

```
Examen 3 por de datos en pandas:

- Series --> Son las series de elementos en una dimensión, como pueden ser las tuplas, las listas o los diccionarios.
- DataFrames --> Son los sets de datos de 2 dimensiones, es decir, tablas
- Index objects --> Se usan para guardar metadato de otras estructuras de datos.

In [16]: import pandas as pn
import numpy as np

SERIES

Se definen como pn.Series(), por ejemplo:

In [3]: Series(pn.Series([1,2,3]))
Series
Out[3]: 0    1
        1    2
        dtype: int64

Podemos crear un diccionario en el que las keys serán los índices y los valores los valores de la serie:

In [4]: dic={'a':100, 'b':200, 'c':300, 'd':400, 'e':800}
serieA=pn.Series(dic)
serieA
Out[4]: a    100
        b    200
        c    300
        d    400
        e    800
        dtype: int64

serieA.values, serieA.index

In [5]: serieA.values #Permite ver el contenido de la serie
Out[5]: array([1, 2, 3], dtype=int64)

In [6]: serieA.index #Permite ver el índice de los elementos de la serie
Out[6]: RangeIndex(start=0, stop=3, step=1)

Cómo llamar a elementos concretos de una serie

In [4]: serie2=pn.Series([1,2,3,4,5], index=['a','b','c','d','e']) #Permite indexar los elementos del array según éste es creado

In [5]: serie2.index #Permite ver los índices de la serie indicada
Out[5]: Index(['a', 'b', 'c', 'd', 'e'], dtype='object')

In [6]: serie2['a', 'b', 'c', 'e', 'd', 'e']
Out[6]: 1
        2
        3
        4
        4
        dtype: int64

In [7]: serie2[['b','c']] #Para acceder a varios datos debes poner doble corchete [[]]
Out[7]: a    1
        b    2
        c    3
        dtype: int64

In [12]: serie2['b'] #Puedes cambiar los valores como en Python
Out[12]: 4

In [13]: serie2[serie2==0] #Puedes introducir Booleanos en los corchetes
Out[13]: a    1
        c    3
        dtype: int64

In [14]: serie3=pn.Series({'Pepe':19, 'Juan':21, 'Luis':30, 'Paco':10}) #también funciona para diccionarios
Out[14]: Pepe    19
        Juan    21
        Luis    30
        Paco    10
        dtype: int64

In [15]: serie3=pn.Series({'Pepe':19, 'Juan':21, 'Luis':30, 'Paco':None}, index=['Pepe', 'Juan', 'Luis', 'Paco']) #también podemos poner ,index= en un diccionario, pero como el object ya es un index de por si es un poco tontería
Out[15]: Pepe    19.0
        Juan    21.0
        Luis    30.0
        Paco    NaN
        dtype: float64

In [16]: pn.isnull() y pn.notnull()
In [8]: pn.isnull(serie2) #te enseña que elementos de la serie son nulos, si es nulo pondrá True y si no lo es pondrá False
Out[8]: a    False
        b    False
        c    False
        d    False
        e    False
        dtype: bool

In [17]: pn.notnull(serie2) #Lo opuesto a pn.isnull(), te muestra los elementos de la serie que no son nulos con True, y los que son nulos con False
Out[17]: a    True
        b    True
        c    True
        d    True
        e    False
        dtype: bool

Un concepto importante de pandas es que a la hora de hacer operaciones aritméticas sobre series hace un alineamiento automático, es decir, busca los índices que son iguales y los fusiona, y si hay algún índice que no coincide se le asigna valor NaN, independientemente de si antes tenía un valor o no.
```

```
DATAFRAMES

Crear un dataframe

Se crean igual que una serie, pero es necesario poner pn.DataFrame() y la tabla se crea como un diccionario en el que la key es el título de la columna (categoría) y el value es una tupla en la que cada valor es una fila

In [16]: data={'state':['Ohio','Ohio','Ohio','Nevada','Nevada'], 'year':[2000,2001,2002,2001,2002], 'pop':[1.5,1.7,3.6,2.4,2.9]}
frame=pn.DataFrame(data)
frame
Out[16]:   state  year  pop
0  Ohio  2000  1.5
1  Ohio  2001  1.7
2  Ohio  2002  3.6
3  Nevada 2001  2.4
4  Nevada 2002  2.9
```

```
Añadiendo ,columns='X','Y','Z'... podemos cambiar el orden o título de las categorías o añadir nuevas categorías, pero éstas al no haberle merteido ningún valor tendrán valor NaN.

