**lista.sort()**# ordena de menor a mayor, usar con (reverse=True) para ordenar de mayor a menor  
**lista.reverse()**# ordena los elementos al revés del orden guardado

#### Quitar elementos de una lista

#### .pop()

**lista.pop(i)#** quita el elemento en indice i y devuelve su valor

#### .remove()

**lista.remove(x)#** quita el primer elemento de la lista con valor x

**lista.clear()**# vacia la lista

**del lista#** borra la lista  
**del lista[i]**# borra el elemento en indice i

### Diccionarios { key : value , }

**diccionario = {x:y}** compuestos por un key(x) unica y un valor(y) (cualquier tipo de datos)  
**dict()**  
**variable = dict(x=y, m=n)** crear un diccionario  
**dicc.copy()** crear una copia  
**len(dicc)** devuelve el no. de elementos (x:y) hay en el diccionario  
**sorted(dicc)** ordena los keys; usar con .items() para ordenar tuplas de los elementos o .values() para ordenar los valores solos

### Diccionarios – Metodos

#### Obtener informacion de un diccionario

**dicc.keys()** devuelve todas las keys  
**dicc.values()** devuelve todos los valores  
**dicc.items()** devuelve tuplas de los key:value  
**in/not in** comprobar si existe una clave  
**dicc.get(x, y)** devuelve el valor asociado al key x, o si no existe devuelve el output y  
**dicc["key"]** devuelve el valor del key (ver abajo que tiene mas usos)

#### Ampliar un diccionario

**.update()**  
**dicc.update({x:y})#** para insertar nuevos elementos  
**dicc["key"] = valor#** para insertar un nuevo key o valor, o cambiar el valor de un key  
**dicc.setdefault(x, y)#** devuelve el value del key x, o si no existe la key x, la crea y asigna el valor y por defecto

#### Quitar elementos de un diccionario

**dicc.pop(x)#** elimina la key x (y lo devuelve)  
**dicc.popitem()**# elimina el ultimo par de key:value  
**dicc.clear()**# vacia el diccionario

### Tuplas (,) inmutables, indexados

tupla = (x,y) tuplas se definen con () y , o solo ,  
tupla1 + tupla2 juntar tuplas

**tuple(lista)** crear tuplas de una lista  
**tuple(dicc)** crear tuplas de los keys de un diccionario  
**tuple(dicc.values())** crear tuplas de los valores  
**tuple(dicc.items())** crear tuplas de los key:valores

**len(tupla)** devuelve el no. de elementos  
**in/not in** comprobar si hay un elemento  
**tupla.index(x)** devuelve el indice de x  
**tupla.count(x)** devuelve el no. de elementos con valor x en la tupla

\*para cambiar el contenido de una tupla hay que convertirla en una lista y luego a tupla\*

### zip()

**zip(iterable1, iterable2)** crea una lista de tuplas de parejas de los elementos de las dos listas (mientras se puede)  
**listzip.sort()** ordena las tuplas del zip por el primer elemento

### Sets { }

#### no permiten duplicados, no tienen orden

**set = {x,y}**  
**set(iterable)** solo permite un argumento iterable; elimina duplicados

**in/not in** comprobar si hay un elemento  
**len(set)** devuelve el no. de elementos

#### Ampliar un set

**set.add(x)#** añadir un elemento  
**set.update(set o lista)#** añadir uno o mas elementos con [] o {} o un variable tipo lista o set

#### Quitar elementos de un set

**set.pop()**# elimina un elemento al azar  
**set.remove(x)#** elimina el elemento x  
**set.discard(x)#** elimina el elemento x (y no devuelve error si no existe)  
**set.clear()**# vacia el set

