Módulo I: Introducción a Python

Temario de la Clase 6 Librería Pandas

- Series
- DataFrames
 - Indexado
 - Selección condicional
 - Máscaras
 - Operaciones vectorizadas
 - Modificación
 - Borrado
- CSV

Pandas



Pandas es una librería de Python especializada en el manejo y análisis de estructuras de datos

import pandas as pd

- Su nombre viene de Panel de Datos.
- Me permite trabajar con datos tabulares, es decir, en forma de tabla.
- A diferencia de los arrays de NumPy, Pandas me permite combinar distintos tipos de variables (numéricas, strings, booleanos) en una misma estructura.
- Permite tratar con bases de datos incompletas.
- Estructuras intuitivas y flexibles, al tratarse de datos etiquetados.
- Está basada en estructuras de NumPy, por lo que hereda sus propiedades de optimización y rapidez.

Una serie de Pandas es un objeto comparable con un array uni-dimensional de NumPy.

Todos los datos son del mismo tipo (como un array).

Se pueden asignar índices para trabajar con sus elementos.

Se pueden crear a partir de una lista o un diccionario.

Se puede definir de la siguiente manera:

```
series = pd.Series(data=None, index=None, dtype=None,
name=None)
```

Con **data** podemos pasar la información que queremos almacenar en la serie, **index** nos permite especificar los índices a utilizar, **dtype** se usa de forma similar a Numpy, y **name** me permite darle un nombre a la serie (independiente del nombre de la variable en la que se guarda).

Para ejemplificar la creación de una serie, voy a trabajar con datos de la población del "Grupo de los 7" (G7). Solamente voy a especificar los valores.

```
# In millions
        g7 pop = pd.Series([35.467, 63.951, 80.940, 60.665, 127.061, 64.511, 318.523])
     ✓ 0.0s
[3]
        g7 pop
        0.0s
[4]
          35.467
          63.951
    1
          80.940
        60.665
        127.061
         64.511
    5
         318.523
    dtype: float64
```

Le puedo dar un nombre a la serie para identificar mejor su contenido. Al igual que en listas puedo acceder a los elementos mediante su posición.

```
g7_pop[0]
        g7 pop.name = 'G7 Population in millions'

√ 0.0s

     ✓ 0.0s
[5]
                                                                       35.467
        g7_pop
                                                                          g7_pop[1]
        0.0s
                                                                 [12]

√ 0.0s

           35,467
          63.951
    1
                                                                      63.951
          80.940
         60.665
         127.061
                                                                          g7 pop.index
          64.511
          318.523

√ 0.0s

                                                                 [13]
    Name: G7 Population in millions, dtype: float64
                                                                      RangeIndex(start=0, stop=7, step=1)
```

Vemos de qué tipo es el objeto creado y que su estructura interna está basada en arrays de NumPy.

```
type(g7_pop)
[67]
      ✓ 0.0s
     pandas.core.series.Series
        g7_pop.dtype
     ✓ 0.0s
[7]
    dtype('float64')
        g7_pop.values
      ✓ 0.0s
[8]
    array([ 35.467, 63.951, 80.94 , 60.665, 127.061, 64.511, 318.523])
        type(g7_pop.values)
      ✓ 0.0s
[9]
    numpy.ndarray
```

A diferencia de las listas, puedo especificar el índice de una serie.

```
g7_pop
     ✓ 0.0s
[15]
    Canada
                     35.467
    France
                     63.951
    Germany
                     80.940
    Italy
                 60.665
                  127.061
    Japan
    United Kingdom 64.511
    United States 318.523
    Name: G7 Population in millions, dtype: float64
```

Puedo crear directamente desde un diccionario o especificar al crear el objeto.

```
pd.Series(
            [35.467, 63.951, 80.94, 60.665, 127.061, 64.511, 318.523],
            index=['Canada', 'France', 'Germany', 'Italy', 'Japan', 'United Kingdom',
               'United States'],
            name='G7 Population in millions')

√ 0.0s

[17]
                                                                       pd.Series({
    Canada
                       35.467
                                                                           'Canada': 35.467,
    France
                       63.951
                                                                           'France': 63.951,
    Germany
                       80.940
    Italy
                                                                           'Germany': 80.94,
                      60.665
                                                                           'Italy': 60.665,
    Japan
                      127.061
                                                                           'Japan': 127.061,
    United Kingdom
                      64.511
                                                                           'United Kingdom': 64.511,
    United States
                      318.523
    Name: G7 Population in millions, dtype: float64
                                                                           'United States': 318.523
                                                                       }, name='G7 Population in millions')
                                                              [16]

√ 0.0s

                                                                   Canada
                                                                                      35.467
                                                                   France
                                                                                      63.951
                                                                   Germany
                                                                                      80.940
                                                                   Italy
                                                                                      60.665
                                                                   Japan
                                                                                     127.061
                                                                   United Kingdom 64.511
                                                                   United States
                                                                                     318.523
                                                                   Name: G7 Population in millions, dtype: float64
```

