

# **Fundamentos de tratamiento de datos con el stack científico de Python**

**© Ediciones Roble S. L.**

# Indice

<b>Fundamentos de tratamiento de datos con el stack científico de Python</b>	<b>3</b>
I. Introducción	3
II. Objetivos	3
III. Gestión de matrices y cálculo estadístico con NumPy	4
IV. Representación gráfica con Matplotlib	25
V. Manipulación y análisis de datos con Pandas	48
5.1. Series	49
5.2. Dataframes	53
5.2.1. Datos de un dataframe	53
5.2.2. Operaciones sobre dataframes	57
VI. Resumen	67
VII. Apéndice	68
VIII. Caso práctico	69
Se pide	69
Solución	69
<b>Recursos</b>	<b>71</b>
Bibliografía	71
Glosario	71

# Fundamentos de tratamiento de datos con el stack científico de Python

## I. Introducción

Actualmente, en el ámbito del análisis de datos, Python se ha convertido en uno de los lenguajes más utilizados. Python es un lenguaje de propósito general que no está diseñado de forma específica para realizar análisis de datos. Sin embargo, se ha definido un conjunto de librerías científicas para este fin que le dotan de toda la funcionalidad necesaria para realizar este tipo de tareas.

En esta unidad se van a estudiar las librerías científicas de Python. En primer lugar, se examinará la librería NumPy, que permite la gestión y manipulación de arrays de datos, así como el cálculo de operaciones estadísticas. La principal novedad de NumPy es la posibilidad de manejar arrays de datos como si fueran tipos de datos básicos.

También se estudiará la librería Matplotlib. Esta permite la representación gráfica de datos usando diferentes formatos de representación y tipos de gráficas. De una manera sencilla se pueden configurar las diferentes gráficas y datos que se representan.

Por último, se describirá la librería Pandas, la cual introduce un conjunto de estructuras de datos nuevas que permite realizar operaciones de manipulación de datos de una forma muy sencilla. Estas operaciones —como selección, limpieza o transformación de datos— suelen ser bastante útiles cuando se están preparando los datos para ser analizados.



En esta unidad, se incluye cada ejercicio en un cuaderno de Jupyter Notebook (.ipynb). Con el objetivo de familiarizarse con los distintos métodos y librerías, se recomienda descargarlos y ejecutarlos. Todos estos cuadernos están en el archivo comprimido que se puede descargar en el siguiente enlace:

[Ejercicios\\_ud\\_stack\\_python](#)

## II. Objetivos

Los objetivos que los alumnos alcanzarán tras el estudio de esta unidad son:



- Conocer la librería NumPy de Python y saber utilizarla para llevar a cabo la gestión de arrays y el cálculo estadístico.
- Conocer la librería Matplotlib de Python y saber utilizarla para realizar representaciones gráficas de datos.
- Conocer la librería Pandas de Python y saber utilizarla para manipular datos y analizarlos.
- Saber utilizar de manera conjunta las librerías científicas para poder realizar un análisis de datos completo.

### III. Gestión de matrices y cálculo estadístico con NumPy

El módulo NumPy (Numerical Python) es una extensión de Python que proporciona funciones y rutinas matemáticas para la manipulación de arrays y matrices de datos numéricos de una forma eficiente.

El elemento esencial de NumPy son unos objetos denominados ndarray, arrays multidimensionales donde todos sus elementos son del mismo tipo y están indexados por una tupla de números positivos.

Su principal ventaja es la eficiencia para manipular vectores y matrices. En este sentido, proporciona funciones que operan sobre ndarrays.

#### Atributos de los ndarray

Cada ndarray tiene un conjunto de atributos que le caracterizan:

##### **ndarray.ndim**

Número de dimensiones del array.

##### **ndarray.dtype**

Describe el tipo de elementos del array.

##### **ndarray.shape**

Dimensiones del array. Se trata de una tupla de enteros que indica el tamaño del array en cada dimensión. Por ejemplo, para una matriz de n filas y m columnas, sus dimensiones serían (n,m) y el número de dimensiones sería 2.

##### **ndarray.size**

Número total de elementos del array (producto de los elementos de la dimensión).

##### **ndarray.itemsize**

El tamaño en bytes de cada elemento del array.

##### **ndarray.data**

Es el buffer conteniendo los elementos actuales del array. Normalmente no se usa este atributo, pues se accede a los elementos directamente mediante los índices.

#### Importación de NumPy

Para usar NumPy, hay que importarlo. Normalmente se importa con un alias:

```
import numpy as np
```

#### Creación de un array

Existen varias formas para crear un array:

**La función array()**

La forma más sencilla de crear un array es utilizando la función **array** y una secuencia de valores (Figura 7.1.) .), cuaderno *Ejemplo\_1.ipynb*. En este caso el tipo del array resultante se deduce del tipo de elementos de las secuencias.

```
import numpy as np
a = np.array( [2,3,4] )
print (a)
```

```
[2 3 4]
```

```
a.dtype
```

```
dtype('int32')
```

**Figura 7.1.** Creación de un array usando la función array.

Cuando a la función array se le pasa como argumento una secuencia de secuencias, genera un array de tantas dimensiones como secuencias (Figura 7.2.), cuaderno *Ejemplo\_2.ipynb*.

```
import numpy as np
a = np.array( [[2,3,4], [5,6,7]] )
a
```

```
array([[2, 3, 4],
       [5, 6, 7]])
```

**Figura 7.2.** Creación de un array con una secuencia de secuencias.

El tipo del array puede ser especificado en el momento de la creación del mismo (Figura 7.3.), cuaderno *Ejemplo\_3.ipynb*.

```
import numpy as np
c = np.array( [[2,3], [5,6]], dtype=complex )
c
```

```
array([[2.+0.j, 3.+0.j],
       [5.+0.j, 6.+0.j]])
```

**Figura 7.3.** Especificación del array en el momento de su creación.

**La función `arange()`**

Creación de un array mediante la función **`arange()`**, que genera una secuencia de números (Figura 7.4.), cuaderno *Ejemplo\_4.ipynb*. Toma como parámetros el rango de los números a generar, y la distancia entre ellos, y genera un array unidimensional.

```
import numpy as np
a = np.arange(1,10,2)
a
array([1, 3, 5, 7, 9])
```

**Figura 7.4.** Creación de un array con el método `arange`.

**Redimensionar el array**

Se puede redimensionar el array que genera **`arange()`** mediante el método **`reshape()`**, que toma como argumentos la dimensiones (Figura 7.5.), cuaderno *Ejemplo\_5.ipynb*. Las dimensiones deben ser consistentes con el número de elementos generados.

```
import numpy as np
a = np.arange(2, 6, 0.3).reshape(2, 7)
a
array([[2. , 2.3, 2.6, 2.9, 3.2, 3.5, 3.8],
       [4.1, 4.4, 4.7, 5. , 5.3, 5.6, 5.9]])
```

**Figura 7.5.** Redimensión de un array.

**Uso de argumentos reales con `arange()`**

Cuando se usan argumentos reales con **`arange()`**, es mejor usar la función **`linspace()`** (Figura 7.6.), cuaderno *Ejemplo\_6.ipynb* que también genera una secuencia de números a partir de un rango dado para crear un array, pero recibe como tercer argumento el número de elementos, en vez de la distancia entre ellos.

```
import numpy as np
a = np.linspace(1,10,8)
a
array([ 1.         ,  2.28571429,  3.57142857,  4.85714286,  6.14285714,
        7.42857143,  8.71428571, 10.         ])
```

**Figura 7.6.** Uso de la función `linspace`.

**Array aleatorio con rand()**

Creación de un array unidimensional con datos aleatorios mediante la función **rand** del módulo **Random** (Figura 7.7.), cuaderno *Ejemplo\_7.ipynb*. La función **rand** devuelve un número aleatorio procedente de una distribución uniforme en el intervalo [0,1].

```
import numpy as np
a = np.random.rand(10)
a
array([0.03136756, 0.56639071, 0.86969584, 0.58186369, 0.23643603,
       0.82588378, 0.52276312, 0.27644873, 0.31443282, 0.62791627])
```

**Figura 7.7.** Creación de un array con datos aleatorios.

También es posible generar arrays multidimensionales si se proporcionan las dimensiones como argumentos (Figura 7.8.), cuaderno *Ejemplo\_8.ipynb*.

```
import numpy as np
a = np.random.rand(3,4) #valores aleatorios 3 filas y 4 columnas
a
array([[0.68768809, 0.1936327 , 0.32025235, 0.26777976],
       [0.12570141, 0.28699652, 0.90407167, 0.54861304],
       [0.52002797, 0.5257531 , 0.18119511, 0.07235105]])
```

**Figura 7.8.** Creación de un array multidimensional.

### Funciones especiales para la creación de arrays

Existen funciones que generan arrays especiales que solo requieren como argumento el tamaño del array (Figura 7.9.), cuaderno *Ejemplo\_9.ipynb*:

- **zeros**: crea un array lleno de ceros.
- **ones**: crea un array lleno de unos.
- **empty**: crea un array cuyo contenido es aleatorio.

```
import numpy as np
a = np.zeros(4)
a
array([0., 0., 0., 0.])

b = np.ones ([1 ,2])
b
array([[1., 1.]])

c = np.empty((2,3))
c
array([[0., 0., 0.],
       [0., 0., 0.]])
```