#### Operaciones con dos Sets

**set1.union(set2)** devuelve la union de los dos sets: todos los elementos menos dupl.  
**set1.intersection(set2)** devuelve los elementos comunes de los dos sets  
**set1.difference(set2)** devuelve los sets que estan en set1 pero no en set2 (restar)  
**set1.symmetric\_difference(set2)** devuelve todos los elementos que no estan en ambos  
**set1.isdisjoint(set2)** comprobar si todos los elementos de dos sets son diferentes  
**set1.issubset(set2)** comprobar si todos los elementos de set1 estan en set2  
**set1.superset(set2)** comprobar si todos los elementos de set2 estan en set1

### input()

• permite obtener texto escrito por teclado del usuario  
**input("el texto que quieres mostrar al usuario")**  
• se puede guardar en un variable  
• por defecto se guarda como un string  
**x = int(input("escribe un número"))** para usar el variable como integer o float se puede convertir en el variable

### Sentencias de control

**if ... elif ... else**  
**if** estableca una condición para que se ejecute el código que esta debajo del if. \*tiene que estar indentado\*  
**elif** para chequear mas condiciones después de un if  
**else** agrupa las condiciones que no se han cumplido; no puede llevar condiciones nuevas

if x > y:  
    print("x es mayor que y")  
elif x == y:  
    print("x es igual que y")  
else:  
    print("x e y son iguales")

**while**  
• repite el código mientras la condición sea True, o sea se parará cuando la condición sea False  
• se pueden incluir condiciones con if... elif... else  
• \*pueden ser infinitos\* (si la condición no llega a ser False)

while x < 5:  
    print("x es mayor que 5")

### For loops

• sirven para iterar por todos los elementos de un variable que tiene que ser un iterable (lista, diccionario, tupla, set, or string)  
• se pueden combinar con if ... elif ... else, while, u otro for loop  
• en diccionarios por defecto intera por las keys; podemos usar dicc.values() para acceder a los valores  
for i in lista:  
    print("hola mundo")

### List comprehension

• su principal uso es para crear una lista nueva de un un for loop en una sola línea de código  
[ **lo que queremos obtener** **iterable** **condición** (opcional) ]

### try ... except

Se usan para evitar que nuestro código se pare debido a un error en el código. Se puede imprimir un mensaje que avisa del error.  
**try:**  
    print("2.split())  
**except:**  
    print("no funciona")

### range()

• nos devuelve una lista de números que por defecto se aumentan de uno en uno empezando por 0  
**range(start:stop:step)**  
• se puede especificar por donde empieza y el limite (que debe ser +1 por que se para uno antes del limite que ponemos como stop)  
• tambien se puede especificar saltos