In [19]: frame2=pn.DataFrame(data, columns=['pop','state','year'])
frame2
Out[19]:   pop  state  year
0  1.5  Ohio  2000
1  1.7  Ohio  2001
2  3.6  Ohio  2002
3  2.4  Nevada 2001
4  2.9  Nevada 2002
```

```
Añadiendo ,index='one','two','three'... podemos cambiar los índices (números que aparecen a la izquierda de la tabla), que por defecto son 0,1,2,3,4...

In [34]: frame3=pn.DataFrame(data, columns=['pop','state','year'], index=['one','two','three','four','five'])
frame3
Out[34]:   pop  state  year
one  1.5  Ohio  2000
two  1.7  Ohio  2001
three 3.6  Ohio  2002
four  2.4  Nevada 2001
five  2.9  Nevada 2002
```

```
frame.values, frame.columns y frame.index

In [21]: frame3.values #te muestra los valores del DataFrame creado
Out[21]: array([[1.5, 'Ohio', 2000],
               [1.7, 'Ohio', 2001],
               [3.6, 'Ohio', 2002],
               [2.4, 'Nevada', 2001],
               [2.9, 'Nevada', 2002]])
dtype=object

In [22]: frame3.columns #te muestra los nombres de las columnas, es decir el título de cada columna, es decir, las categorías
Out[22]: Index(['pop', 'state', 'year'], dtype='object')

In [23]: frame3.index #te muestra los índices (por defecto son 0,1,2,3,4...)
Out[23]: Index(['one', 'two', 'three', 'four', 'five'], dtype='object')
```

```
Para llamar a elementos de un dataframe se hace de la misma forma que para llamar a elementos de una serie, pero teniendo en cuenta que primero llamas a un diccionario y luego a una tupla indexada, por ello:

In [24]: frame3['pop'] #Llamamos a la columna pop al completo
Out[24]: one    1.5
        two    1.7
        three  3.6
        four    2.4
        five    2.9
        Name: pop, dtype: float64

In [35]: frame3['pop'][ 'three' ] #Llamamos especificamente al elemento de la columna pop con índice three
Out[35]: 3.6
```

```
Coger una fila completa en función del índice con frame.loc[indice]

In [26]: frame3[frame3.loc['three']] #te permite meter en una variable la fila de índice 'three' con sus correspondientes categorías, que pasan a ser los índices de esta serie creada
frame3
Out[26]: pop      2.4
        state  Ohio
        year   2001
        Name: three, dtype: object

Añadir categorías (columnas) a un dataframe ya creado

Para ello debemos de llamar a un dataframe con un índice que no exista, de forma que se creará esta nueva categoría y se llenará con los valores que especifiquemos:
```

```
In [37]: import numpy as np
frame3['categoría_añadida']=np.nan
frame3
Out[37]:   pop  state  year  categoría_añadida
one  1.5  Ohio  2000                NaN
two  1.7  Ohio  2001                NaN
three 3.6  Ohio  2002                NaN
four  2.4  Nevada 2001                NaN
five  2.9  Nevada 2002                NaN
```

```
IMPORTANTE: CAMBIAR EL VALOR DE UN ELEMENTO CONCRETO DE LA TABLA

TABLA.LOC[INDICE,'CATEGORIA'] Es muy importante que pongas coma y no separes los elementos en 2 corchetes, porque si lo haces no te va a dejar cambiarlo.

