

DATA SELECTION

```
In [18]: import pandas as pd
import numpy as np

In [19]: datos = pd.date_range('1/1/2008', periods=8)
df = pd.DataFrame({'a':datos, 'b':1, 'c':2, 'd':1})
df

Out[19]:
   a      b  c  d
0 2008-01-01  1  2  1
1 2008-01-02  1  2  1
2 2008-01-03  1  2  1
3 2008-01-04  1  2  1
4 2008-01-05  1  2  1
5 2008-01-06  1  2  1
6 2008-01-07  1  2  1
7 2008-01-08  1  2  1

Out[19]: Series([], Freq: D, Name: b, dtype: float64)
```

Llamar a un valor de una dataframe

Podemos seleccionar un elemento de los datos de distintas formas, es decir, para llamar a un elemento concreto de un data frame o series tenemos distintos elementos. Ahora nos centramos en los dataframes, pero para las series se hace de la misma forma:

a) Llamaros a una columna y después a un índice

Por ejemplo si queremos el sexto elemento de la columna A:

```
In [19]: df['A'][datos[5]]

Out[19]: 1.2806685948989247

b) Llamaros a un índice y después a la columna con loc[]
```

```
In [11]: df.loc[datos[5], 'A']

Out[11]: 1.2806685948989247
```

EJEMPLO: Cambiar el orden de 2 columnas

Para ello simplemente llamamos a 2 columnas, y las igualamos a las columnas al revés. Es decir:

```
In [12]: df[['A','B']] = df[['B','A']]

Out[12]:
   A      B  C  D
0 2008-01-01  1.000000  0.000000  0.000000  1.889646
1 2008-01-02  1.566900  0.203986  1.576848  1.008718
2 2008-01-03  1.422871  0.442027  0.393688  0.793887
3 2008-01-04  0.964930  0.458422  0.060079  0.508323
4 2008-01-05  0.564930  0.403422  0.090079  0.508323
5 2008-01-06  0.564930  0.203986  0.179829  0.469625
6 2008-01-07  1.444701  0.217499  0.274414  0.698038
7 2008-01-08  0.967001  1.365496  0.688785  0.775784
```

Llamados por índice, loc[] y iloc[]

Basicamente hacen lo mismo pero loc[] únicamente funciona con el nombre del índice mientras que iloc[] solo funciona con la posición en número:

```
In [13]: from IPython.display import Image
Image('loc_iloc.png')

Out[13]:


|                 | loc                | iloc               |
|-----------------|--------------------|--------------------|
| Input           | 5 or 'a'           | 5                  |
| List of values  | ['a', 'b', 'c']    | [4, 3, 0]          |
| Slice object    | a:b                | 1:7                |
| Boolean array   | [True, False, ...] | [True, False, ...] |
| Callable object | Function           | Function           |


```

Llamar elementos a partir de booleanos

Basicamente, metes la condición dentro del corchete y lo que de verdad es, es decir, cuando se selecciona el valor correspondiente. Por ejemplo, si hacemos un dataframe loc[dataframe.columns>0], lo que hace es comparar todos los elementos de la columna a 0, dando por ejemplo [True, True, False, True]. Esto hará que se impriman la primera, segunda y cuarta fila por parte del loc.

```
In [14]: df.loc[df.A > 0]

Out[14]:
   A      B  C  D
0 2008-01-02  1.566900  0.203986  1.576848  1.008718
1 2008-01-03  1.422871  0.442027  0.393688  0.793887
2 2008-01-04  0.964930  0.458422  0.060079  0.508323
3 2008-01-05  0.564930  0.403422  0.090079  0.508323
4 2008-01-06  0.564930  0.203986  0.179829  0.469625
5 2008-01-07  1.444701  0.217499  0.274414  0.698038
6 2008-01-08  0.967001  1.365496  0.688785  0.775784
```

IMPORTANTE: Tienes que usar dataframe.columna, no vale poner directamente el nombre de la columna, porque > o < o = no pueden ir entre un string y un int.