Python Cheat Sheet 2
Funciones
<b>Definir una funcion:</b> <code>def nombre_funcion(parametro1, parametro2, ...):     return valor_del_return</code>
<b>Lllamar una funcion:</b> <code>nombre_funcion(argumento1, argumento2, ...)</code>
<b>return:</b> es opcional, pero sin return devuelve None parametros por defecto: – siempre deben ser lo ultimo
<b>*args:</b> una tupla de argumentos sin limite <b>**kwargs:</b> 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>
<b>Lllamar una funcion con *args:</b> <code>nombre_funcion(argumento, argumento, argumento, ...)</code> o <code>nombre_funcion(*[lista_o_tupla_de_args])</code>
<b>Lllamar una funcion con **kwargs:</b> <code>nombre_funcion(**diccionario)</code>
Clases
<b>Definir una clase:</b> <code>class NombreClase:</code>
<code>    def __init__(self, atributo1, atributo2):         self.atributo1 = atributo1         self.atributo2 = atributo2         self.atributo_por_defecto = ‘valor’</code>
<code>    def nombre_funcion1(self, parametros)         self.atributo += 1         return f“el nuevo valor es {self.atributo}”</code>
<b>Definir una clase hija:</b> <code>class NombreClaseHija(NombreClaseMadre):     def __init__(self, atributo1, atributo2):         super().__init__(atributo_hereditado1, ...)</code>
<code>    def nombre_funcion_hija (self, parametros):</code>
<b>Crear un objeto de la clase:</b> <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
<b>Operadores communes de regex</b> + 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
<b>Sintaxis básica de regex</b> \\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
<b>Métodos Regex</b> <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)
<b>Importar y usar modulos y sus funciones</b> <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
<b>Libreria os</b> <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
<b>Libreria shutil</b> from shutil import rmtree <code>rmtree(‘carpeta’)</code> borra la carpeta y subcarpetas
<b>Abrir y cerrar ficheros</b> Primero hay que guardar la ruta del archivo: ubicacion_carpeta = os.getcwd() nombre_archivo = “text.txt” <code>ubicacion_archivo = ubicacion_carpeta + “/” + nombre_archivo</code>  <code>f = open(ubicacion_archivo)</code> abrir un archivo en variable f <code>f.close()</code> cerrar un archivo * IMPORTANTE * <code>with open(ubicacion_archivo) as f:</code> <code>    codigo e.g. variable = f.read()</code> abre el archivo solo para ejecutar el codigo indicado (y despues lo deja)
<b>Encoding</b> <code>from locale import getpreferredencoding</code> <code>getpreferredencoding()</code> 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()
<b>mode: argumento opcional al abrir un archivo</b> <code>r</code> – read <code>w</code> – write - sobreescribe <code>x</code> – exclusive creation, sólo crearlo si no existe todavía <code>a</code> – appending, añadir texto al archivo sin manipular el texto que ya habia hay que anadir otra letra: <code>t</code> – texto – leer en texto <code>b</code> – bytes – leer en bytes (no se puede usar con encoding)
<code>f = open(ubicacion_archivo, mode = “rt”)</code>
<b>Leer ficheros</b> <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
<b>Escribir en ficheros</b> <code>with open(ubicacion_archivo, “w”) as f:</code> <code>    f.write(“Texto que va en el fichero.”)</code> para escribir <code>with open(ubicacion_archivo, “a”) as f:</code> <code>    f.write(“Texto que va en el fichero.”)</code> para anadir texto <code>f.writelines(‘lista’)</code> 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
<b>Conectar a una base de datos</b> <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
<b>Realizar queries</b> <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
<b>Argumentos cursor:</b> <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
<b>Obtener resultados de una query</b> <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
<b>Pandas dataframe with SQL</b> <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
<b>Crear y alterar una base de datos</b> <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>
<b>Insertar datos</b> <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>
<b>Insertar múltiples filas a una tabla</b> <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
<b>Eliminar registros</b> <code>variable_query = “DROP TABLE nombre_tabla”</code>
<b>Añadir errores</b> importar errorcode y usar try/except: <code>try:</code> <code>    accion</code> <code>except mysql.connector.Error as err:</code> <code>    print(err)</code> <code>    print("Error Code:", err.errno)</code> <code>    print("SQLSTATE", err.sqlstate)</code> <code>    print("Message", err.msg)</code>