**Figura 7.9.** Funciones especiales.

### Operaciones sobre arrays

En NumPy se puede operar aritméticamente sobre los arrays como si se tratara de escalares. Se caracteriza por:



## Operaciones sobre la misma posición

Las operaciones actúan sobre los elementos de la misma posición y, como resultado, crean un nuevo array (Figura 7.10.), cuaderno *Ejemplo\_10.ipynb*:

```
import numpy as np
a = np.array([34,12,3,4])
a
```

```
array([34, 12,  3,  4])
```

```
b = np.array([2,3,5,6])
b
```

```
array([2, 3, 5, 6])
```

```
c=a-b
c
```

```
array([32,  9, -2, -2])
```

```
d=a*10
d
```

```
array([340, 120,  30,  40])
```

**Figura 7.10.** Operaciones aritméticas entre dos arrays.

## Arrays de diferente tipo

Cuando se opera con arrays de diferente tipo, el tipo del array resultante corresponde al más general (Figura 7.11.), cuaderno *Ejemplo\_11.ipynb*:

```
import numpy as np
a = np.array([2,3,6,7])
a
```

```
array([2, 3, 6, 7])
```

```
b=np.linspace(2,3,4)
b
```

```
array([2.          , 2.33333333, 2.66666667, 3.          ])
```

```
c=a+b
c
```

```
array([ 4.          ,  5.33333333,  8.66666667, 10.          ])
```

```
c.dtype
```

```
dtype('float64')
```

**Figura 7.11.** Operaciones entre arrays de diferente tipo.

## Funciones matemáticas universales

En NumPy existe un conjunto de funciones matemáticas denominadas **funciones universales**: sin, cos, exp... (Figura 7.12.), cuaderno *Ejemplo\_12.ipynb*. Estas funciones operan elemento a elemento y generan como resultado un nuevo array.

```
import numpy as np
a = np.array([-7, 60, -9])
a

array([-7, 60, -9])

b=np.square(a)
b

array([ 49, 3600, 81], dtype=int32)

c=np.sqrt(np.abs(a))
c

array([2.64575131, 7.74596669, 3.          ])
```

**Figura 7.12.** Funciones universales sobre arrays.

En particular es muy útil la función **dot()**, que permite realizar la multiplicación matricial (Figura

7.13.), cuaderno *Ejemplo\_13.ipynb*:

```
import numpy as np
a = [[1, 0], [0, 1]]
b = [[4, 1], [2, 2]]

c=np.dot(a,b)
c

array([[4, 1],
       [2, 2]])
```

**Figura 7.13.** Multiplicación matricial.

### El argumento axis

Existen algunas operaciones implementadas como métodos de la clase `ndarray` y, por defecto, se aplican a todos los elementos del array. Sin embargo, es posible especificar la dimensión sobre la que se quiere aplicar la operación mediante el argumento **axis** (Figura 7.14.), cuaderno *Ejemplo\_14.ipynb* que toma el valor 0 (columnas) o 1 (filas).

```
import numpy as np
a = np.random.rand(3,4)
a

array([[0.58584044, 0.62487763, 0.18664195, 0.18512665],
       [0.22869104, 0.80880337, 0.76676978, 0.84193634],
       [0.22640765, 0.48739203, 0.84777263, 0.8589042 ]])

a.sum()

6.6491636921061925

a.sum(axis=0)

array([1.04093912, 1.92107303, 1.80118436, 1.88596719])
```

Figura 7.14. Uso del argumento axis.

### Operadores += y \*=

Los operadores `+=` y `*=` modifican los arrays en vez de crear uno nuevo (Figura 7.15.), cuaderno *Ejemplo\_15.ipynb*.

```
import numpy as np
a = np.array([2,3,5,6])
a

array([2, 3, 5, 6])

a+=3
a

array([5, 6, 8, 9])

a*=a
a

array([25, 36, 64, 81])
```

Figura 7.15. Modificación de arrays.

### Acceso a arrays

El acceso a los arrays se realiza de forma indexada, de manera que se pueden seleccionar subrangos y se puede iterar sobre sus elementos. Téngase en cuenta lo siguiente:

- Cuando se accede a un array por índice, se indica un número por cada dimensión.
- Cuando se elige un subrango se indica por cada dimensión a:b:c, donde a indica dónde se comienza, b dónde se termina y c el salto entre elementos.

### Acceso en arrays unidimensionales

Cuando trabajamos con arrays de una dimensión, el acceso a los elementos se realiza de forma similar al de las listas o tuplas de elementos (Figura 7.16.), cuaderno *Ejemplo\_16.ipynb*.

```
import numpy as np
arr = np.arange(2,20)
arr

array([ 2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16, 17, 18,
        19])

arr[0]

2

arr[1::3]

array([ 3,  6,  9, 12, 15, 18])
```

**Figura 7.16.** Acceso indexado a un array unidimensional

### Diferencia con las listas

Una diferencia importante con las listas es que las particiones de un ndarray mediante la notación [inicio:fin:paso] son vistas del array original. Todos los cambios realizados en las vistas se reflejan en el array original (Figura 7.17.), cuaderno *Ejemplo\_17.ipynb*:

```
import numpy as np
arr = np.arange(2,20)
arr

array([ 2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16, 17, 18,
        19])

b = arr[-2:]
b[:] = 0
arr

array([ 2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16, 17,  0,
        0])
```

**Figura 7.17.** Modificaciones de un array

### Acceso a elementos en array multidimensional

El acceso a los elementos de un array bidimensional se realiza indicando los índices separados por una coma (Figura 7.18.), cuaderno *Ejemplo\_18.ipynb*.

```
import numpy as np
b = np.array([
    [ 0, 1, 2, 3],
    [10,11,12,13],
    [20,21,22,23],
    [30,31,32,33],
    [40,41,42,43]
])
b
```

```
array([[ 0,  1,  2,  3],
       [10, 11, 12, 13],
       [20, 21, 22, 23],
       [30, 31, 32, 33],
       [40, 41, 42, 43]])
```

```
b[2,3] #Elemento de la fila 3 y columna 4
```

```
23
```

```
b[0:5,1] #Elementos de todas las filas de la columna 2
```

```
array([ 1, 11, 21, 31, 41])
```

**Figura 7.18.** Acceso de un array de varias dimensiones.

Cuando se accede a un array y se indican menos índices que el número de dimensiones, se consideran todos los valores de las dimensiones no indicadas (Figura 7.19.), cuaderno *Ejemplo\_19.ipynb*.

```
b[0] #Toda la fila 1
```

```
array([0, 1, 2, 3])
```

**Figura 7.19.** Acceso a un array multidimensional indicando menos dimensiones

### Acceso mediante máscaras

Otra forma de acceso a partes de un array de NumPy es mediante un array de booleanos denominado máscara (Figura 7.20.), cuaderno *Ejemplo\_20.ipynb*.

```
import numpy as np
n = np.arange(10000)
n

array([ 0,  1,  2, ..., 9997, 9998, 9999])

( n > 20 ) & ( n < 44 ) #Comprobar si se cumple la condición en los elementos del array
array([False, False, False, ..., False, False, False])

mascara = ( n > 20 ) & ( n < 44 ) #Se crea la máscara con valor booleano
mascara

array([False, False, False, ..., False, False, False])

n[mascara] #Se usa máscara para recuperar los valores que cumplen con la condición
array([21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37,
       38, 39, 40, 41, 42, 43])
```

**Figura 7.20.** Acceso a un array mediante una máscara

## Iteración sobre arrays

Hay varias formas de iterar sobre los elementos de un array:

### Usando un bucle for

Usando un **for** (Figura 7.21.), cuaderno *Ejemplo\_21.ipynb*.

```
import numpy as np
a = np.arange(10)
a
```

array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])

```
for i in a:
    print (i)
```

```
0
1
2
3
4
5
6
7
8
9
```

Figura 7.21. Recorrido de un array usando un for.

### Usando el iterador flat

Usando el iterador **flat**, que permite iterar sobre todos los elementos del array: aplanar el array y lo convierte en un array unidimensional (Figura 7.22.), cuaderno *Ejemplo\_22.ipynb*.

```
import numpy as np
a = np.arange(10)
a
```

array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])

```
for i in a.flat:
    print (i)
```

```
0
1
2
3
4
5
6
7
8
9
```

Figura 7.22. Recorrido de un array usando el iterador flat.