Filtrar números de una tabla, poniendo NaN a los que no cumplan la condición

El comando es dataframe[dataframe > 34], lo que hace la condición del corchete es crear un dataframe de booleanos, y después los valores True conservan su valor mientras que los False pasan a perder el valor(NaN).

Por ejemplo, filtrar del dataframe los valores mayores de 0:

```
In [15]: df[df > 0]

Out[15]:
   A      B  C  D
0 2008-01-01  NaN  NaN  NaN  1.889646
1 2008-01-02  1.566900  0.203986  1.576848  1.008718
2 2008-01-03  NaN  NaN  NaN  0.793887
3 2008-01-04  NaN  NaN  NaN  0.508323
4 2008-01-05  NaN  1.964107  0.811359  0.469625
5 2008-01-06  0.669350  NaN  0.179829  NaN
6 2008-01-07  1.444701  NaN  NaN  NaN
7 2008-01-08  NaN  1.365496  0.688785  NaN
```

Dar la vuelta a las columnas y a las filas

df_ej5 = pd.DataFrame(np.random.randn(8,4), index=datos, columns=['A','B','C','D'])

df_ej5.loc[:,:] #esto de la vuelta a las columnas

df.loc[:,:] #esto de la vuelta a las filas

df.loc[:,:] .reset_index(drop = True) #devuelve las filas pero no los índices

```
Out[15]:
   A      B  C  D
0 0.567001  1.365486  0.488785 -0.775784
1 1.444701  1.273499 -0.044715 -0.688784
2 0.669350 -1.288661  0.179829 -0.462134
3 0.550484 -1.944037  0.811359  0.469625
4 0.564930  0.458422  0.060079  0.508323
5 1.422871  0.442027  0.393688 -0.793887
6 1.565900  0.203986 -1.576848  1.008718
7 0.669350  0.203986  0.641780  1.889646
```

Mostrar los valores de la columna B que son más grandes que la columna A y más pequeños que la columna C

```
In [17]: datos = pd.date_range('1/1/2008', periods=8)
df = pd.DataFrame(np.random.randn(8,4), index=datos, columns=['A','B','C','D'])

print(df)
df.loc[(df.B>df.A) & (df.B<df.C), 'B'] #si quieres poner varias condiciones
#tienes que poner cada una entre
paréntesis y el and es &, el or so-
#poniendo así sera | para no sé
```

```
Out[17]:
   A      B  C  D
0 2008-01-01  1.281312  0.187812  0.278424  0.329489
1 2008-01-02  0.288862  1.014449  0.526872  0.114387
2 2008-01-03  0.453267  1.807189  0.388187  0.591719
3 2008-01-04 -1.249429  0.852616  0.123124  0.344837
4 2008-01-05  0.345910  0.477611  0.175903  0.899586
5 2008-01-06 -0.776188 -1.187224  0.488126  1.172228
6 2008-01-07  0.185135  0.999801  1.058420  0.127483
7 2008-01-08 -0.529417  0.210373  1.742735  1.469252
Freq: D, Name: B, dtype: float64
```

SELECTION BY DF-QUERY(CONDITION) (muy útil)

dataframe.query('condition') es una herramienta muy útil que nos muestra los valores que cumplen la condición que se indique, simplemente por columnas. Es decir, si ponemos df.query('a>0') nos mostrará por filas los valores de la columna a que sean mayores que 0, es el equivalente a df.loc[df.A>0.]

DATA MANIPULATION IN PANDAS

RENAMING

dataframe.rename() permite renombrar índices o columnas, y ello se hace a través de un diccionario. Es decir, si queremos por ejemplo renombrar las columnas de la tabla df:

```
In [18]: df.rename(columns={'A':'One','B':'Two','C':'Three','D':'Four'})

Out[18]:
   One  Two  Three  Four
0 2008-01-01 -1.251115  1.027412 -0.778824  0.329489
1 2008-01-02  0.209022 -1.039489  0.616072 -0.114387
2 2008-01-03 -0.404914  1.407709  0.388187  0.591719
3 2008-01-04 -1.348420  0.852616 -0.131124  0.344837
4 2008-01-05 -1.319116  0.477611 -0.173063  0.899586
5 2008-01-06 -0.776188 -1.187224  0.488126  1.172228
6 2008-01-07  0.180130 -1.369850 -1.084442  0.727483
7 2008-01-08 -0.529417  0.210373 -1.742735  1.469252
```