Python Cheat Sheet 3				
Pandas	<div><div>Crear DataFrames</div><div>df = pd.DataFrame(data, index, columns)</div><div>data: NumPy Array, diccionario, lista de diccionarios</div><div>index: indice que por defecto se asigna como 0-(n-1), n siendo el número de filas;</div><div>index = [lista] para asignar “etiquetas” (nombres de filas)</div><div>column: nombre de las columnas; por defecto 0-(n-1);</div><div>columns =[lista] para poner mas nombres</div></div>	<div><div>df.head(n)</div> devuelve las primeras n lineas del dataframe</div> <div><div>df.tail(n)</div> devuelve las últimas n lineas del dataframe</div> <div><div>df.sample(n)</div> devuelve n filas aleatorias de nuestro dataframe, o uno por defecto</div> <div><div>df.shape</div> devuelve el número de filas y columnas</div> <div><div>df.dtypes</div> devuelve el tipo de datos que hay en cada columna</div> <div><div>df.columns</div> devuelve los nombres de las columnas</div> <div><div>df.describe</div> devuelve un dataframe con un resumen de los principales estadísticos de las columnas numéricas</div> <div><div>df.info()</div> devuelve un resumen sobre el no. de columnas, nombres de columnas, numero de valores no nulos y los tipos de datos</div> <div><div>df[“nombre_columna”.unique()]</div> o <div>df.nombre_columna.unique()</div> devuelve un array con los valores únicos de la columna</div> <div><div>df[“nombre_columna”.value_counts()]</div> o <div>df.nombre_columna.value_counts()</div> devuelve una serie con el recuento de valores únicos en orden descendente</div> <div><div>df.duplicated().sum()</div> devuelve el numero de filas duplicadas</div> <div>Eliminar filas duplicadas</div> <div><div>df.drop_duplicates(inplace = True, ignore_index=True)</div> elimina filas duplicadas; ignore_index para no tener el indice en cuenta</div> <div>Metodos de estadistica</div> <div><div>df[‘columna’].mean()</div>   <div>mode()</div>   <div>median()</div>   <div>var()</div>   <div>std()</div> calcula la media/moda/mediana/variación/desviación estándar de los valores de una columna</div> <div><div>df[‘columna1’].corr(df[‘columna2’])</div> calcula la correlacion entre dos variables</div> <div><div>matriz_correlacion = df.corr()</div> crea una matriz mostrando las correlaciones entre todos los variables</div> <div><div>df_crosstab = pd.crosstab(df[‘columna1’], df[‘columna2’], normalize = True, margins = True)</div></div> <div>normalize muestra los valores en porcentajes (por uno)</div> <div>margins muestra los totales y subtotales</div> <div><div>media_ponderada = np.average(df[‘columna’], weights = w)</div> calcula la media ponderada según los pesos</div> <div><div>percentil_n = np.percentile(df[‘columna’], n)</div> saca el valor en el percentil n</div> <div><div>q3, q1 = np.percentile(df[“columna”], [75, 25])</div> saca los tercer y primer cuartiles</div> <div>Sidetable: frecuencias de datos</div> <div><div>df.stb.freq([‘columna’])</div> devuelve un dataframe con informacion sobre la frecuencia de ocurrencia de cada categoría de un variable categorica</div> <div>parametros:</div> <div>thresh = n limita los valores mostrados a los más frecuentes hasta un umbral de n% cumulative y agrupando los restantes bajo la etiqueta “other”</div> <div>other_label = ‘etiqueta’ cambia la etiqueta ‘other’</div> <div>value = ‘columna’ ordena los resultados por la columna especificada</div> <div><div>df.stb.freq([‘columna1’, ‘columna2’])</div> combina dos columnas y devuelve las frecuencias de las subcategories</div> <div><div>df.stb.missing([‘columna’])</div> devuelve informacion sobre la frecuencia de datos nulos</div>	<div>Tipos de datos en Pandas:</div> <div><div>- object</div><div>- int64</div><div>- float64</div><div>- datetime, timedelta[ns]</div><div>- category</div><div>- bool</div></div> <div><div>df.dtypes</div> devuelve el tipo de datos que hay en cada columna</div> <div><div>df_tipo = df.select_dtypes(include = “tipo”)</div> crea un dataframe de las columnas del tipo de datos especificado</div> <div><div>df[‘columna’] = df[‘columna’].astype(‘tipo’, copy = True, errors = ‘ignore’)</div> convierte una columna en el tipo de dato especificado</div> <div><div>copy = True</div> devuelve una copia</div> <div><div>copy = False</div> *cuidado: los cambios en los valores pueden propagarse a otros objetos pandas*</div> <div><div>errors = ignore</div> omita excepciones; en caso de error devuelve el objeto original</div> <div><div>errors = raise</div> permite que se generen excepciones</div> <div><div>pd.options.display.max_columns = None</div> ejecutar antes del df.head() para poder ver todas las columnas</div> <div><div>pd.set_option("display.precision", 2)</div></div> <div>Outliers</div> <div><div>Calcular tres desviaciones estandares:</div><div><div>media = df.column.mean()</div><div><div>desviacion = df.column.std()</div></div><div><div>lcb = media - desviacion * 3</div><div><div>ucb = media + desviacion * 3</div></div></div><div>Eliminar Outliers</div><div><div>outlier_step = 1.5 * IQR</div> calcular outlier step</div><div><div>outliers_data = df[(df[‘columna’] &lt; Q1 - outlier_step)   (df[‘columna’] &gt; Q3 + outlier_step)]</div> identificar datos fuera del rango del maximo hasta el minimo</div><div><div>lista_outliers_index = list(outliers_data.index)</div> crear una lista de los indices de las filas con outliers</div><div><div>if outliers_data.shape[0] &gt; 0:</div><div><div>dicc_indices[key] = (list(outliers_data.index))</div> crear un diccionario de los indices de las filas con nulos; se puede hacer iterando por columnas</div></div><div><div>valores = dicc_indices.values()</div> sacar todos los valores e.g. todos los indices</div><div><div>valores = {indice for sublista in valores for indice in sublista}</div> set comprehension para eliminar duplicados</div><div><div>df_sin_outliers = df.drop(df.index[list (valores)])</div> crear nuevo dataframe sin outliers</div><div>Reemplazar Outliers</div><div><div>for k, v in dicc_indices.items():</div><div><div>media = df[k].mean()</div><div><div>for i in v:</div><div><div>df.loc[i,k] = media</div> reemplazar outliers por la media</div></div></div></div></div></div>	<div>Identificar nulos</div> <div><div>df.isnull()</div> o <div>df.isna()</div> devuelve True o False según si cada valor es nulo o no</div> <div><div>df.isnull().sum()</div> o <div>df.isna().sum()</div> devuelve una serie con el número de valores nulos por columnas</div> <div><div>df_% nulos = ((df.isnull().sum() / df.shape[0] * 100).reset_index())</div><div><div>df_% nulos.columns = [‘columna’, ‘% nulos’]</div> crea un dataframe de los porcentajes de los valores nulos</div></div> <div>Eliminar nulos</div> <div><div>df.dropna(inplace = True, axis=b, subset=[lista_de_columnas], how=)</div> quitar nulos</div> <div><div>how = ‘any’   ‘all’</div> 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</div> <div><div>subset</div> una columna o lista de columnas</div> <div>Tipos de nulos</div> <div><div>np.nan</div> significa “not a number”; es un tipo numérico</div> <div><div>None</div> valores nulos en columnas tipo string</div> <div><div>NaT</div> valores nulos tipo datetime</div> <div><div>valores texto: “n/a”, “NaN”, “nan”, “null”</div> strings que normalmente se convierten automaticamente a np.nan</div> <div><div>99999</div> o <div>00000</div> integers que se pueden convertir a nulos</div> <div>Reemplazar nulos</div> <div><div>df = pd.read_csv(‘archivo.csv’, na_values = [‘n/a’])</div><div><div>.fillna(np.nan)</div> reemplaza los strings ‘n/a’ con np.nan al cargar el dataframe</div><div><div>df.fillna(df[value=n, axis=b, inplace=True])</div> reemplazar todos los NaN del dataframe con el valor que especifiquemos</div><div><div>df[‘columna’].fillna(df[‘columna’].median, axis=b, inplace=True)</div> reemplazar los nulos de una columna por la mediana de esa columna</div><div><div>value=n</div> por defecto NaN; es el valor por el que queremos reemplazar los valores nulos que puede ser un escalars, diccionario, serie o dataframe</div><div><div>axis</div> por defecto 0 (filas)</div><div><div>df.replace(valor_nulo, valor_nuevo, inplace=True, regex=False)</div> reemplazar los nulos por el valor nuevo</div><div>Imputacion de nulos</div><div><div>from sklearn.impute import SimpleImputer</div><div><div>imputer = SimpleImputer(strategy=‘mean’, missing_values = np.nan)</div> inicia la instancia del metodo, especificando que queremos reemplazar los nulos por la media</div><div><div>imputer = imputer.fit(df[‘columna1’])</div> aplicamos el imputer</div><div><div>df[‘media_columna1’] = imputer.transform(df[[‘price’]])</div> rellena los valores nulos según como hemos especificado</div><div><div>from sklearn.experimental import enable_iterative_imputer</div><div><div>from sklearn.impute import IterativeImputer</div><div><div>imputer = IterativeImputer(n_nearest_features=n, imputation_order=‘ascending’)</div> crea la instancia</div><div><div>n_nearest_features</div> por defecto None; numero de columnas a utilizar para estimar los valores nulos</div><div><div>imputation_order</div> por defecto ascendente; el orden de imputacion</div><div><div>imputer.fit(df_numericas)</div> aplicamos el imputer</div><div><div>df_datos_trans = pd.DataFrame(imputer.transform(df_numericas), columns = df_numericas.columns)</div> crea un dataframe de los datos transformados; metemos estas columnas en el dataframe original</div><div><div>from sklearn.impute import KNNImputer</div><div><div>imputerKNN = KNNImputer(n_neighbors=5)</div> crea la instancia</div><div><div>imputerKNN.fit(df_numericas)</div><div><div>df_knn_imp = pd.DataFrame(imputerKNN.transform(df_numericas), columns = numericas.columns)</div> crea un dataframe de los datos transformados; metemos estas columnas en el dataframe original</div></div></div></div></div></div></div>