## Manipulaciones sobre un array

Sobre un array se pueden realizar diferentes tipos de manipulaciones:

### Cambio de tamaño

Para ello se puede aplanar la matriz mediante la función **ravel()** y, a continuación, usar el método **shape()** o **resize()** para establecer las nuevas dimensiones (Figura 7.23.), cuaderno *Ejemplo\_23.ipynb*.

```
import numpy as np
a = np.arange(12)
a
array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11])

a.reshape(3,4)
array([[ 0,  1,  2,  3],
       [ 4,  5,  6,  7],
       [ 8,  9, 10, 11]])

a.ravel()
array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11])

a.shape=(4,3)
a
array([[ 0,  1,  2],
       [ 3,  4,  5],
       [ 6,  7,  8],
       [ 9, 10, 11]])
```

**Figura 7.23.** Cambio de tamaño.



## Fusión de dos arrays

Fusionar dos arrays (Figura 7.24.), cuaderno *Ejemplo\_24.ipynb* mediante la función **concatenate()** verticalmente (atributo `axis=0`) o bien horizontalmente (atributo `axis=1`).

```
import numpy as np
a = np.arange(4).reshape(2,2)
a
```

```
array([[0, 1],
       [2, 3]])
```

```
b = np.arange(5,9,1).reshape(2,2)
b
```

```
array([[5, 6],
       [7, 8]])
```

```
np.concatenate((a,b),axis=0) #Fusión vertical
```

```
array([[0, 1],
       [2, 3],
       [5, 6],
       [7, 8]])
```

```
np.concatenate((a,b),axis=1) #Fusión horizontal
```

```
array([[0, 1, 5, 6],
       [2, 3, 7, 8]])
```

**Figura 7.24.** Fusión de dos arrays.

## División de un array

Dividir un array en partes iguales usando las funciones **hsplit()** y **vsplit()**, indicando la columna o fila por donde se divide dependiendo la función (Figura 7.25.), cuaderno *Ejemplo\_25.ipynb*.

```
import numpy as np
a = np.arange(12).reshape(2,6)
a

array([[ 0,  1,  2,  3,  4,  5],
       [ 6,  7,  8,  9, 10, 11]])
```

Figura 7.25. División de un array.

```
np.hsplit(a,3) #División por la columna 2

[array([[0, 1],
       [6, 7]]), array([[2, 3],
       [8, 9]]), array([[4, 5],
       [10, 11]])]
```

Alternativamente se podría descomponer el array, lo cual se hace por filas (Figura 7.26.), cuaderno

```
import numpy as np
a = np.arange(10).reshape(2,5)
a

array([[0, 1, 2, 3, 4],
       [5, 6, 7, 8, 9]])
```

*Ejemplo\_26.ipynb*.

```
b,c=a #División horizontal
b
```

```
array([0, 1, 2, 3, 4])
```

```
c
```

```
array([5, 6, 7, 8, 9])
```

Figura 7.26. División de un array por filas,

## Asignaciones y copias

Las asignaciones y las llamadas a funciones no hacen copia del array o de sus datos (Figura 7.27.), cuaderno *Ejemplo\_27.ipynb*.

```
import numpy as np
a = np.array( [ [1,2],[3,4] ] )
```

```
b=a
b.shape=(1,4)
a
```

Figura 7.27. Ejemplo de asignación.

```
array([[1, 2, 3, 4]])
```

```
b
```

```
array([[1, 2, 3, 4]])
```

**El método view()**

Permite crear un nuevo array que toma los datos del array original, pero sin tratarse del mismo. Sin embargo, el array original se ve afectado por las operaciones que se hagan sobre la vista (Figura 7.28.), cuaderno *Ejemplo\_28.ipynb*.

```
import numpy as np
a = np.array( [ [1,2],[3,4] ] )
```

```
c = a.view()
c
array([[1, 2],
       [3, 4]])
```

```
c.shape=(1,4)
c
array([[1, 2, 3, 4]])
```

```
a
array([[1, 2],
       [3, 4]])
```

**Figura 7.28.** Uso del método view().

**El método copy()**

Permite realizar una copia independiente del array original (Figura 7.29.), cuaderno *Ejemplo\_29.ipynb*:

```
import numpy as np
a = np.array( [ [1,2],[3,4] ] )
```

```
d=a.copy()
d
array([[1, 2],
       [3, 4]])
```

```
d[1,1]=999
d
array([[ 1,  2],
       [ 3, 999]])
```

```
a
array([[1, 2],
       [3, 4]])
```

**Figura 7.29.** Copia de un array.

**El método sort()**

Es posible ordenar los arrays usando el método **sort()**. En el caso de arrays multidimensionales, hay que indicar la dimensión sobre la que se ordena (Figura 7.30.), cuaderno *Ejemplo\_30.ipynb*.

```
import numpy as np
a = np.array([9,3,5,2,5,1])
a
```

```
array([9, 3, 5, 2, 5, 1])
```

```
a.sort()
a
```

```
array([1, 2, 3, 5, 5, 9])
```

```
b = np.array([[4,1],[2,3]])
b
```

```
array([[4, 1],
       [2, 3]])
```

```
b.sort(0) #Ordenación por columnas
b
```

```
array([[2, 1],
       [4, 3]])
```

**Figura 7.30.** Ordenación de arrays.

**Los métodos any() y all()**

Se puede operar sobre arrays de booleanos mediante los métodos **any** and **all**. Permiten chequear si algún valor es cierto o si todos los valores son ciertos respectivamente (figura 3.31.), cuaderno *Ejemplo\_31.ipynb*.

```
import numpy as np
b = np.array([False,True,False,False])
b.any()
```

```
True
```

**Figura 7.31.** Operación sobre arrays.

**Operaciones estadísticas**

Por último, es apropiado comentar que el módulo NumPy proporciona métodos que permiten realizar otras operaciones matemáticas, tales como obtener el mínimo elemento de un array, el máximo, álgebra lineal, operaciones sobre conjuntos...

Se van a mostrar algunas de las operaciones estadísticas que facilita:

### Suma, mínimo, máximo

Suma, mínimo, máximo (Figura 7.32.), cuaderno *Ejemplo\_32.ipynb*.

```
import numpy as np
a = np.arange(8)
a
array([0, 1, 2, 3, 4, 5, 6, 7])
```

```
np.sum(a) #Suma todos los elementos
```

```
28
```

**Figura 7.32.** Suma, máximo y mínimo.

```
print(np.min(a), np.argmin(a)) #Valor mínimo y su posición
```

```
0 0
```

```
print(np.max(a), np.argmax(a)) #Valor máximo y su posición
```

```
7 7
```

### Media

Es la suma de todos los elementos dividida por el número de elementos. En Python puede calcularse usando la función `numpy.mean` (Figura 7.33.), cuaderno *Ejemplo\_33.ipynb*.

```
import numpy as np
a = np.arange(8)
a
array([0, 1, 2, 3, 4, 5, 6, 7])
```

```
np.mean(a) #Media del array, considerando todos su valores
```

```
3.5
```

```
np.mean(a, axis=0) #Media por columnas
```

```
3.5
```

**Figura 7.33.** Media.

## Varianza

Es una medida de dispersión de una variable aleatoria definida como la esperanza del cuadrado de dicha variable respecto a su media. En Python puede calcularse usando la función `numpy.var` (Figura 7.34.), cuaderno *Ejemplo\_34.ipynb*.

```
import numpy as np
a = np.arange(8)
a

array([0, 1, 2, 3, 4, 5, 6, 7])

np.var(a, axis=0) #Varianza por columnas
5.25
```

Figura 7.34. Varianza.

## Desviación estándar

La desviación estándar es una medida de dispersión de una variable que se define como la raíz cuadrada de la varianza de la variable. En Python se puede calcular usando la función `numpy.std` (Figura 7.35.), cuaderno *Ejemplo\_35.ipynb*:

```
import numpy as np
a = np.array([[1, 2], [3, 4]])
a

array([[1, 2],
       [3, 4]])

np.std(a, axis=0) #Desviación estandar (por columnas)
array([1., 1.])
```

Figura 7.35. Desviación estándar.

## Mediana

Representa el valor de la variable de posición central en un conjunto de datos ordenados. En Python puede calcularse usando la función `numpy.median` (Figura 7.36.), cuaderno *Ejemplo\_36.ipynb*

```
In [1]: 1 import numpy as np
        2 a = np.array([[1, 2], [3, 4]])
        3 a

Out[1]: array([[1, 2],
               [3, 4]])
```

---

```
In [2]: 1 np.median(a, axis=1) #Mediana por filas

Out[2]: array([1.5, 3.5])
```

Figura 7.36. Mediana.

## Correlación

Dadas dos variables aleatorias, este coeficiente indica si están relacionadas o no. En Python puede calcularse usando `numpy.corrcoef`. La función devuelve los coeficientes de correlación (Figura 7.37.), cuaderno *Ejemplo\_37.ipynb*.

```
"""
Se calcula el coeficiente de correlación teniendo
en cuenta que cada fila representa una variable
"""
import numpy as np
a = np.array([[2,3,4,5], [4,5,6,6]])
a

array([[2, 3, 4, 5],
       [4, 5, 6, 6]])

np.corrcoef(a) #correlación de la transpuesta

array([[1.          , 0.94387981],
       [0.94387981, 1.          ]])
```

Figura 7.37. Correlación.

## Covarianza

Determina si existe una dependencia entre dos variables aleatorias. En Python se puede calcular usando `numpy.cov` (Figura 7.38.), cuaderno *Ejemplo\_38.ipynb*:

```
import numpy as np
a = np.array([[2,3,4,5], [4,5,6,6]])
a

array([[2, 3, 4, 5],
       [4, 5, 6, 6]])

np.cov(a[:,0]) #covarianza de la primera columna
array(2.)
```

Figura 7.38. Covarianza.

## Métodos aleatorios

Métodos aleatorios (Figura 7.39.), cuaderno *Ejemplo\_39.ipynb*.

```
import numpy.random as r
s = r.rand(10) #Genera los números aleatorios de una normal (0,1)
s

array([0.41440616, 0.72875733, 0.69914009, 0.38068869, 0.78980851,
       0.48717792, 0.43402754, 0.9362674 , 0.42987006, 0.82364173])

t=r.normal(size=(5,1)) #Genera un array con números pertenecientes a una normal
t

array([[ -0.87920083],
       [ 0.29948192],
       [ 0.27530868],
       [-0.42021168],
       [ 0.07699345]])

k=r.normal(4,5) #Genera un número perteneciente a una normal (4,5)
k

-1.554950120683019
```

Figura 7.39. Métodos aleatorios







Pincha [aquí](#) para descargar los archivos necesarios para seguir el tutorial.