Series: indexado

Así como en los diccionarios y las listas, puedo usar el índice para acceder a un elemento.

```
g7_pop
     ✓ 0.0s
[19]
                        35.467
    Canada
     France
                        63.951
                       80.940
    Germany
    Italy
                      60.665
    Japan
                      127.061
    United Kingdom
                      64.511
    United States
                      318.523
    Name: G7 Population in millions, dtype: float64
        g7_pop['Canada']
     ✓ 0.0s
[20]
    35.467
        g7_pop['Japan']

√ 0.0s

     127.061
```

Series: iloc

El comando **iloc** me permite utilizar los índices numéricos de los elementos. El último caso me devuelve una serie nueva.

```
g7_pop.iloc[0]
[22]
      ✓ 0.0s
    35.467
        g7 pop.iloc[-1]

√ 0.0s

[23]
     318.523
        g7_pop.iloc[0:2]
[27]
      ✓ 0.0s
    Canada
            35.467
            63.951
     France
    Name: G7 Population in millions, dtype: float64
```

Series: máscaras y slicing

También puedo trabajar con máscaras y slicing, en este caso el elemento superior también es incluido en el resultado.

```
g7_pop[['Italy', 'France']]
[24]
     ✓ 0.0s
    Italv
            60.665
    France 63.951
    Name: G7 Population in millions, dtype: float64
       g7 pop.iloc[[0, 1]]
[25]

√ 0.0s

            35.467
    Canada
    France 63.951
    Name: G7 Population in millions, dtype: float64
       g7_pop['Canada': 'Italy']
[26]
     ✓ 0.0s
             35.467
    Canada
    France 63.951
    Germany 80.940
    Italy
            60.665
    Name: G7 Population in millions, dtype: float64
```

Series: selección condicional

Se pueden aplicar condiciones booleanas al igual que en NumPy.

```
g7_pop
                                                                    g7_pop[g7_pop > 70]
[35]
      ✓ 0.0s

√ 0.0s

                                                            [37]
     Canada
                        35.467
                                                                Germany
                                                                                   80.940
     France
                        63.951
                                                                Japan
                                                                                  127.061
     Germany
                        80.940
                                                                United States
                                                                                  318.523
     Italy
                       60.665
                                                                Name: G7 Population in millions, dtype: float64
     Japan
                       127.061
    United Kingdom
                      64.511
    United States
                       318.523
    Name: G7 Population in millions, dtype: float64
                                                                    g7 pop.mean()
                                                            [38]
                                                                  ✓ 0.0s
                                                                107.30257142857144
        g7 pop > 70

√ 0.0s

                       False
     Canada
                                                                    g7 pop[g7 pop > g7 pop.mean()]
     France
                       False
                                                                  ✓ 0.0s
     Germany
                       True
     Italy
                       False
                                                                Japan
                                                                                  127.061
                                                                United States
                                                                                 318.523
     Japan
                       True
    United Kingdom
                       False
                                                                Name: G7 Population in millions, dtype: float64
     United States
                       True
    Name: G7 Population in millions, dtype: bool
```

Series: selección condicional

Se pueden aplicar condiciones booleanas al igual que en NumPy.

```
g7_pop[(g7_pop > 80) | (g7_pop < 40)]
 ✓ 0.0s
Canada
                 35.467
Germany
                80.940
Japan
                127.061
United States
              318.523
Name: G7 Population in millions, dtype: float64
   g7_pop[(g7_pop > 80) & (g7_pop < 200)]
 ✓ 0.0s
           80.940
Germany
Japan
          127.061
Name: G7 Population in millions, dtype: float64
```