Ahora hacemos lo mismo para los índices:

```
In [19]: df = pd.DataFrame(np.random.randn(8,4), index=[1,2,3,4,5,6,7,8], columns=['A','B','C','D'])
df.rename(index={'a':'1','b':'3','c':'4','d':'6','e':'6','f':'7','g':'8','h':'8'})

Out[19]:
   A      B  C  D
a -0.940880  0.946700  0.890988  0.294990
b -1.548842  0.041309  0.317861  1.436887
c -0.002279  0.653818  0.087386  1.408887
d -0.372120  0.548446  1.352986  1.341535
e -0.981832  0.068498  2.677317  0.205802
f -0.029465  1.134713  0.337710  1.597972
g 0.197687  1.327264  1.296326  0.318335
h -1.057785  0.802463 -0.018320 -0.032486
```

REINDEXING

dataframe.reindex(nuevo_indice,1_nuevo_indice, 2_). Permite cambiar o añadir nuevos índices, ESA ES LA DIFERENCIA ENTRE INDEX EN LA TUPLA Y ESTOS SE AÑADIRÁN EN ORDEN DE DEFENEAQ. También funciona para columnas, pero es necesario especificarlo con dataframe.reindex(columns= nueva_columna, 1_nueva_columna, 2_)

```
In [20]: a = df.reindex([1,2,3,4,5,6,7,8,9,10])
a

Out[20]:
   A      B  C  D
1 -0.940880  0.946700  0.890988  0.294990
2 -1.548842  0.041309  0.317861  1.436887
3 -0.002279  0.653818  0.087386  1.408887
4 -0.372120  0.548446  1.352986  1.341535
5 -0.981832  0.068498  2.677317  0.205802
6 -0.029465  1.134713  0.337710  1.597972
7 0.197687  1.327264  1.296326  0.318335
8 -1.057785  0.802463 -0.018320 -0.032486
9 NaN  NaN  NaN  NaN
10 NaN  NaN  NaN  NaN
```

DROPPING ENTRIES

drop(ind_de_fila_a_eliminar,1_indice_de_fila_a_eliminar, 2_)/drop(columna_a_eliminar,1_tubo_columna_a_eliminar, 2_axis=1)

Parando de la siguiente tabla:

```
In [21]: df.drop(df.DataFrame(np.arange(16).reshape((4,4))), index=['Ohio','Colorado','Utah','New York'], columns=['one','two','three','four'])
df_ohio

Out[21]:
   one  two  three  four
Ohio   0   1   2   3
Colorado  4   5   6   7
Utah   8   9  10  11
New York 12  13  14  15
```

Eliminamos las filas 'Ohio' y 'Colorado'

```
In [22]: df.drop(drop(['Ohio','Colorado']))

Out[22]:
   one  two  three  four
Utah   8   9  10  11
New York 12  13  14  15
```

Eliminamos las columnas 'two' y 'four':

```
In [23]: df.drop(drop(['two','four'],axis=1))

Out[23]:
   one  three
Ohio   0   2
Colorado  4   6
Utah   8  10
New York 12  14
```

ARITHMETIC AND DATA ALIGNMENT

Queremos hacer operaciones aritméticas sobre 2 tablas y juntarlas, por ejemplo si queremos sumar 2 tablas, hay que tener en cuenta que si hacemos df1+df2 ambas tablas deben tener el mismo tamaño, porque los valores que no tengan una pareja tomarán el valor de NaN. Por tanto era evitar ésto tenemos los operadores aritméticos de pandas:

```
In [24]: from IPython.display import Image
Image('operadores-aritmeticos.png')

Out[24]:


| Method | Description                   |
|--------|-------------------------------|
| add    | Method for addition (+)       |
| sub    | Method for subtraction (-)    |
| div    | Method for division (/)       |
| mul    | Method for multiplication (*) |