## IV. Representación gráfica con Matplotlib



Matplotlib es una librería de Python para realizar gráficos. Se caracteriza porque es fácil de usar, flexible y se puede configurar de múltiples maneras.

### Importación del módulo

Para importar Matplotlib desde cualquier programa de Python, se hará de la siguiente manera: **import matplotlib.pyplot as plt**.

Y para usar un comando de Matplotlib se usará el estilo siguiente: **plt.comando()**.

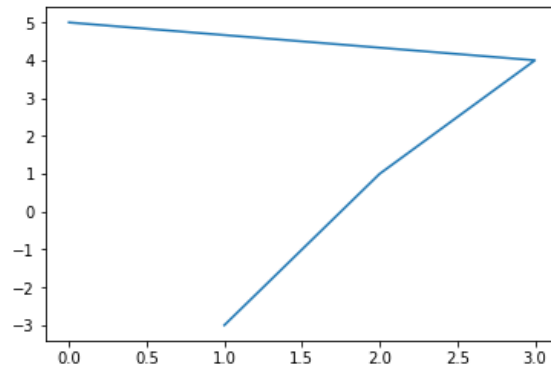
### El comando plot()

El principal comando de Matplotlib es **plot()**, el cual permite representar gráficamente una lista de pares de valores (x, y) sobre un eje de coordenadas uniendo cada punto de acuerdo al orden en que aparecen.

En el siguiente ejemplo, se representa la lista de pares de valores: (1,-3), (2,1), (3,4), (0,5) (Figura 7.40.), cuaderno *Ejemplo\_40.ipynb*. El comando `%matplotlib inline` permite ver un gráfico dentro de un *notebook*.



```
%matplotlib inline
import matplotlib.pyplot as plt
plt.plot([1,2,3,0],[-3,1,4,5])
plt.show()
```



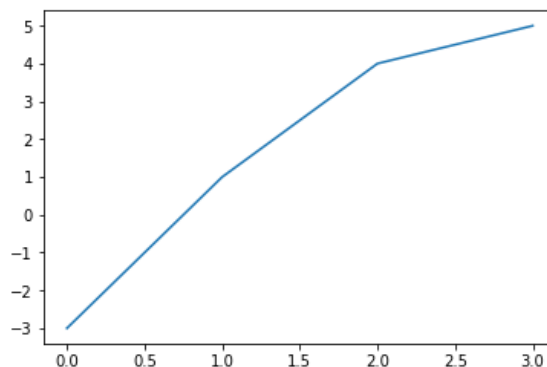
**Figura 7.40.** Uso de plot().

Se puede omitir la lista de valores para el eje x, y. En este caso se usa como valores, por defecto, la lista de valores que van desde el 0 (primer valor) hasta el n-1, donde n es el número de valores proporcionados para el eje y.

En el siguiente ejemplo se proporcionan únicamente valores para el eje y, por lo que se representa la secuencia de pares (0,-3), (1,1), (2,4), (3,5) (Figura 7.41.),cuaderno *Ejemplo\_41.ipynb*.



```
%matplotlib inline
import matplotlib.pyplot as plt
plt.plot([-3,1,4,5])
plt.show()
```



**Figura 7.41.** Representación omitiendo valores.

En los anteriores ejemplos, se ha usado el comando **show()**, el cual sirve para abrir la ventana que contiene la imagen generada.

## Presentar datos de secuencias

### Secuencia única

Para representar datos, puede ser útil usar funciones que generen secuencias tales como **range(i,j,k)** o **arange(i, j, k)** (Figura 7.42.), cuaderno *Ejemplo\_42.ipynb*.

```
%matplotlib inline
import matplotlib.pyplot as plt
x = range(8)
plt.plot(x, [xi**2 for xi in x])
plt.show()
```

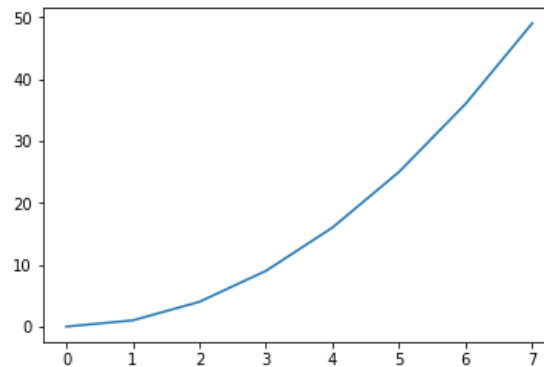
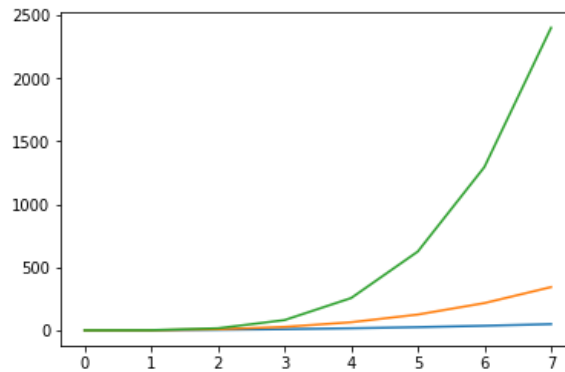


Figura 7.42. Representación usando range().

### Secuencias múltiples

En algunas ocasiones, puede ser interesante mostrar varias representaciones sobre un mismo gráfico. Esto se puede realizar con el comando **plot()**, de formas distintas. La manera más fácil de representarlo es invocar a **plot()** para que dibuje cada una de las funciones y, en último lugar, invocar al comando **show()** (Figura 7.43.), cuaderno *Ejemplo\_43.ipynb*.

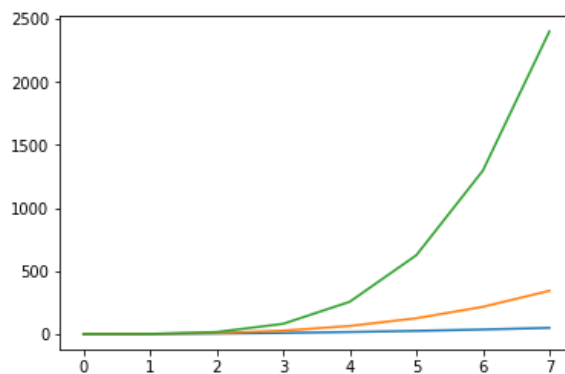
```
%matplotlib inline
import matplotlib.pyplot as plt
x = range(8)
plt.plot(x, [xi**2 for xi in x])
plt.plot(x, [xi**3 for xi in x])
plt.plot(x, [xi**4 for xi in x])
plt.show()
```



**Figura 7.43.** Representación de varias funciones en una misma gráfica.

Otra forma de realizar la misma operación es utilizando NumPy (Figura 7.44.), cuaderno *Ejemplo\_44.ipynb*.

```
%matplotlib inline
import matplotlib.pyplot as plt
x = range(8)
plt.plot(x, [xi**2 for xi in x], x, [xi**3 for xi in x], x, [xi**4 for xi in x])
plt.show()
```



**Figura 7.44.** Representación de varias funciones en una misma gráfica.

Obsérvese que, al superponer diferentes representaciones en la misma gráfica, Matplotlib utiliza automáticamente diferentes colores para cada representación.

### Elementos decorativos

A continuación, se va a mostrar cómo añadir algunos elementos decorativos que suelen aparecer en un gráfico:

## Ejes de coordenadas

Por defecto, se establecen los valores que aparecen en los ejes de coordenadas, de manera que puedan mostrarse en la gráfica los puntos que son dibujados. Para configurar estos valores se usa la función **axis()**.

### axis() sin parámetros

Cuando se ejecuta sin parámetros, devuelve la escala actual utilizada en los ejes en forma de una tupla de 4 valores que indican el límite inferior y superior de la escala de valores usada para el eje x y para el eje y, respectivamente (Figura 7.45.), cuaderno *Ejemplo\_45.ipynb*.

```
%matplotlib inline
import matplotlib.pyplot as plt
import numpy as np
x = np.arange(8.0)
plt.plot(x, [xi**2 for xi in x])
plt.axis()
```

(-0.3500000000000003, 7.35, -2.45, 51.45)

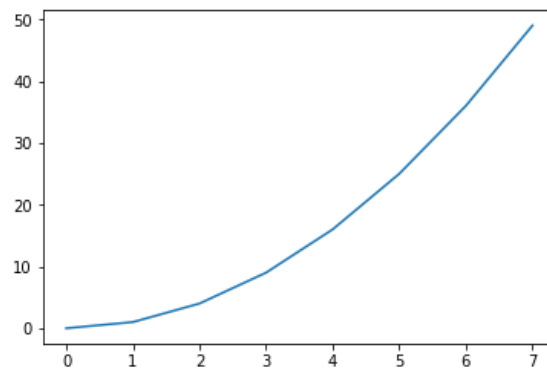


Figura 7.45. Función axis().