Series: selección condicional

```
g7_pop
[43]
    ✓ 0.0s
    Canada
                       35.467
     France
                       63.951
    Germany
                     80.940
     Italy
                     60.665
                    127.061
     Japan
    United Kingdom 64.511
    United States 318.523
    Name: G7 Population in millions, dtype: float64
\triangleright \vee
        g7_pop.mean(), g7_pop.std()
[44]
     ✓ 0.0s
     (107.30257142857144, 97.24996987121581)
        g7_pop[(g7_pop > g7_pop.mean() - g7_pop.std() / 2) & (g7_pop < g7_pop.mean() + g7_pop.std() / 2)]
[42]
   ✓ 0.0s
    France
                       63.951
     Germany
                       80.940
     Italy
                  60.665
     Japan
             127.061
    United Kingdom 64.511
     Name: G7 Population in millions, dtype: float64
```

Series: operaciones vectorizadas

```
g7 pop * 1000000
                                                                      np.log(g7_pop)
                                                              [48]
                                                                    ✓ 0.0s
[47]
     ✓ 0.0s
                                                                   Canada
                                                                                     3.568603
    Canada
                        35467000.0
                                                                   France
                                                                                     4.158117
    France
                        63951000.0
                                                                   Germany
    Germany
                                                                                     4.393708
                       80940000.0
                                                                   Italy
                                                                                     4.105367
    Italy
                      60665000.0
                                                                   Japan
                                                                                     4.844667
    Japan
                      127061000.0
                                                                   United Kingdom
    United Kingdom
                      64511000.0
                                                                                   4.166836
                                                                   United States
    United States
                                                                                     5.763695
                     318523000.0
                                                                   Name: G7 Population in millions, dtype: float64
    Name: G7 Population in millions, dtype: float64
                                                                      np.exp(g7 pop)
        g7 pop + 30
                                                                    ✓ 0.0s
[54]
     ✓ 0.0s
                                                              [53]
                       65.467
                                                                   Canada
                                                                                      2.530011e+15
    Canada
                                                                                      5.936991e+27
    France
                       93.951
                                                                   France
    Germany
                      110.940
                                                                   Germany
                                                                                      1.418389e+35
                                                                                      2.220623e+26
                                                                   Italy
    Italy
                      90.665
                                                                   Japan
                                                                                      1.520167e+55
    Japan
                      157.061
                                                                   United Kingdom 1.039373e+28
    United Kingdom
                      94.511
                      348.523
                                                                   United States
                                                                                     2.151698e+138
    United States
                                                                   Name: G7 Population in millions, dtype: float64
    Name: G7 Population in millions, dtype: float64
```

Series: modificación

Hasta ahora vimos que al realizar operaciones sobre una serie de Pandas me devuelve una serie nueva. Veamos como modificar datos de un objeto existente.



Series: modificación

```
g7_pop[g7_pop < 70]
[64]
     ✓ 0.0s
    Canada
                40.500
    France
                63.951
    Italy
              60.665
    United Kingdom 64.511
    Name: G7 Population in millions, dtype: float64
       g7_pop[g7_pop < 70] = 99.99
[65]
     ✓ 0.0s
       g7_pop
[66]
     ✓ 0.0s
    Canada
                     99.990
    France
                    99.990
                 80.940
    Germany
    Italy
                 99.990
    Japan
              127.061
    United Kingdom 99.990
    United States 500.000
    Name: G7 Population in millions, dtype: float64
```

Data Frames

Los Data Frames son estructuras de Pandas más complejos que las Series. En lugar de tener una sola lista de valores, puedo tener varias columnas, y cada columna corresponde un mismo tipo de dato.

La información contenida en un data frame puede ser ingresada manualmente, pero lo mas común es importarla desde algún archivo que trabaje con datos tabulares, como ser el CSV.

		G7 Sta	ts		
	Population	GDP	Surface	HDI	Continent
Canada	35.467	1,785,387.00	9,984,670	0.913	America
France	63.951	2,833,687.00	640,679	0.888	Europe
Germany	80.94	3,874,437.00	357,114	0.916	Europe
Italy	60.665	2,167,744.00	301,336	0.873	Europe
Japan	127.061	4,602,367.00	377,930	0.891	Asia
United Kingdom	64.511	2,950,039.00	242,495	0.907	Europe
United States	318.523	17,348,075.00	9,525,067	0.915	America

Data Frames: creación

	Population	GDP	Surface Area	HDI	Continent
0	35.467	1785387	9984670	0.913	America
1	63.951	2833687	640679	0.888	Europe
2	80.940	3874437	357114	0.916	Europe
3	60.665	2167744	301336	0.873	Europe
4	127.061	4602367	377930	0.891	Asia
5	64.511	2950039	242495	0.907	Europe
6	318.523	17348075	9525067	0.915	America

Data Frames: indices

Puedo tener un resumen de los datos en el DF.