```

```
In [25]: df1 = pd.DataFrame(np.arange(12).reshape((3,4)), columns=list('abcd'))
df2 = pd.DataFrame(np.arange(12).reshape((4,5)), columns=list('abcde'))

print(df1.add(df2, fill_value=0))
print(df1)
print(df1.sub(df2, fill_value=0))
print(df1)
print(df1.div(df2, fill_value=0))
print(df1.mul(df2, fill_value=0))
```

```
#Es importante el fill_value=0, porque si no lo pones es como si
#hicieras por ejemplo df1*df2

   a  b  c  d  e
0  0  0  0  0  0
1  0  1  0  1  0
2  0  2  0  2  0
3  0  3  0  3  0
4  0  4  0  4  0
5  0  5  0  5  0
6  0  6  0  6  0
7  0  7  0  7  0
8  0  8  0  8  0
9  0  9  0  9  0
10 0 10 0 10 0
11 0 11 0 11 0
12 0 12 0 12 0
13 0 13 0 13 0
14 0 14 0 14 0
15 0 15 0 15 0
```

FUNCTION APPLICATION AND MAPPING

dataframe.apply(function, axis=1) -> Se ejecuta la función entre paréntesis usando como input el dataframe especificado, PERO SE EJECUTA POR COLUMNAS

dataframe.apply(function, axis=2) -> Se ejecuta la función entre paréntesis usando como input el dataframe especificado, PERO SE EJECUTA POR FILAS

```
In [26]: f = lambda x: x.max()-x.min() #la función

frame = pd.DataFrame(np.random.randn(4,3), columns=['b','d','e'], index=['utah','Ohio','Texas','Oregon'])

print(frame.apply(f)) #Se busca el mayor valor de la columna y se le resta el
#menor de la misma, siendo el output una serie con cada
#columna

print(frame.apply(f, axis=1)) #Se busca el mayor valor de la fila y se le resta el
#menor de la misma, siendo el output una serie con cada
#fila
```

```
b      2.76397
d      2.46926
e      1.48440
utah      1.74891
Ohio      1.50756
Texas     0.53977
Oregon    1.54571
dtype: float64
```

SORTING AND RANKING

```
In [27]: df.sort DataFrame({'A':[0,1,0,1], 'b':[4,7,-3,2], index=[0,1,2,3]})

Sorting

Out[27]:
   A      b
0 0 4
1 1 7
2 0 -3
3 1 2
```

```
In [28]: df5.sort_values(by=['a','b']) #órdenar respecto a las columnas 'a' y 'b'

Out[28]:
   a  b
2 0 -3
0 0 4
3 1 2
1 1 7
```

```
In [29]: df5.sort_values(by=['a','b']) #órdenar respecto a las columnas 'a' y 'b'

Out[29]:
   a  b
2 0 -3
0 0 4
3 1 2
1 1 7
```

Ranking

El ranking asigna un valor de 1 al número de datos no nulos en una serie/habla. El ranking se puede aplicar a una serie o a un dataframe.

Existen distintos métodos de hacer ranking, que se marcan por .rank(method=)

```
In [30]: from IPython.display import Image
Image('metodos_ranking.png')

Out[30]:


| Method    | Description                                                       |
|-----------|-------------------------------------------------------------------|
| 'average' | Default: assign the average rank to each entry in the equal group |
| 'min'     | Use the minimum rank for the whole group                          |
| 'max'     | Use the maximum rank for the whole group                          |
| 'first'   | Assign ranks in the order the values appear in the data           |


```

```
In [31]: obj.rank(Series([7,-5,7,4,2,6,4,1])

obj.rank()

Out[31]:
0 6.5
1 6.5
2 6.5
3 6.5
4 3.8
5 2.8
6 4.5
dtype: float64
```