### Cambio de escalas

Para cambiar las escalas usadas, se invoca la función **axis()** mediante una lista de 4 valores que representan el valor mínimo del eje x, el valor máximo del eje x, el valor mínimo del eje y, y el valor máximo del eje y (Figura 7.46.), cuaderno *Ejemplo\_46.ipynb*.

```
%matplotlib inline
import matplotlib.pyplot as plt
import numpy as np
x = np.arange(8.0)
plt.plot(x, [xi**2 for xi in x])
plt.axis([1,8,-3,8])
plt.show()
```

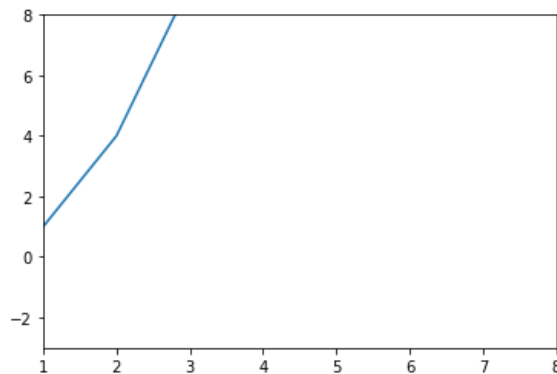


Figura 7.46. Cambio de escala.

### Configuración de un solo eje

La función **axis()** se puede invocar para configurar solo uno de los ejes. Para ello, habrá que proporcionarle los valores correspondientes de la escala utilizando la notación de parámetros argumentales.



Por ejemplo, si se quiere modificar solo el eje y, de manera que varíe entre 4 y 10, se invocaría de la siguiente forma: `plt.axis(ymin=4,ymax=10)`.



Existen dos funciones **xlim()** e **ylim()** que permiten controlar las escalas de los ejes x e y respectivamente, de manera independiente. Se invocan con una lista de 2 valores que representan respectivamente el límite inferior y superior de la escala de un eje.

### Gestión de los valores sobre los ejes

Para gestionar los valores que se colocan sobre los ejes y la localización de ellas, se usan las funciones **plt.xticks()** e **plt.yticks()**, que permiten manipular estos elementos en los ejes x e y respectivamente. En ambos casos, toman como parámetros dos listas de la misma longitud que contienen las localizaciones de los valores en los ejes y, por otra parte, los valores que se desean colocar en esas posiciones.

Esta última lista se puede omitir, de manera que se tomarían como valores los usados en las posiciones (Figura 7.47.), cuaderno *Ejemplo\_47.ipynb*..

```
1 %matplotlib inline
2 import matplotlib.pyplot as plt
3 import numpy as np
4 x = np.arange(8)
5 plt.plot(x, [xi**2 for xi in x])
6 plt.xticks(range(8), ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h'])
7 plt.yticks(range(0,16,2))
8 plt.show()
```

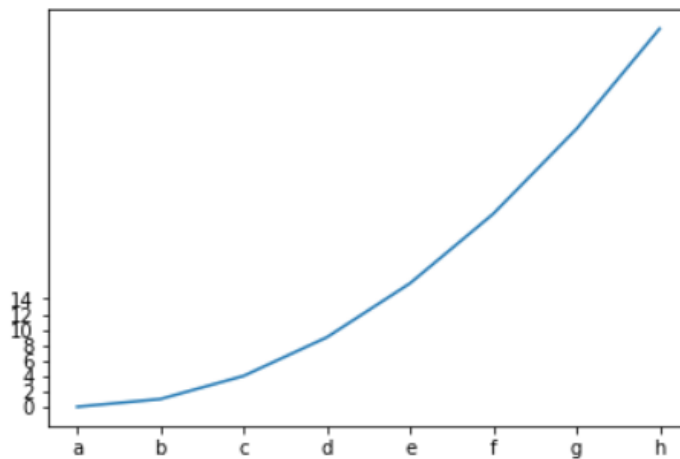


Figura 7.47. Valores de los ejes.

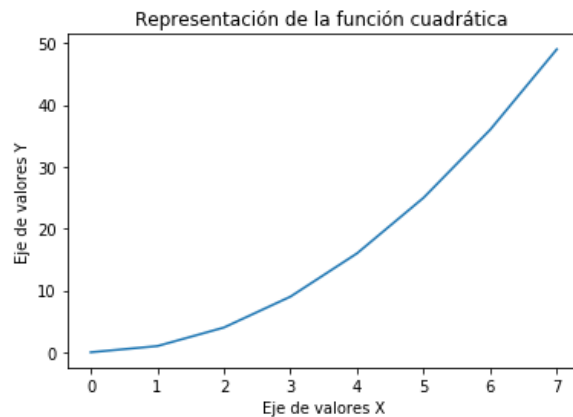
### Añadir etiquetas a los ejes

Para añadir etiquetas a los ejes x e y, se usan las funciones **plt.xlabel()** y **plt.ylabel()** respectivamente, que toman como parámetros los valores de las etiquetas que se quieran añadir a los ejes.

**Título**

Para añadir un título se utiliza la función **plt.title()**, que toma como parámetro el valor del título (Figura 7.48.), cuaderno *Ejemplo\_48.ipynb*.

```
%matplotlib inline
import matplotlib.pyplot as plt
x = range(8)
plt.plot(x, [xi**2 for xi in x])
plt.title("Representación de la función cuadrática ")
plt.xlabel("Eje de valores X")
plt.ylabel("Eje de valores Y")
plt.show()
```



**Figura 7.48.** Gráfica con título y valores en los ejes.

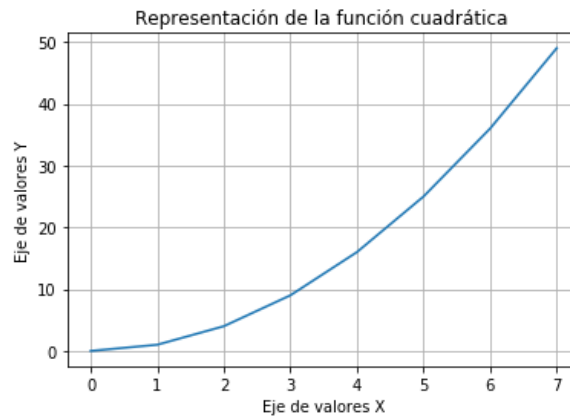


## Cuadrícula

Para añadir una cuadrícula a una gráfica basta con usar la función **plt.grid()**, la cual toma como parámetro el valor True (habilitar la cuadrícula) o False (deshabilitar la cuadrícula).

En la Figura 7.49., aparece representada, cuaderno *Ejemplo\_49.ipynb*..

```
%matplotlib inline
import matplotlib.pyplot as plt
x = range(8)
plt.plot(x, [xi**2 for xi in x])
plt.grid(True)
plt.title("Representación de la función cuadrática ")
plt.xlabel("Eje de valores X")
plt.ylabel("Eje de valores Y")
plt.show()
```



**Figura 7.49.** Gráfica con cuadrícula.

## Leyenda

Para añadir una leyenda a la representación de un conjunto de puntos, se invoca la función **plt.plot()**, que dibuja la secuencia con el argumento **label**, al cual se le asigna el valor que debe tomar la leyenda —esta asignación debe aparecer a continuación de los parámetros que aparecen sin formato de argumentos—.

Por último, se invoca la función **plt.legend()** sin argumentos para que dibuje las leyendas (Figura 7.50.), cuaderno *Ejemplo\_50.ipynb*.

```
%matplotlib inline
import matplotlib.pyplot as plt
x = range(8)
plt.plot(x, [xi for xi in x], label="Lineal")
plt.plot(x, [xi**2 for xi in x], label="Cuadrática")
plt.plot(x, [xi**3 for xi in x], label="Cúbica")
plt.grid(True)
plt.title("Representación de la función cuadrática ")
plt.xlabel("Eje de valores X")
plt.ylabel("Eje de valores Y")
plt.show()
```

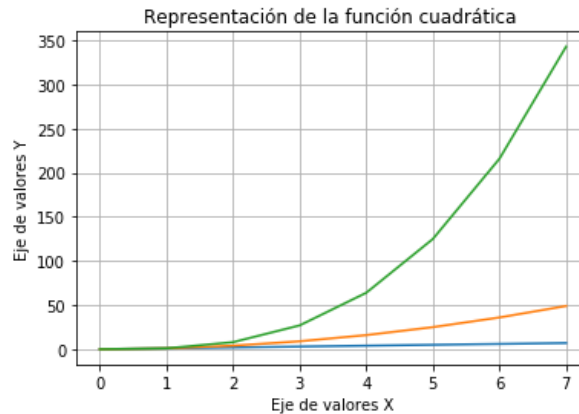


Figura 7.50. Leyenda.

## Formato de puntos

Para configurar el formato para cada par de puntos x,y, se añade un tercer argumento adicional a la función **plt.plot()**. El formato es una cadena que se obtiene concatenando conjuntos de caracteres que representan el color, el estilo y la forma, respectivamente (tabla 6.1.).