	Population	GDP	Surface Area	HDI	Continent
Canada	35.467	1785387	9984670	0.913	America
France	63.951	2833687	640679	0.888	Europe
Germany	80.940	3874437	357114	0.916	Europe
Italy	60.665	2167744	301336	0.873	Europe
Japan	127.061	4602367	377930	0.891	Asia
United Kingdom	64.511	2950039	242495	0.907	Europe
United States	318.523	17348075	9525067	0.915	America

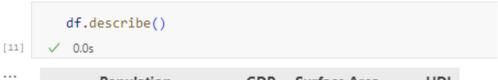
```
df.info()
 ✓ 0.0s
<class 'pandas.core.frame.DataFrame'>
Index: 7 entries, Canada to United States
Data columns (total 5 columns):
 # Column
                Non-Null Count Dtype
  Population 7 non-null
                              float64
 1 GDP
          7 non-null
                              int64
 2 Surface Area 7 non-null
                              int64
          7 non-null float64
    HDT
  Continent 7 non-null
                              object
dtypes: float64(2), int64(2), object(1)
memory usage: 336.0+ bytes
```

Data Frames: dimensiones

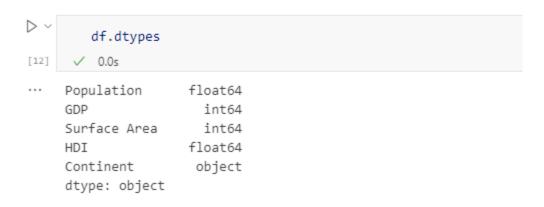
Para conocer las dimensiones del data frame:

```
df.size
[9]
     ✓ 0.0s
   35
        df.shape
[10]
     ✓ 0.0s
   (7, 5)
        df.columns
     ✓ 0.0s
[6]
    Index(['Population', 'GDP', 'Surface Area', 'HDI', 'Continent'], dtype='object')
        df.index
     ✓ 0.0s
    Index(['Canada', 'France', 'Germany', 'Italy', 'Japan', 'United Kingdom',
            'United States'],
          dtype='object')
```

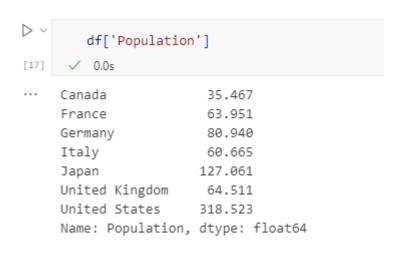
Data Frames: información estadística



	Population	GDP	Surface Area	HDI
count	7.000000	7.000000e+00	7.000000e+00	7.000000
mean	107.302571	5.080248e+06	3.061327e+06	0.900429
std	97.249970	5.494020e+06	4.576187e+06	0.016592
min	35.467000	1.785387e+06	2.424950e+05	0.873000
25%	62.308000	2.500716e+06	3.292250e+05	0.889500
50%	64.511000	2.950039e+06	3.779300e+05	0.907000
75%	104.000500	4.238402e+06	5.082873e+06	0.914000
max	318.523000	1.734808e+07	9.984670e+06	0.916000



Data Frames: filtros





··· pandas.core.series.Series

	df[['Population',	'GDP']]
[21]	✓ 0.0s	

	Population	GDP
Canada	35.467	1785387
France	63.951	2833687
Germany	80.940	3874437
Italy	60.665	2167744
Japan	127.061	4602367
United Kingdom	64.511	2950039
United States	318.523	17348075

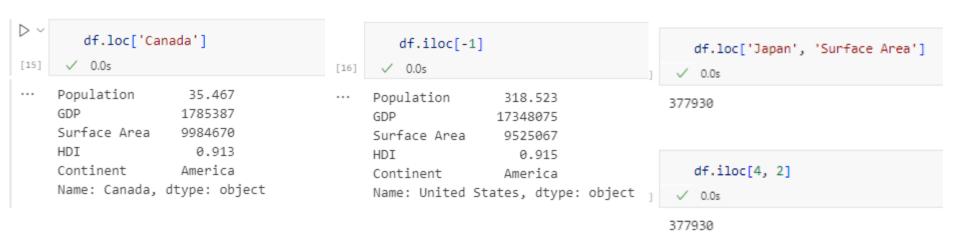
df[:3]				
✓ 0.0s				

	Population	GDP	Surface Area	HDI	Continent
Canada	35.467	1785387	9984670	0.913	America
France	63.951	2833687	640679	0.888	Europe
Germany	80.940	3874437	357114	0.916	Europe