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

Subsets: loc e iloc
<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> <div><b>df.loc[“etiqueta_fila”, “etiqueta_columna”]</b> devuelve el contenido de un campo en una columna de una fila  <b>df.loc[“etiqueta_fila”,:]</b> devuelve los valores de todas las columnas de una fila  <b>df.loc[:,“etiqueta_columna”]</b> devuelve los valores de todas las filas de una columna  <b>df.iloc[indice_fila, indice_columna]</b> devuelve el contenido de un campo en una columna de una fila  <b>df.iloc[indice_fila, :]</b> devuelve los valores de todas las columnas de una fila  <b>df.iloc[:,indice_columna]</b> devuelve el contenido de un campo en una columna de una fila  <b>df.loc[[lista_etiquetas_filas], [lista_etiquetas_columnas]]</b> devuelve el contenido de varias filas / varias columnas  <b>df.loc[[lista_indices_filas], [lista_indices_columnas]]</b> 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  <b>df.loc[df.etiqueta &gt; x]</b> seleccionar datos basado en una condición usando operadores comparativos  <b>df.loc[(df.etiqueta &gt; x) &amp; (df.etiqueta == y)]</b> seleccionar datos que tienen que cumplir las dos condiciones (and)  <b>df.loc[(df.etiqueta &gt; x)   (df.etiqueta == y)]</b> seleccionar datos que tienen que deben cumplir una de las dos condiciones (or)  <b>df.iloc[list(df.etiqueta &gt; x), :]</b> iloc no acepta una Serie booleana; hay que convertirla en lista  <b>variable_df.head(n)</b> 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><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><b>df_filtrado = df[df[“nombre_columna”] == valor]</b> 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><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><b>df_filtrado = df[(df[“columna1”] == valor) &amp; (df[“columna2”] == valor) &amp; (df[“columna3”] &gt; n valor)]</b> extrae las filas donde los valores de las columnas cumplan las condiciones en parentesis</div>
<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> <div><b>df_filtrado = df[(df[“columna1”] == valor)   (df[“columna1”] == valor)</b> extrae las filas donde los valores de las columnas cumplan con una condición u otra</div>
<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> <div><b>df_filtrado = ~(df[df[“columna1”] == valor])</b> extrae las filas donde los valores de las columnas NO cumplan con la condición</div>