También podemos hacer un ranking de dataframes, pero éste se puede hacer por filas o por columnas. Por defecto se hará por filas, por lo que deberemos especificar con axis=1 para que se haga por columnas:

```
In [32]: frame.sort DataFrame({'b':[4,3,7,-3,2], 'a':[0,1,0,1], 'c':[-2,5,8,-2,5]})
print(frame)

print(frame.rank(method='first',axis=1)) #Ranking por columnas
print(frame.rank(method='first') #Ranking por filas
```

```
b      a      c
0  4  0  -2
1  3  1  5
2  7  0  8
3 -3  1  5
4  2  0  1
5  5  1  0
6  1  3  0
7  0  1  0
8  0  3  0
9  0  2  0
10 0  2  0
11 0  2  0
12 0  2  0
13 0  2  0
14 0  2  0
15 0  2  0
```

EJERCICIOS

EJERCICIO 1: Write a Pandas program to replace all the NaN values with Zero's in a column of a dataframe.

```
In [33]: df = pd.DataFrame(np.random.randn(8,4), index=datos, columns=['A','B','C','D'])
df.name(df.name+'@CASERO LA TABLA DE SERVIDA DE ENUNCIADO')

df_nan[df_nan.isnull()]=0
df_nan
```

```
Out[33]:
   A      B  C  D
0 2008-01-01  0.312890  0.422461  0.223761  0.000000
1 2008-01-02  0.000000  0.000000  0.000000  1.030336
2 2008-01-03  0.976438  0.402716  0.000000  0.000000
3 2008-01-04  0.000000  0.637970  0.538376  0.000000
4 2008-01-05  0.000000  0.000000  0.000000  0.654955
5 2008-01-06  0.000000  0.000000  0.000000  0.000000
6 2008-01-07  0.879621  0.000000  0.807574  1.533723
7 2008-01-08  0.000000  1.122834  0.247540  0.000000
```

EJERCICIO 2: Write a Pandas program to convert index in a column of the given dataframe.

```
In [27]: df_ej2 = pd.DataFrame(np.random.randn(8,4), index=datos, columns=['A','B','C','D'])
print(df_ej2)

df_ej2['nueva_columna']=df_ej2.index
df_ej2
```

```
Out[27]:
   A      B  C  D  nueva_columna
0 2008-01-01  0.535105  0.384661  0.530176  0.171163
1 2008-01-02  0.489148  0.584518  0.358421  0.138270
2 2008-01-03  0.489939  0.165032  1.735788  1.578797
3 2008-01-04  1.385771  0.995058  0.125143  0.817820
4 2008-01-05  0.271717  0.272448  0.384922  0.651685
5 2008-01-06  1.372558  0.763122  1.137484  1.854456
6 2008-01-07  0.621777  0.199946  0.144280  2.264885
7 2008-01-08  0.774640  0.887146  1.157412  1.855113
```

```
Out[27]:
   A      B  C  D
0 2008-01-01  0.535105  0.384661  0.530176  0.171163
1 2008-01-02  0.489148  0.584518  0.358421  0.138270
2 2008-01-03  0.489939  0.165032  1.735788  1.578797
3 2008-01-04  1.385771  0.995058  0.125143  0.817820
4 2008-01-05  0.271717  0.272448  0.384922  0.651685
5 2008-01-06  1.372558  0.763122  1.137484  1.854456
6 2008-01-07  0.621777  0.199946  0.144280  2.264885
7 2008-01-08  0.774640  0.887146  1.157412  1.855113
```

EJERCICIO 3: Write a Pandas program to count the NaN values in one or more columns of a dataframe.

```
In [28]: df_ej3 = pd.DataFrame(np.random.randn(8,4), index=datos, columns=['A','B','C','D'])
df_ej3[df_ej3['C'].isnull()]
print(df_ej3)

NaN
df_ej3[df_ej3['C'].isnull()]
for i in df_ej3.index:
    if df_ej3.loc[i, 'C'] == True:
        count = count + 1

5/9
```

```
Out[28]:
   A      B  C  D
0 2008-01-01  0.138296  1.347074  1.389182  NaN
1 2008-01-02  0.250521  0.584713  0.817820  NaN
2 2008-01-03  0.639996  0.848488  0.817820  NaN
3 2008-01-04  NaN  NaN  1.385165  0.767852
4 2008-01-05  NaN  NaN  2.026519  0.899586
5 2008-01-06  NaN  0.487218  NaN  NaN
6 2008-01-07  0.218746  0.995058  NaN  0.888114
7 2008-01-08  0.389563  0.995058  NaN  0.788788
```