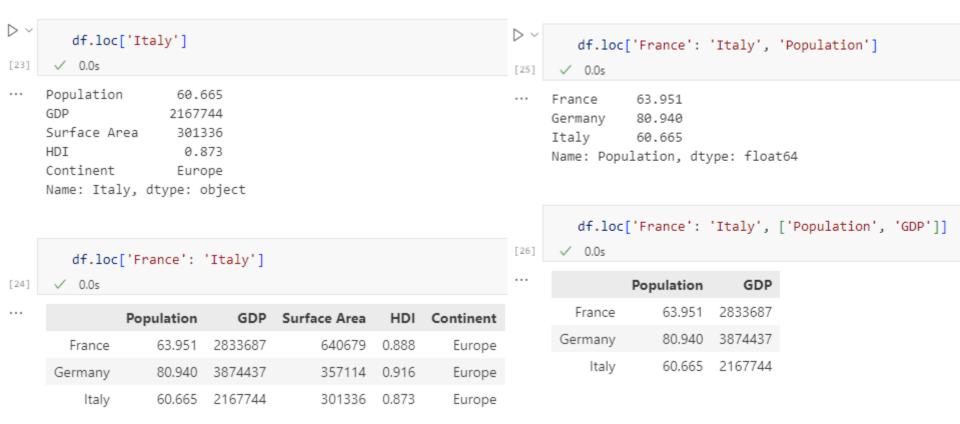
Data Frames: loc y iloc

Con **loc** puedo filtrar según el nombre del índice o de la columna.

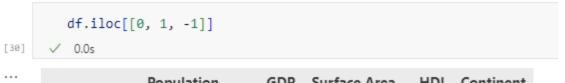
Con loc puedo filtrar según el valor numérico del índice o de la columna.



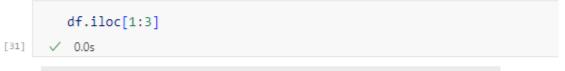
Data Frames: loc



Data Frames: iloc



Continent	HDI	Surface Area	GDP	Population	
America	0.913	9984670	1785387	35.467	Canada
Europe	0.888	640679	2833687	63.951	France
America	0.915	9525067	17348075	318.523	United States



	Population	GDP	Surface Area	HDI	Continent
France	63.951	2833687	640679	0.888	Europe
Germany	80.940	3874437	357114	0.916	Europe

	df.iloc[1:3, 3]
[32]	✓ 0.0s

··· France 0.888 Germany 0.916

Name: HDI, dtype: float64

[33]	df.ilo ✓ 0.0s	c[1:3, [0,	3]]
		Population	HDI
	France	63.951	0.888
	Germany	80.940	0.916

	df.iloc[1:3, 1:3]	
[34]	✓ 0.0s	

	GDP	Surface Area
France	2833687	640679
Germany	3874437	357114

Data Frames: selección condicional

```
df

√ 0.0s

[35]
                      Population
                                       GDP
                                             Surface Area
                                                            HDI
                                                                 Continent
             Canada
                          35.467
                                   1785387
                                                 9984670
                                                           0.913
                                                                    America
                          63.951
                                   2833687
                                                           0.888
              France
                                                  640679
                                                                     Europe
            Germany
                          80.940
                                   3874437
                                                  357114
                                                           0.916
                                                                     Europe
                                   2167744
                Italy
                          60.665
                                                  301336
                                                           0.873
                                                                     Europe
                         127.061
                                   4602367
                                                           0.891
                                                                       Asia
               Japan
                                                  377930
                                   2950039
      United Kingdom
                          64.511
                                                  242495
                                                           0.907
                                                                     Europe
                         318.523
        United States
                                 17348075
                                                 9525067 0.915
                                                                    America
        df['Population'] > 70

√ 0.0s

[36]
                        False
     Canada
     France
                        False
     Germany
                         True
     Italy
                        False
     Japan
                         True
     United Kingdom
                        False
     United States
                         True
     Name: Population, dtype: bool
```

Data Frames: selección condicional

```
df.loc[df['Population'] > 70, 'Population']

0.0s
```

Germany 80.940 Japan 127.061 United States 318.523

Name: Population, dtype: float64

```
df.loc[df['Population'] > 70, ['Population', 'GDP']]

$\square$ 0.0s
```