Por ejemplo si queremos cambiar el año de la fila con índice five:

In [28]: frame3.loc['five','year']=2001
frame3
Out[28]:   pop  state  year  categoría_añadida
one  1.5  Ohio  2000                None
two  1.7  Ohio  2001                None
three 3.0  Ohio  2002                None
four  2.4  Nevada 2001                None
five  2.9  Nevada 2001                None
```

```
Forma alternativa de crear dataframes

Podemos establecer un data frame con un diccionario de diccionarios, en el que el primer key es el título de la columna, y el value es un diccionario en el que a su vez distinguimos keys (que serán el índice de cada fila) y los valores serán los valores de la columna. Por ello podemos crear:

In [29]: frame3=pn.DataFrame({'pop':{'one':1.2,'two':2.3,'three':4.5,'four':5.6,'five':6.7}, 'state':{'one':'Ohio','two':'Ohio','three':'Ohio','four':'Nevada','five':'Nevada'}, 'year':{'one':2000,'two':2001,'three':2002,'four':2001,'five':2004}}
frame3
Out[29]:   pop  state  year
one  1.2  Ohio  2000
two  2.3  Ohio  2001
three 4.5  Ohio  2002
four  5.6  Nevada 2001
five  6.7  Nevada 2002
```

```
EJERCICIO

• You are working on a project for data analysis related with the automotive sector. In order to do some tests you plan to define some test data
• Think about the structure of DataFrame for storing the relevant data of the car models of a company
• Once defined the structure, fill in the DataFrame with some test data using any of the methods explained for DataFrame creation
• Change the DataFrame index and assign a new one, for example changing the data type (e.g. int to string)
```

```
In [28]: frame_coches=pn.DataFrame({'matricula':['2345AVV','3333AAA','4444VVV','4444BBB'], 'marca':['Renault','Skoda','Peugeot','Seat'], 'gama':['1,2,2,2'], 'color':['Azul','Verde','Verde','Gris'], 'precio':[40,10,20,15], 'mano':['1,1,1,1]})
frame_coches
Out[28]:   matricula  Marca  Gama  Color  Precio  Mano
0  2345AVV  Renault    1  Azul    40    1
1  3333AAA  Skoda     2  Verde   10    1
2  4444VVV  Peugeot   2  Verde   20    1
3  4444BBB  Seat      2  Gris    15    1
```

```
In [31]: frame_coches.pn.DataFrame(columns=['Matricula','Marca','Gama','Color','Precio','Mano'])
frame_coches
Out[31]:   Matricula  Marca  Gama  Color  Precio  Mano
```

```
In [32]: frame_coches['Matricula']=['2345AVV','3333AAA','4444VVV','4444BBB']
frame_coches['Marca']=['Renault','Skoda','Peugeot','Seat']
frame_coches['Gama']=['1,2,2,2']
frame_coches['Color']=['Azul','Verde','Verde','Gris']
frame_coches['Precio']=[40,10,20,15]
frame_coches['Mano']=[1,1,1,1]
frame_coches
```

```
Out[32]:   Matricula  Marca  Gama  Color  Precio  Mano
0  2345AVV  Renault    1  Azul    40    1
1  3333AAA  Skoda     2  Verde   10    1
2  4444VVV  Peugeot   2  Verde   20    1
3  4444BBB  Seat      2  Gris    15    1
```

```
In [33]: frame_coches.pn.Series(index=['a','b','c','d'])
frame_coches
Out[33]:   Matricula  Marca  Gama  Color  Precio  Mano
a  2345AVV  Renault    1  Azul    40    1
b  3333AAA  Skoda     2  Verde   10    1
c  4444VVV  Peugeot   2  Verde   20    1
d  4444BBB  Seat      2  Gris    15    1
```

```
In [34]: serie8=pn.Series([2, 4, 0, 0, 10])
serie8=pn.Series([1, 3, 5, 7, 10])
suma=serieA+serie8
resta=serieA-serie8
multiplicacion=serieA*serie8
potenciacion=serieA**serie8
print(suma,resta,multiplicacion,potencia)
```

```
0    3
1    7
2   11
3   15
4   20
dtype: int64
0    1
1    1
2    1
3    1
dtype: int64
0    2
1   16
2   30
3   50
dtype: int64
0    64
1   776
2  2897152
3  1000000000
dtype: int64
```

```
Ejercicio 2: Crear una serie a partir de este diccionario

In [35]: dic={'a':100, 'b':200, 'c':300, 'd':400, 'e':800}
values=[]
keys=[]
for i in dic.values():
    values.append(i)
for i in dic.keys():
    keys.append(i)
serieA=pn.Series(values,index=keys)
serieA=pn.Series(dic)
serieA
```

```
Out[35]: a    100
        b    200
        c    300
        d    400
        e    800
        dtype: int64
```

```
Ejercicio 3: Write a Pandas program to convert the first column of a DataFrame as a Series and then convert that Series to an Array

In [36]: data_frame=pn.DataFrame(columns=['col1','col2','col3'])
data_frame['col1']=[1,2,3,4,7,11]
data_frame['col2']=[1,3,5,9,5,9,9]
data_frame['col3']=[1,5,6,12,1,11]
columns=data_frame['col1']
columns
Out[36]: 1
        2
        3
        4
        7
        11
        dtype: int64
array([1, 2, 3, 4, 7, 11])
array([1, 2, 3, 4, 7, 11])
for i in array_python:
    array_final_python.append(i)
array_final_python
```

```
Out[36]: [1, 2, 3, 4, 7, 11]
```

```
Ejercicio 4: Write a Pandas program to get the items of a given series not present in another given series.