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

## Cambiar columnas

<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> <div><b>lista_columnas = df.columns.to_list()</b> crea una lista de los nombres de las columnas del dataframe  <b>df.set_index([“nombre_columna”, inplace = True)</b> establece el índice utilizando uno o mas columnas; puede sustituir o ampliar un índice existente <b>inplace = True</b> los cambios sobrescriben sobre el df  * cuando una columna se cambia a índice ya no es columna *  <b>df.reset_index(inplace = True)</b> 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><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><b>df.rename(columns = {“nombre_columna”: “nombre_nueva”}, inplace = True)</b> cambia los nombres de una o mas columnas  ejemplo de dict comprehension para crear diccionario sobre las columnas existentes de un dataframe: <b>diccionario = {col : col.upper() for col in df.columns}</b>  <b>df.rename(columns = diccionario, inplace = True)</b> cambia los nombres de las columnas según el diccionario</div>
---

### Eliminar columnas

<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> <div><b>df.drop(columns = [“columna1”, “columna2”], axis = b, inplace=True)</b> eliminar una o mas columnas o filas segun lo que especificamos</div>
--

### Reordenar columnas

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

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

## Apply

<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> <div><b>apply()</b> toma una función como argumento y la aplica a lo largo de un eje del DataFrame <b>df[‘columna_nueva’] = df[‘col_1’].apply(función)</b> crea una columna nueva con los valores de otra columna transformados según la función indicada <b>df[‘columna_nueva’] = df[‘col_1’].apply(lambda x: x.método() if x &gt; 1)</b> crea una columna nueva con los valores de otra columna transformados según la lambda indicada <b>df[‘columna_nueva’] = df.apply(lambda nombre: función(nombre[‘columna1’], nombre[‘columna2’]), axis = b)</b> crea una columna nueva usando una función que coge dos parámetros (columna 1 y columna2) <b>df.applymap(funcion, na_action=None, **kwargs)</b> acepta y devuelve un escalar a cada elemento de un dataframe; se tiene que aplicar a todo el DataFrame <b>df[‘columna’] = df[‘columna’].map(mapa, na_action = ‘ignore’)</b> reemplaza valores de la columna según el mapa, que puede ser un diccionario o una serie; solo se puede aplicar a una columna en particular.  <b>apply()</b> con datetime <b>df[‘columna_fecha’] = df[‘columna_fecha’].apply(pd.to_datetime)</b> cambia una columna de datos tipo fecha en el formato datetime <b>def sacar_año(x):</b>     <b>return x.strftime(“%Y”)</b>  <b>df[‘columna_año’] = (df[‘columna_fecha’].apply(sacar_año)</b> crea una columna nueva del año solo usando un método de la librería datetime; (“%B”) para meses</div>
---