	Population	GDP
Germany	80.940	3874437
Japan	127.061	4602367
United States	318.523	17348075

df.loc[df['Population'] > 70]

	Population	GDP	Surface Area	HDI	Continent
Germany	80.940	3874437	357114	0.916	Europe
Japan	127.061	4602367	377930	0.891	Asia
United States	318.523	17348075	9525067	0.915	America

Data Frames: borrado de datos

[63] V 0.0s

Estos métodos devuelven un data frame nuevo, para modificar el DF existente se usa el parámetro inplace.

```
df.drop('Canada')
                                                           df.drop(['Population', 'HDI'], axis=1)
df.drop(['Canada', 'Japan'])
                                                           df.drop(['Population', 'HDI'], axis=1)
df.drop(columns=['Population', 'HDI'])
                                                           df.drop(['Population', 'HDI'], axis='columns')
df.drop(['Italy', 'Canada'], axis=0)
                                                           df.drop(['Canada', 'Germany'], axis='rows')
                                   df.drop(columns='Language', inplace=True)
```

Data Frames: borrado de datos

Estos métodos devuelven un data frame nuevo, para modificar el DF existente se usa el parámetro inplace.

HDI Continent

Europe

Europe

Europe

Europe

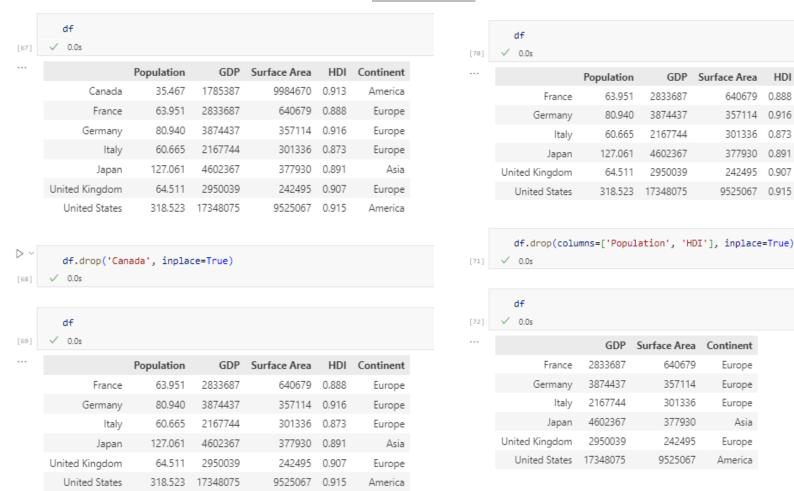
America

Asia

0.888

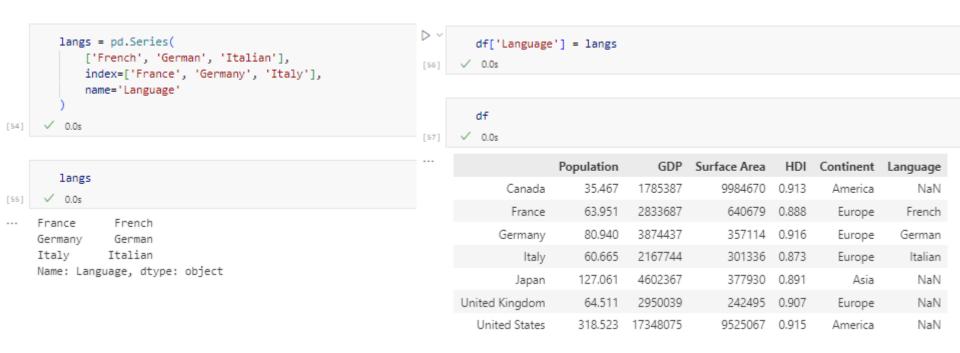
0.891

0.907



Data Frames: añadir una nueva columna

Puedo crear una nueva columna a partir de una serie. En caso de haber datos faltantes, Panda los rellena con NaN.



Data Frames: renombrar una columna o índice

Con el comando rename puedo cambiar el nombre de una columna o un índice existente.