In [22]: serie1=pn.Series([1,2,3,4,7,11])
serie2=pn.Series([4,5,6,7,9,6,11,0])
lista=[]
long=len(serie2)
for i in serie1:
    cont=0
    for j in serie2:
        if i==j:
            cont=cont+1
    if cont==long:
        lista.append(i)
serie3=pn.Series(lista)
serie3
```

```
Out[22]: 0    1
        1    2
        2    3
        dtype: int64
```

```
Ejercicio 5: Write a Pandas program to compute the minimum, 25th percentile, median, 75th, and maximum of a given series.

In [38]: import numpy as np
serie2=pn.Series([4,5,6,7,9,5,11,0])
minu=pn.Series(2)
print(minu)

print(max(serie2))

print(np.percentile(serie2,25))
print(np.percentile(serie2,50))
print(np.percentile(serie2,75))
```

```
0
11
4.75
5.5
7.5
```

```
Ejercicio 6: Write a Pandas program to display most frequent value in a given series and replace everything else as 'Other' in the series.

In [39]: serie3=pn.Series([4,5,6,7,9,5,11,0,12,34,3,4,5,6,2,1,3,4,3,4,3,5,6,6,6,6,7,45,6,7,4,1,1])
max=0
for i in serie:
    cont=0
    for j in serie:
        if i==j:
            cont=cont+1
    if cont > max:
        max=cont
        most_frequent=i

print(serie)

0    Other
1    Other
2    6
3    Other
4    Other
5    Other
6    Other
7    Other
8    Other
9    Other
10  Other
11  Other
12  6
13  Other
14  Other
15  Other
16  Other
17  Other
18  Other
19  Other
20  Other
21  Other
22  6
23  Other
24  6
25  Other
26  Other
27  Other
28  6
29  Other
30  Other
31  Other
32  Other
dtype: object
```

```
Ejercicio 7: Write a Pandas program to calculate the frequency counts of each unique value of a given series.

In [12]: serie7=pn.Series([4,5,6,7,9,5,11,0,12,34,3,4,5,6,2,1,3,4,3,4,3,5,6,6,6,6,7,45,6,7,4,1,1])
lista_unicos=[]
long=len(serie)

for i in serie:
    cont=0
    for j in serie:
        if i==j:
            if cont==long-1: #es por que al revisar la serie entera se compara consigo mismo
                lista_unicos.append(int(i)/long)

print(lista_unicos)

[0.2727272727272727, 0.3333333333333333, 0.0, 0.36363636363636365, 1.0303030303030303, 0.06060606060606061, 1.3636363636363635]
```

```
Ejercicio 8: Write a Pandas program to get the first 3 rows of a given DataFrame.

In [17]: exam_data = {'name': ['Anastasia', 'Dima', 'Katherine', 'James', 'Emily', 'Michael', 'Matthew', 'Laura', 'Kevin', 'Jonas'],
                    'score': [12.5, 9, 16.5, np.nan, 9, 20, 14.5, np.nan, 8, 19],
                    'attempts': [1, 3, 2, 3, 2, 3, 5, 1, 2, 3],
                    'quality': ['yes', 'no', 'yes', 'no', 'no', 'yes', 'yes', 'no', 'no', 'yes']}
labels = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j']

tabla=pn.DataFrame(exam_data,index=labels)

tabla.loc[['a','b','c']]

Out[17]:   name  score  attempts  quality
a  Anastasia   12.5         1      yes
b  Dima        9.0         3       no
c  Katherine   16.5         2      yes
```

```
Ejercicio 9: Write a Pandas program to count the number of rows and columns of a DataFrame.