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

## datetime

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

## APIs

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

### Codigos de respuesta de HTTP

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







EDA Exploratory Data Analysis	ETL: Extract, Transform, Load			
<b>Análisis exploratorio de datos</b>	<b>Extraccion</b> <ul style="list-style-type: none"><li>- obtener datos crudos y almacenarlos<ul style="list-style-type: none"><li>- Tablas de bases de datos SQL o NoSQL</li><li>- Ficheros de texto plano</li><li>- Emails</li><li>- Información de páginas web</li><li>- Hojas de cálculo</li><li>- Ficheros obtenidos de API's</li><li>- Otros tipos de datos.</li></ul></li></ul> <b>Transformación</b> <ul style="list-style-type: none"><li>- procesar los datos, unificarlos, limpiarlos, validarlos, filtrarlos, etc.<ul style="list-style-type: none"><li>- Formetear fechas</li><li>- Reordenar filas o columnas</li><li>- Unir o separar datos</li><li>- Combinar las fuentes de datos</li><li>- Limpiar y estandarizar los datos</li><li>- Verificar y validar los datos</li><li>- Eliminar duplicados o datos erroneos</li><li>- Filtrado, realización de calculos o agrupaciones</li></ul></li></ul> <b>Carga</b> <ul style="list-style-type: none"><li>- cargar los datos en su formato de destino, el tipo de lo cual dependerá de la naturaleza, el tamaño y la complejidad de los datos. Los sistemas más comunes suelen ser:<ul style="list-style-type: none"><li>- Ficheros csv</li><li>- Ficheros json</li><li>- Bases de datos</li><li>- Almacenes de datos (Data Warehouse)</li><li>- Lagos de datos (Data Lakes)</li></ul></li></ul> <b>APIs</b> <p>es un tipo de conexión entre ordenadores o diferentes programas del ordenadores, haciendo uso de funciones o subrutinas que ofrece una bibloteca para ser utilizado por otro software</p>			
<p>El Análisis Exploratorio de Datos se refiere al proceso de realizar una serie de investigaciones iniciales sobre los datos que tenemos para poder descubrir patrones, detectar anomalías, probar hipótesis y comprobar suposiciones con la ayuda de estadísticas y representaciones gráficas.</p>				
<b>1. Entender las variables</b>				
<ul style="list-style-type: none"><li>- que variables temenos</li><li><code>.head()</code>, <code>.tail()</code>, <code>.describe()</code>, <code>.info()</code>, <code>.shape</code></li><li>- que tipos de datos</li><li><code>.dtypes()</code>, <code>.info()</code></li><li>- si temenos nulos o duplicados</li><li><code>.isnull().sum()</code></li><li><code>.duplicated().sum()</code></li><li>- que valores unicos temenos</li><li><code>.unique()</code>, <code>.value_counts()</code></li></ul> <p>librería sidetable: <code>stb.freq()</code> devuelve el <code>value_counts</code> de variables categóricas, mas el porcentaje, cuenta cumulativa y porcentaje cumulativa <code>stb.missing()</code> tabla de cuenta de nulos y el porcentaje del total</p>				
<b>2. Limpiar el dataset</b>				
<ul style="list-style-type: none"><li>- quitar duplicados (filas o columnas)</li><li>- cambiar nombres de columnas</li><li>- cambiar tipo de datos de columnas</li><li>- ordenar columnas</li><li>- separar columna en dos con <code>str.split()</code></li><li>- crear intervalos con <code>pd.cut()</code></li><li>- crear porcentajes o ratios</li><li>- decidir como tratar outliers: mantenerlos, eliminarlos, o reemplazarlos con la media, mediana o moda; o aplicar una imputacion</li></ul> <ul style="list-style-type: none"><li>- decidir como tratar nulos:<ul style="list-style-type: none"><li>- eliminar filas o columnas con nulos <code>drop.na()</code></li><li>- imputar valores perdidos:<ul style="list-style-type: none"><li>- reemplazarlos con la media, mediana o moda usando <code>.fillna()</code> o <code>.replace()</code></li><li>- imputer con metodos de machine learning usando la libreria sklearn: Simple-Imputer, Iterative-Imputer, o KNN Imputer</li></ul></li></ul></li></ul>				
<b>3. Analizar relaciones entre variables</b>				
<p>Analizar relaciones entre las variables</p> <ul style="list-style-type: none"><li>- para encontrar patrones, relaciones o anomalias</li></ul> <p>Relaciones entre dos variables numéricas:</p> <ul style="list-style-type: none"><li>- <code>scatterplot</code></li><li>- <code>regplot</code> – <code>scatterplot</code> con línea de regresion</li><li>- matriz de correlación y <code>heatmap</code></li><li>- <code>joinplot</code> – permite emparejar dos gráficas – una histograma con <code>scatter</code> o <code>reg plot</code> por ejemplo</li></ul> <p>Relaciones entre dos variables categóricas:</p> <ul style="list-style-type: none"><li>- <code>countplot</code></li></ul> <p>Relaciones entre variables numéricas y categóricas:</p> <ul style="list-style-type: none"><li>- <code>swarmplot</code></li><li>- <code>violinplot</code></li><li>- <code>pointplot</code></li><li>- <code>boxplot</code></li></ul>				