	Population	GDP	Surface Area	Human Development Index	Continent	Language
Canada	35.467	1785387	9984670	0.913	America	English
France	63.951	2833687	640679	0.888	Europe	English
Germany	80.940	3874437	357114	0.916	Europe	English
Italy	60.665	2167744	301336	0.873	Europe	English
Japan	127.061	4602367	377930	0.891	Asia	English
UK	64.511	2950039	242495	0.907	Europe	English
USA	318.523	17348075	9525067	0.915	America	English

Data Frames: añadir nuevas filas

Puedo usar el comando concat para concatenar dos DF distintos.

		df					
7]	~	0.0s					
			Population	GDP	Surface Area	HDI	Continent
		Canada	35.467	1785387	9984670	0.913	America
		France	63.951	2833687	640679	0.888	Europe
		Germany	80.940	3874437	357114	0.916	Europe
		Italy	60.665	2167744	301336	0.873	Europe
		Japan	127.061	4602367	377930	0.891	Asia
	Uni	ited Kingdom	64.511	2950039	242495	0.907	Europe
		United States	318.523	17348075	9525067	0.915	America
~					: 1409670, 'G	DP': 35	291000}, i
7]	(df3 = pd.Dat df2 = pd.con df2 0.0s			: 1409670, 'Gi	DP': 35	291000}, i
	(df2 = pd.con df2					Continent
']	(df2 = pd.con df2	cat([df, df	31)		HDI	
7]	(df2 = pd.con df2 0.0s	Population	GDP	Surface Area	HDI 0.913	Continent
7]	(df2 = pd.con df2 0.0s Canada	Population 35,467	GDP 1785387	Surface Area 9984670.0	HDI 0.913 0.888	Continent America
']	(df2 = pd.con df2 0.0s Canada France	Population 35.467 63.951	GDP 1785387 2833687	Surface Area 9984670.0 640679.0	HDI 0.913 0.888 0.916	Continent America Europe
7]	(df2 = pd.con df2 0.0s Canada France Germany	Population 35,467 63,951 80,940	GDP 1785387 2833687 3874437	Surface Area 9984670.0 640679.0 357114.0	HDI 0.913 0.888 0.916 0.873	Continent America Europe Europe
']	~	df2 = pd.con df2 0.0s Canada France Germany	Population 35.467 63.951 80.940 60.665	GDP 1785387 2833687 3874437 2167744	Surface Area 9984670.0 640679.0 357114.0 301336.0	HDI 0.913 0.888 0.916 0.873 0.891	Continent America Europe Europe Europe
	Uni	df2 = pd.con df2 0.0s Canada France Germany Italy Japan	Population 35,467 63,951 80,940 60,665 127,061	GDP 1785387 2833687 3874437 2167744 4602367	Surface Area 9984670.0 640679.0 357114.0 301336.0 377930.0	HDI 0.913 0.888 0.916 0.873 0.891 0.907	Continent America Europe Europe Europe Asia

Data Frames: añadir nuevas filas

En caso de solo tener índices numéricos, los debo ignorar al hacer la

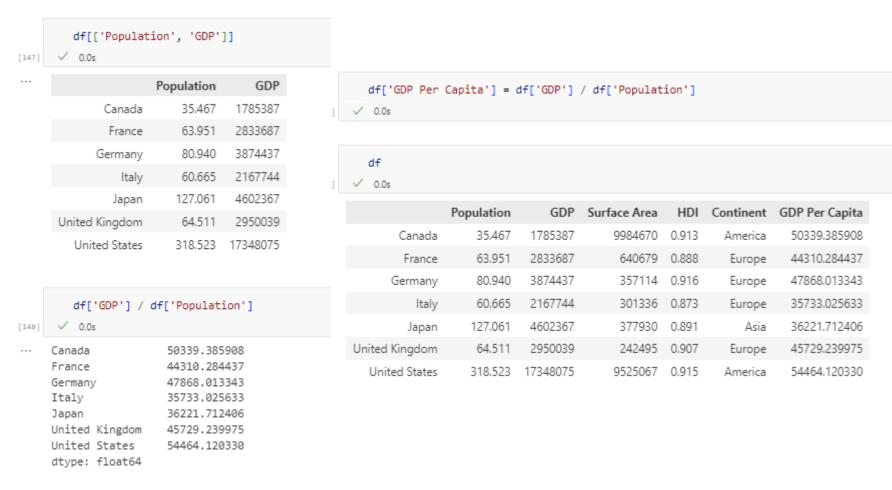
concatenación.

	Name	Maths	Science
0	Martha	87	83
1	Tim	91	99
2	Rob	97	84
3	Georgia	95	76

	Name	Maths	Science
0	Martha	87	83
1	Tim	91	99
2	Rob	97	84
3	Georgia	95	76
4	Amy	89	93

Data Frames: calcular una nueva columna

Puedo crear columnas nuevas a partir de columnas existentes.