In [2]: exam_data = {'name': ['Anastasia', 'Dima', 'Katherine', 'James', 'Emily', 'Michael', 'Matthew', 'Laura', 'Kevin', 'Jonas'],
                    'score': [12.5, 9, 16.5, np.nan, 9, 20, 14.5, np.nan, 8, 19],
                    'attempts': [1, 3, 2, 3, 2, 3, 5, 1, 2, 3],
                    'quality': ['yes', 'no', 'yes', 'no', 'no', 'yes', 'yes', 'no', 'no', 'yes']}
labels = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j']

tabla=pn.DataFrame(exam_data,index=labels)

print('El número de líneas es:')
print(len(tabla.name))

print()

print('El número de columnas es:')
print(len(tabla.loc['a']))

El número de líneas es:
10
-----
4
10
```

```
Ejercicio 10: Write a Pandas program to select the rows where the score is missing, i.e. is NaN

In [2]: exam_data = {'name': ['Anastasia', 'Dima', 'Katherine', 'James', 'Emily', 'Michael', 'Matthew', 'Laura', 'Kevin', 'Jonas'],
                    'score': [12.5, 9, 16.5, np.nan, 9, 20, 14.5, np.nan, 8, 19],
                    'attempts': [1, 3, 2, 3, 2, 3, 5, 1, 2, 3],
                    'quality': ['yes', 'no', 'yes', 'no', 'no', 'yes', 'yes', 'no', 'no', 'yes']}
labels = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j']

tabla=pn.DataFrame(exam_data,index=labels)

NaN=np.nan
for i in labels:
    fila=tabla.loc[i]
    for j in fila.isnull():
        if j==True:
            print(fila)

tabla[tabla.isnull().any(axis=1)]

name      James
score      NaN
attempts    3
quality     no
name      Laura
score      NaN
attempts    1
quality     yes
name      h, dtype: object

Out[2]:   name  score  attempts  quality
d  James  NaN         3       no
h  Laura  NaN         1       no
```

```
Ejercicio 11: Write a Pandas program to select the rows the score is between 15 and 20

In [4]: exam_data = {'name': ['Anastasia', 'Dima', 'Katherine', 'James', 'Emily', 'Michael', 'Matthew', 'Laura', 'Kevin', 'Jonas'],
                    'score': [12.5, 9, 16.5, np.nan, 9, 20, 14.5, np.nan, 8, 19],
                    'attempts': [1, 3, 2, 3, 2, 3, 5, 1, 2, 3],
                    'quality': ['yes', 'no', 'yes', 'no', 'no', 'yes', 'yes', 'no', 'no', 'yes']}
labels = ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j']

tabla=pn.DataFrame(exam_data,index=labels)

lista_indices=tabla.index
for i in lista_indices:
    if len(tabla.loc[i]) <=20:
        print(lista_indices)
        print('fin')

name      Katherine
score      20
attempts    2
quality     yes
name      c, dtype: object

name      Michael
score      20
attempts    9
quality     yes
name      f, dtype: object

name      Jonas
score      19
attempts    3
quality     yes
name      j, dtype: object
```

```
Ejercicio 12: Write a Pandas program to change the score in row 'f' to 11.5

In [ ]:   tabla.loc['f','score']=11.5
```

```
Ejercicio 13: Write a Pandas program to group by the first column and get second column as lists in rows.

In [ ]:   col1  col2
0  C1  1
1  C1  2
2  C2  3
3  C2  3
4  C2  4
5  C3  6
6  C2  5

In [75]: dic={'col1':['C1','C1','C2','C2','C3','C2'],'col2':[1,2,3,3,4,6,9]}
tabla=pn.DataFrame(dic,index=[0,1,2,3,4,5,6])
tabla

Out[75]:   col1  col2
0    C1      1
1    C1      2
2    C2      3
3    C2      3
4    C2      4
5    C3      6
6    C2      5

In [ ]: 
```