Cadena	Código
blue	b
cyan	c
green	g
black	k
magenta	m
red	r
white	w
yellow	y

Código	Significado
-	Línea sólida
--	Línea discontinua
-.	Línea discontinua con puntos
:	Línea de puntos

Código	Significado
.	Punto
,	Píxel
0	Círculo
V	Triángulo hacia abajo
^	Triángulo hacia arriba
<	Triángulo hacia la izquierda
>	Triángulo hacia la derecha
1	Trípode hacia abajo
2	Trípode hacia abajo
3	Trípode hacia la izquierda
4	Trípode hacia la derecha
s	Cuadrado
p	Pentágono
*	Estrella
h	Hexágono
H	Hexágono rotado
+	Signo +
x	Cruz
D	Diamante
d	Diamante delgado
	Línea vertical
-	Línea horizontal

**Tabla 6.1.** Formato de los puntos representados.

En el siguiente ejemplo (Figura 7.51.), cuaderno *Ejemplo\_51.ipynb*, se muestra un formateado de puntos.

```
%matplotlib inline
import matplotlib.pyplot as plt
x = range(8)
plt.plot(x, [xi for xi in x], 'b-d', label="Lineal")
plt.plot(x, [xi**2 for xi in x], 'y--p', label="Cuadrática")
plt.plot(x, [xi**3 for xi in x], 'r:s', label="Cúbica")
plt.legend()
plt.xlabel("Eje de valores X")
plt.ylabel("Eje de valores Y")
plt.show()
```

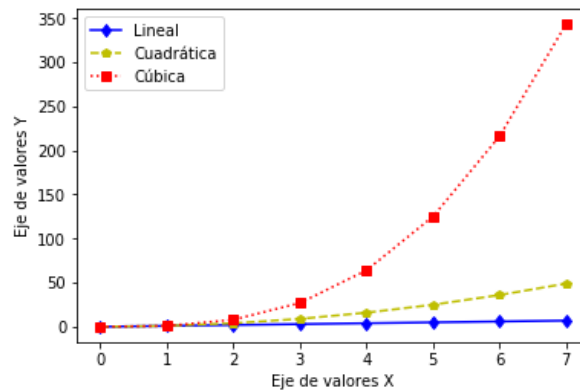


Figura 7.51. Formateado de varios puntos.

## Texto

Para añadir texto, se puede usar la función **plt.text()**, que toma como parámetros un par de puntos que indican la posición donde se quiere escribir y la cadena que se quiere escribir. La posición es relativa a los ejes utilizados para representar los datos (Figura 7.52.), cuaderno *Ejemplo\_52.ipynb*.

```
%matplotlib inline
import matplotlib.pyplot as plt
x = range(8)
plt.plot(x, [xi for xi in x])
plt.text(1,3,"Función lineal")
plt.show()
```

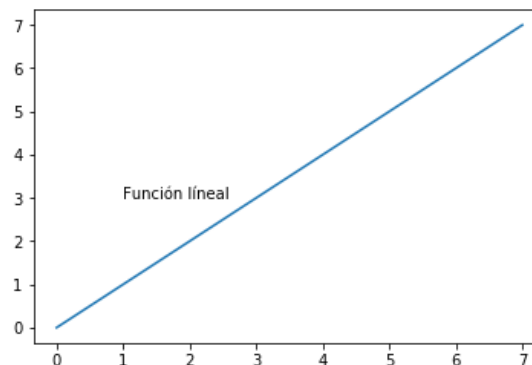
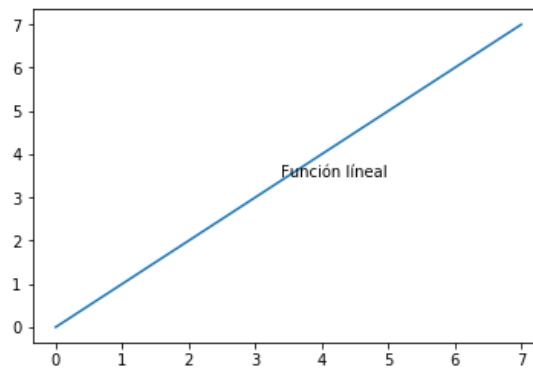


Figura 7.52. Texto sobre la gráfica.

También es posible añadir texto a un dibujo con la función `plt.figtext()`, expresando las coordenadas de forma relativa a la figura dibujada. Las coordenadas varían entre 0 y 1 —la posición (0,0) representa la esquina inferior izquierda y (1,1) la esquina superior derecha—.

Se muestra en la Figura 7.53, cuaderno *Ejemplo\_53.ipynb*:

```
%matplotlib inline
import matplotlib.pyplot as plt
x = range(8)
plt.plot(x, [xi for xi in x])
plt.figtext(0.5,0.5,"Función lineal")
plt.show()
```



**Figura 7.53.** Situando texto sobre la gráfica.

## Anotación

Para añadir una anotación, se usa la función **plt.annotate()**, la cual toma como argumentos (Figura 7.54.), cuaderno *Ejemplo\_54.ipynb*:

- Texto de la anotación.
- Argumento **xy**: posición de la función en la que se quiere añadir la anotación expresada en coordenadas respecto a los ejes.
- Argumento **xytext**: indica la posición de la anotación que se quiere añadir expresada en coordenadas respecto a los ejes.
- Argumento **arrowprops**: propiedades de la flecha que une la anotación con el punto anotado expresado como un diccionario width (ancho de la flecha), headlength (longitud de la flecha reservada a la cabeza de la misma), headwidth (ancho de la base de la flecha), facecolor (color de la superficie de la flecha), shrink (desplazamiento de la flecha con respecto al punto anotado y la anotación expresado como un porcentaje).

```
%matplotlib inline
import matplotlib.pyplot as plt
x = range(8)
plt.plot(x, [xi for xi in x])
plt.grid(True)
plt.annotate("Punto (2,2)", xy=(2,2), xytext=(3,2),
arrowprops=dict(facecolor="green", shrink=0.05, headlength=2))
plt.show()
```

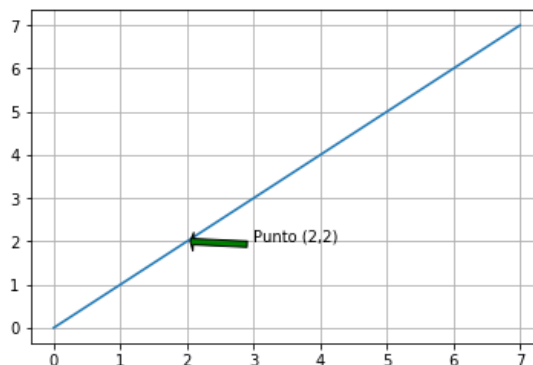


Figura 7.54. Anotación sobre la gráfica.

## Tipos gráficos más usados con Matplotlib

A continuación, se van a revisar algunos de los tipos de gráficos más usados con Matplotlib:

## Histogramas

Los histogramas son un tipo de gráficos que permiten visualizar las frecuencias de aparición de un conjunto de datos en forma de barras. La superficie de cada barra define una categoría que es proporcional a la frecuencia de los valores representados.

Para dibujar un histograma, se usa la función **plt.hist()**, que toma como argumentos los valores que se van a considerar para crear las categorías y el número de categorías que se quieren considerar. Si no se indica ningún valor para el número de categorías, toma por defecto el valor 10 (Figura 7.55.), cuaderno *Ejemplo\_55.ipynb*.

```
%matplotlib inline
import matplotlib.pyplot as plt
import numpy as np
y = np.random.randn(1000)
plt.hist(y,15)
plt.show()
```

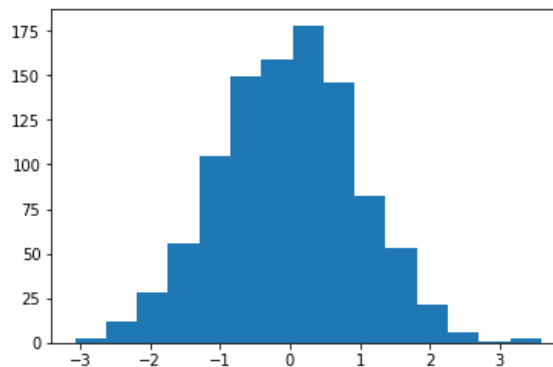


Figura 7.55. Historiograma

## Diagramas de barras

Los diagramas de barras son un tipo de gráficos que representan los datos como barras rectangulares, horizontales o verticales, cuyas dimensiones son proporcionales a los valores que representan los datos. Este tipo de gráficos permite comparar dos o más valores.



### Diagrama de barras verticales

Para dibujar un diagrama de barras se utiliza la función **plt.bar()**, que tiene como parámetros (Figura 7.56.), cuaderno *Ejemplo\_55.ipynb*:

- **Argumento left.** Es una lista de las coordenadas del eje x donde quedará situada gráficamente la esquina izquierda de una barra.
- **Argumento height.** Es una lista con la altura de las barras.
- **Argumento width.** Representa el ancho de la barra, que por defecto es 0.8.
- **Argumento color.** Representa el color de las barras.
- **Argumentos xerr, yerr.** Representa barras de error sobre el diagrama de barras.
- **Argumento bottom.** Representa las coordenadas inferiores en el eje y de las barras, en el caso de que se quiera empezar a dibujar por encima del eje x.

```
%matplotlib inline
import matplotlib.pyplot as plt
left=range(5)
p1 = plt.bar(left, height=[20,35,30,35,27], width=0.4, color='green', yerr=[2,3,4,1,2])
p2 = plt.bar(left, height=[25,32,34,20,25], width=0.4, color='red', bottom=[20,35,30,35,27], yerr=[3,5,2,3,3])
plt.legend(["Secuencia 1", "Secuencia 2"])
plt.xticks(range(5),["A","B","C","D","E"])
plt.show()
```

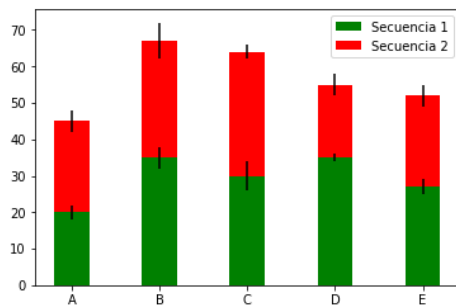


Figura 7.56. Diagrama de barras verticales.