En ciencia de datos y Machine Learning es común que los datos se transfieran en un tipo de archivo desprovisto de formato (a diferencia, por ejemplo, de un archivo Excel) para mejorar la compatibilidad entre librerías y programas.

Un formato de archivo de este tipo es el de extensión CSV (comma separated values). Es parecido a un archivo .XLSX en el cual la información tabular está contenida en varias filas pero en una sola columna, donde se utiliza una coma (,) para significar la separación de datos (similar a los elementos en una lista de Python).

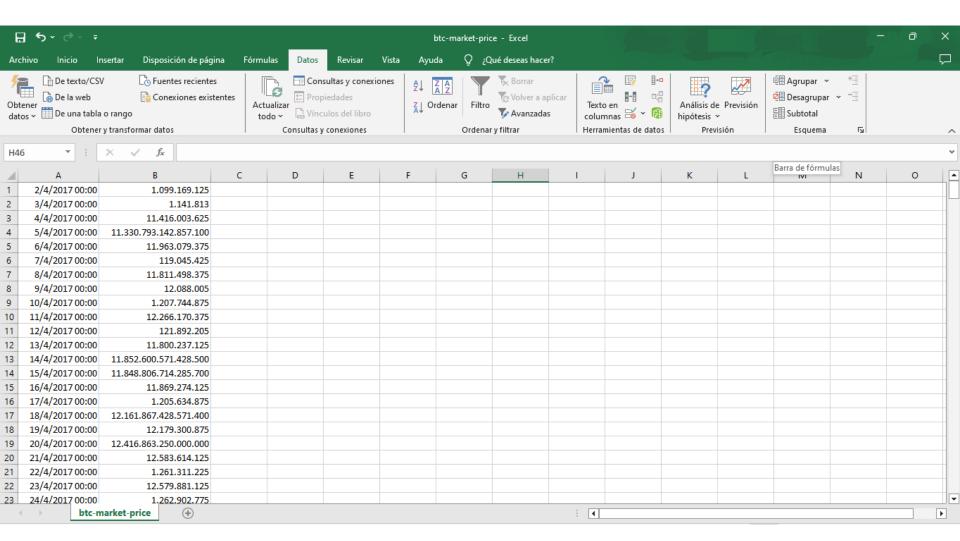
Por lo general, en la primer fila del archivo, se tiene una descripción del contenido de las columnas, lo que se conoce como header.

- A la hora de importar se pueden definir varios parámetros para que el DF tenga el formato que deseamos.
- La primera fila del mismo se interpeta como la que contiene información de los encabezados, por lo que se debe ignorar si no los tuviere.
- También puedo establecer una columna como la que define los índices.
- Puedo establecer el delimitador de columnas (coma, punto y coma).
- Si quiero ignorar algunas filas o columnas.
- Establecer tipos de datos.

pandas.read_csv

```
pandas.read csv(filepath_or_buffer, *, sep= NoDefault.no_default,
delimiter=None, header='infer', names= NoDefault.no default, index col=None,
usecols=None, dtype=None, engine=None, converters=None, true values=None,
false values=None, skipinitialspace=False, skiprows=None, skipfooter=0,
nrows=None, na_values=None, keep_default_na=True, na_filter=True,
verbose = NoDefault.no default, skip blank lines = True, parse dates = None,
infer datetime format = NoDefault.no default,
keep_date_col= NoDefault.no default, date_parser= NoDefault.no default,
date format=None, dayfirst=False, cache dates=True, iterator=False,
chunksize=None, compression='infer', thousands=None, decimal='.',
lineterminator=None, quotechar='"', quoting=0, doublequote=True,
escapechar=None, comment=None, encoding=None, encoding errors='strict',
dialect=None, on_bad_lines='error', delim_whitespace= NoDefault.no default.
low_memory=True, memory_map=False, float_precision=None,
                                                                     [source]
storage options=None, dtype backend= NoDefault.no default)
```

Read a comma-separated values (csv) file into DataFrame.



```
df = pd.read_csv(
             'data/btc-market-price.csv',
            header=None,
             names=['Timestamp', 'Price'],
             parse_dates=True

√ 0.0s

[38]
        df.head()
[31]
      ✓ 0.0s
                                   Price
                Timestamp
         2017-04-02 00:00:00
                            1099.169125
         2017-04-03 00:00:00
                            1141.813000
        2017-04-04 00:00:00
                           1141.600363
      3 2017-04-05 00:00:00 1133.079314
      4 2017-04-06 00:00:00 1196.307937
```