EJERCICIO 4: Write a Pandas program to add a prefix or suffix to all columns of a given DataFrame

```
In [29]: df_ej4 = pd.DataFrame(np.random.randn(8,4), index=datos, columns=['A','B','C','D'])

pref = lambda x:
columns = df_ej4.columns
pref_col = pref + columns
pref_col = pref_col + ' valores'

df_ej4.rename(df_ej4.index(columns+pref_col)) #Se haría así pero el reindex no se va para columnas

EJERCICIO 5: Write a Pandas program to reverse order (rows, column) of a given DataFrame
```

```
In [30]: df_ej5 = pd.DataFrame(np.random.randn(8,4), index=datos, columns=['A','B','C','D'])
df_ej5.loc[:,:]

Out[30]:
   A      B  C  D
0 2008-01-08 -0.132026 -1.788441  1.819314  1.229088
1 2008-01-07  2.308084  0.297705  0.866629  2.102239
2 2008-01-06  0.434028  0.422461  0.131617  0.529688
3 2008-01-05  0.581159  0.361130  0.426023 -0.975931
4 2008-01-04  1.321511  0.003311  0.250729  0.797861
5 2008-01-03  1.054389  0.069318  0.724151 -1.223632
6 2008-01-02  0.365919  1.301306  0.565771  0.859550
7 2008-01-01  1.709282  0.575747  0.222841  0.154718
```

EJERCICIO 6: Escribir un programa que genere y muestre por pantalla una serie con los datos de las ventas indexadas por los años, antes y después de aplicarles un descuento del 10%.

```
In [39]: df = pd.DataFrame({'Ventas': [13434, 20069, 3000, 345555, 23254, 2345433]}, index=[2009, 2001, 2002, 2003, 2004, 2005])

inicio = df['Ventas'].iloc[0]
for i in df.index:
    df['2º Descuento'] = df['Ventas'] * 0.9
    df['1º Descuento'] = df['2º Descuento'] * 0.9
    df['3º Descuento'] = df['1º Descuento'] * 0.9
    df['4º Descuento'] = df['3º Descuento'] * 0.9
    df['5º Descuento'] = df['4º Descuento'] * 0.9
    df['6º Descuento'] = df['5º Descuento'] * 0.9
    df['7º Descuento'] = df['6º Descuento'] * 0.9
    df['8º Descuento'] = df['7º Descuento'] * 0.9
    df['9º Descuento'] = df['8º Descuento'] * 0.9
    df['10º Descuento'] = df['9º Descuento'] * 0.9
    df['11º Descuento'] = df['10º Descuento'] * 0.9
    df['12º Descuento'] = df['11º Descuento'] * 0.9
    df['13º Descuento'] = df['12º Descuento'] * 0.9
    df['14º Descuento'] = df['13º Descuento'] * 0.9
    df['15º Descuento'] = df['14º Descuento'] * 0.9
    df['16º Descuento'] = df['15º Descuento'] * 0.9
    df['17º Descuento'] = df['16º Descuento'] * 0.9
    df['18º Descuento'] = df['17º Descuento'] * 0.9
    df['19º Descuento'] = df['18º Descuento'] * 0.9
    df['20º Descuento'] = df['19º Descuento'] * 0.9
    df['21º Descuento'] = df['20º Descuento'] * 0.9
    df['22º Descuento'] = df['21º Descuento'] * 0.9
    df['23º Descuento'] = df['22º Descuento'] * 0.9
    df['24º Descuento'] = df['23º Descuento'] * 0.9
    df['25º Descuento'] = df['24º Descuento'] * 0.9
    df['26º Descuento'] = df['25º Descuento'] * 0.9
    df['27º Descuento'] = df['26º Descuento'] * 0.9
    df['28º Descuento'] = df['27º Descuento'] * 0.9
    df['29º Descuento'] = df['28º Descuento'] * 0.9
    df['30º Descuento'] = df['29º Descuento'] * 0.9
    df['31º Desc
```