### Diagrama de barras horizontales

Los mismos datos se podrían haber representado horizontalmente utilizando la función **plt.barh()**. Esta función tiene parámetros similares a la anterior (Figura 7.57.), cuaderno *Ejemplo\_57.ipynb*.

- **Argumento bottom.** Es una lista de las coordenadas del eje y que indican dónde se situarán gráficamente cada una de las barras del gráfico.
- **Argumento width.** Es una lista con la altura de las barras.
- **Argumento height.** Representa la altura de la barra, que por defecto es 0.8.
- **Argumento color.** Representa el color de las barras.
- **Argumentos xerr, yerr.** Representa las barras de error sobre el diagrama de barras.
- **Argumento left.** Representa las coordenadas inferiores en el eje x de las barras, en el caso de que se quiera empezar a dibujar por encima del eje y.

```
%matplotlib inline
import matplotlib.pyplot as plt
bottom=range(5)
p1 = plt.barh(bottom, width=[20,35,30,35,27], height=0.4, color='green', xerr=[2,3,4,1,2])
p2 = plt.barh(bottom, width=[25,32,34,20,25], height=0.4, color='red', left=[20,35,30,35,27], xerr=[3,5,2,3,3])
plt.legend(["Secuencia 1", "Secuencia 2"])
plt.yticks(range(5), ["A", "B", "C", "D", "E"])
plt.show()
```

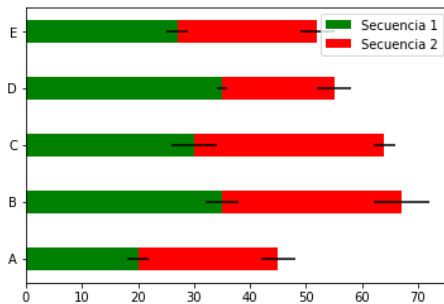


Figura 7.57. Diagrama de barras horizontal.

### Diagramas circulares

Se trata de un gráfico en forma de círculo que se encuentra dividido en sectores, de manera que la superficie ocupada por el sector representa proporcionalmente a la cantidad descrita.

Para representar gráficos circulares se utiliza la función **plt.pie()**, que toma como entrada un array X con valores. Así, para cada valor x del array X, la función genera sectores que son proporcionales a  $x/\text{Suma}(X)$ . Si la suma fuera menor que 1, se representan los valores directamente. Los sectores se dibujan empezando desde el eje horizontal, en el lado derecho y en sentido contrario a las manecillas del reloj.

## Parámetros

La función tiene los siguientes parámetros (Figura 7.58.), cuaderno *Ejemplo\_58.ipynb*:

- **explode**: es un array de la misma longitud que el array de valores X, cuyos valores indican la fracción del radio que va a separar el sector del centro del círculo.
- **colors**: es la lista de colores que se usarán cíclicamente para los sectores. En caso de no indicarlos, sigue una progresión automática de azul, verde, rojo, cian, magenta...
- **labels**: es una lista de las etiquetas asociadas a cada sector.
- **labeldistance**: es la distancia respecto al centro del gráfico en la que se dibuja una etiqueta.
- **autopct**: esta función o cadena de formateo permite etiquetar los sectores con los valores numéricos que representan.
- **pctdistance**: es la distancia respecto al centro del gráfico donde se dibujan los valores numéricos.
- **shadow**: dibuja un sombreado en los sectores y el gráfico.

```
%matplotlib inline
import matplotlib.pyplot as plt
x = [4,7,5,3,9,15]
labels = ['Madrid', 'Barcelona', 'Valencia', 'Zaragoza', 'Sevilla', 'Granada']
colors = ['yellow', 'blue', 'red', 'green', 'brown', 'orange']
explode = [0.2, 0.2, 0, 0, 0.1, 0.1]
plt.pie(x, labels=labels, labeldistance=1.3, explode=explode, autopct='%1.1f%%',
        colors=colors, pctdistance=0.5, shadow=True);
plt.show()
```

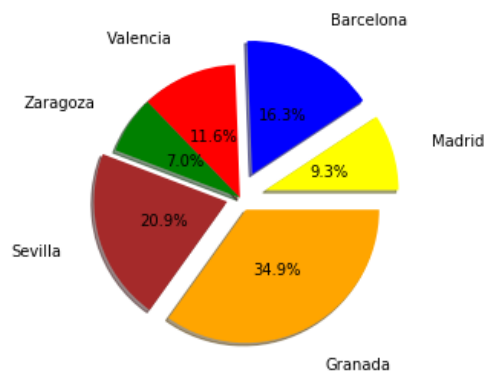


Figura 7.58. Diagrama circular.

## Diagramas de dispersión

Un diagrama de dispersión dibuja los valores de dos conjuntos de datos como una colección de puntos desconectados. Las coordenadas de cada punto se obtienen de cada conjunto (la coordenada x del primer conjunto y la coordenada y del segundo conjunto). Este tipo de dibujos facilita el descubrimiento de correlaciones u otro de tipo de relaciones entre los puntos considerados.

Para representar diagramas de dispersión se utiliza la función **plt.scatter()** que toma como entrada dos arrays unidimensionales de la misma longitud.

## Parámetros

Los principales argumentos que admite la función son:

- **s**: establece el tamaño de los puntos en pixel\*pixel. Se puede indicar un único valor, en cuyo caso se aplicará a todos los puntos, o bien especificar un array de la misma longitud que los arrays de entrada. En este último caso, se especifica un tamaño particular para cada punto.
- **c**: color de los puntos. Se puede indicar un único valor, en cuyo caso se aplicará a todos los puntos, o bien especificar un array de la misma longitud que los arrays de entrada. En este último caso, se especifica un color particular para cada punto. Para especificar los colores se pueden usar los códigos de colores.
- **marker**: especifica el tipo de punto que se va a dibujar de acuerdo a la tabla 6.2.:

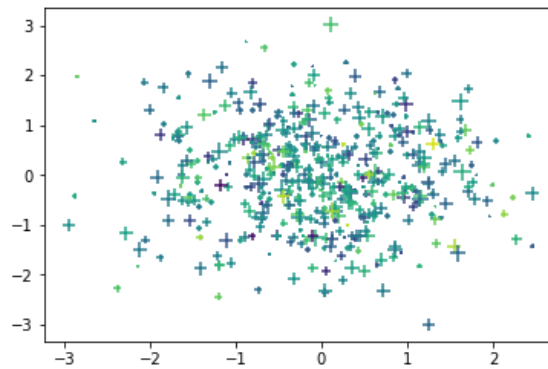
Cadena	Código
o	Círculo
V	Triángulo hacia abajo
^	Triángulo hacia arriba
<	Triángulo hacia la izquierda
>	Triángulo hacia la derecha
s	Cuadrado
p	Pentágono
h	Hexágono
+	Signo +
x	Cruz
d	Diamante
8	Octógono

**Tabla 6.2.** Tabla de configuración.

**Ejemplo**

En la Figura 7.59, cuaderno *Ejemplo\_59.ipynb*, se muestra un ejemplo de diagrama de dispersión.

```
%matplotlib inline
import matplotlib.pyplot as plt
import numpy as np
x = np.random.randn(1000)
y = np.random.randn(1000)
colores = np.random.randn(1000)
tam = 50*np.random.randn(1000)
plt.scatter(x,y,s=tam,c=colores,marker='+')
plt.show()
```



**Figura 7.59.** Diagrama de dispersión.

## Diagramas de contornos

Las líneas de contorno para funciones de dos variables son curvas en las que la función toma un valor constante, es decir,  $f(x,y)=C$  siendo  $C$  una constante. La densidad de las líneas indica la pendiente de la función.

Para dibujar un diagrama de contornos, existen las funciones **plt.contour()**, que toman como entrada un array de 2 dimensiones y, opcionalmente, el número niveles de líneas de contorno (Figura 7.60.), cuaderno *Ejemplo\_60.ipynb*. Las líneas se representan en diferentes colores, de los valores más pequeños a los más grandes —desde el azul oscuro para áreas de bajo volumen hasta el rojo para áreas de alto volumen—.

```
%matplotlib inline
import matplotlib.pyplot as plt
import numpy as np
valores = np.random.rand(21,31)
plt.contour(valores,5)
plt.show()
```

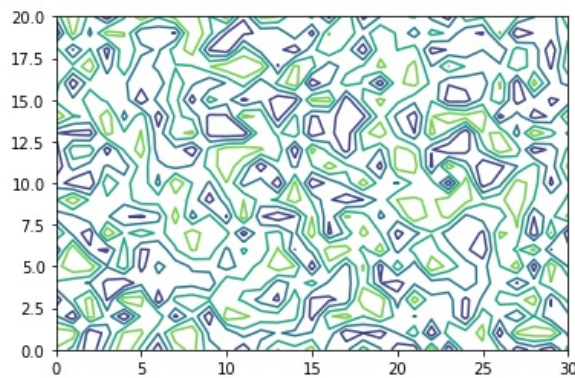


Figura 7.60. Diagrama de contorno.

### Diagrama de caja

Un diagrama de caja es un gráfico que representa los cuartiles de un conjunto de datos. Este tipo de diagramas se construyen mediante la función **plt.boxplot** (Figura 7.61.), cuaderno *Ejemplo\_61.ipynb*.

```
%matplotlib inline
import matplotlib.pyplot as plt
import numpy as np
valores = np.random.rand(100)
plt.boxplot(valores)
plt.show()
```

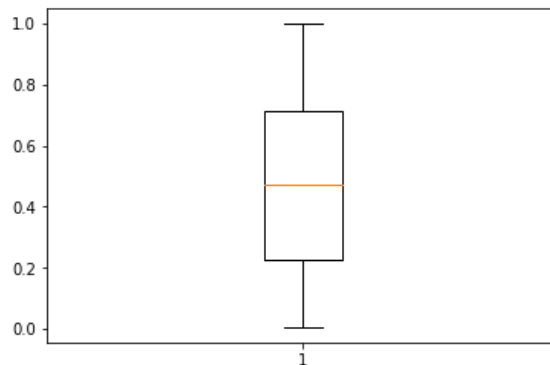


Figura 7.61. Diagrama de caja.

### Conjunto de gráficos

En Matplotlib existe la posibilidad de crear gráficas formadas por un conjunto de gráficos. A estos gráficos se les denomina subgráficos. Para ello, se utiliza la función **plt.figure()**, la cual genera un objeto de tipo **Figure**, que es un contenedor donde se pueden añadir subgráficos utilizando el método **add\_subplot()**.

La invocación del método devuelve un objeto de tipo **axe**, que es un área donde se puede dibujar.

#### Parámetros

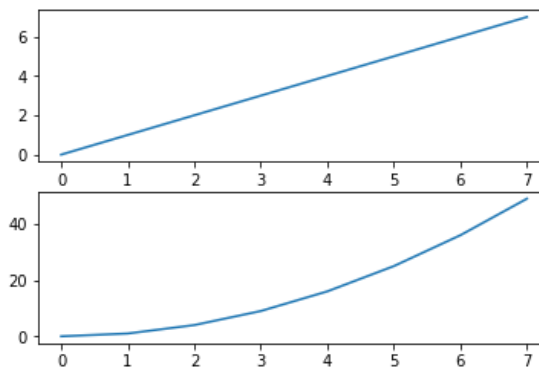
Los parámetros del método permiten especificar de manera matricial dónde se insertarán los subgráficos:

- **numrows**: especifica el número de filas.
- **numcols**: especifica el número de columnas.
- **fignum**: especifica un número que varía entre 1 y  $\text{numrows} \times \text{numcols}$ , que indica el subgráfico actual, avanzando de izquierda a derecha y de arriba abajo. Así, 1 indica la esquina superior izquierda y  $\text{numrows} \times \text{numcols}$  indica la esquina inferior derecha.

**Ejemplo**

Los objetos de tipo **axe** pueden invocar la función **plt.plot()** para representar puntos (Figura 7.62.), cuaderno *Ejemplo\_62.ipynb*.

```
%matplotlib inline
import matplotlib.pyplot as plt
fig = plt.figure()
x = range(8)
ax1=fig.add_subplot(211)
ax1.plot(x, [xi for xi in x])
ax2=fig.add_subplot(212)
ax2.plot(x, [xi**2 for xi in x])
plt.show()
```



**Figura 7.62.** Ejemplo de subgráfico.

## V. Manipulación y análisis de datos con Pandas

Pandas es una librería construida sobre NumPy que ofrece estructuras de datos de alto nivel que facilitan el análisis de datos desde Python. Se van a estudiar dos estructuras de datos:

**1**

**Series**

**2**

**DataFrame**

Primero, importación de librerías. Antes de trabajar con estas estructuras se importan las librerías necesarias:

```
import numpy as np
```

```
import pandas as pd
```

```
from pandas import *
```



## 5.1. Series

Una serie es un objeto, como un array, que está formado por un array de datos y un array de etiquetas denominado índice.

### Tipos de arrays de datos

El array de datos puede ser de 3 tipos:

#### ndarray

Si el array de datos en un ndarray, entonces el índice debe ser de la misma longitud que el array de datos (Figura 7.63.), cuaderno ejercicio *Ejemplo\_63.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
s=Series([1,2,3,4,5], index=['a','b','c','d','e'])
s
```

```
a    1
b    2
c    3
d    4
e    5
dtype: int64
```

Figura 7.63. Ejemplo de ndarray.

Si ningún índice es pasado, entonces se crea uno formado por valores que van desde  $[0, \dots, \text{len}(\text{data})-1]$  (Figura 7.64.), cuaderno ejercicio *Ejemplo\_64.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
s=Series([1,2,3,4,5])
s
```

```
0    1
1    2
2    3
3    4
4    5
dtype: int64
```

Figura 7.64. Ejemplo de ndarray sin índice.

## Diccionario

Si se pasa un índice, los valores del diccionario son asociados a los valores del índice (Figura 7.65.), cuaderno ejercicio *Ejemplo\_65.ipynb*

```
import numpy as np
import pandas as pd
from pandas import *
d={'a':0., 'b':1., 'c':2.}
s=Series(d, index=['b', 'c', 'd', 'a'])
s
```

```
b    1.0
c    2.0
d    NaN
a    0.0
dtype: float64
```

**Figura 7.65.** Creación de un ndarray usando un diccionario.

Si no se proporciona un índice, entonces se construye a partir de las claves ordenadas del diccionario (Figura 7.66.), cuaderno ejercicio *Ejemplo\_66.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
d={'a':0., 'b':1., 'c':2.}
s=Series(d)
s
```

```
a    0.0
b    1.0
c    2.0
dtype: float64
```

**Figura 7.66.** Creación de ndarray usando un diccionario sin índice.

## Valor escalar

En este caso, se debe proporcionar un índice y el valor será repetido tantas veces como la longitud del índice (Figura 7.67.), cuaderno ejercicio *Ejemplo\_67.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
s=Series(5., index=['a', 'b', 'c', 'd', 'e'])
s
```

```
a    5.0
b    5.0
c    5.0
d    5.0
e    5.0
dtype: float64
```

**Figura 7.67.** Creación de ndarray un valor escalar.

## Características de las series

Algunas características de las series son:

#### Similares a ndarray

Las series actúan de forma similar a ndarray, siendo un argumento válido de funciones de NumPy (Figura 7.68.), cuaderno ejercicio *Ejemplo\_68.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
s=Series([1,2,3,4,5], index=['a','b','c','d','e'])
s[0]
```

1

```
s[s > s.median()]
```

```
d    4
e    5
dtype: int64
```

```
s[[4,3,1]]
```

```
e    5
d    4
b    2
dtype: int64
```

**Figura 7.68.** Uso de ndarray con NumPy.

**Diccionario de tamaño fijo**

Las series actúan como un diccionario de tamaño fijo en el que se pueden gestionar los valores a través de los índices (Figura 7.69.), cuaderno ejercicio *Ejemplo\_69.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
s=Series([1,2,3,4,5], index=['a','b','c','d','e'])
s['a']
```

1

```
s['e']
```

5

```
s
```

```
a    1
b    2
c    3
d    4
e    5
dtype: int64
```

```
'e' in s
```

True

**Figura 7.69.** Uso de series como diccionarios.

**Permiten operaciones vectoriales**

Sobre las series se pueden realizar operaciones vectoriales (Figura 7.70.) cuaderno ejercicio *Ejemplo\_70.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
s=Series([1,2,3,4,5], index=['a','b','c','d','e'])
s+s
```

```
a    2
b    4
c    6
d    8
e   10
dtype: int64
```

**Figura 7.70.** Operaciones vectoriales sobre

series.

### Alineamiento

La principal diferencia entre series y ndarray es que las operaciones entre series alinean automáticamente los datos basados en las etiquetas, de forma que pueden realizarse cálculos sin tener en cuenta si las series sobre las que se operan tienen las mismas etiquetas.

Obsérvese que el resultado de una operación entre series no alineadas será la unión de los índices. Si una etiqueta no se encuentra en una de las series, entonces se le asocia el valor nulo. Los datos y el índice de una serie tienen un atributo name que puede asociarle un nombre (Figura 7.71.), cuaderno ejercicio *Ejemplo\_71.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
s=Series([1,2,3,4,5], index=['a','b','c','d','e'])
s.name = "datos"
s
```

a	1
b	2
c	3
d	4
e	5

Name: datos, dtype: int64

Figura 7.71. Uso del atributo name.

## 5.2. Dataframes

Un dataframe es una estructura que contiene una colección ordenada de columnas, cada una de las cuales puede tener valores de diferentes tipos. Está formado por datos, opcionalmente un índice (etiquetas de las filas) y un conjunto de columnas (etiquetas de las columnas). En caso de no existir un índice, se genera a partir de los datos.

### 5.2.1. Datos de un dataframe

### Diccionario de ndarrays

Los arrays (Figura 7.72.), cuaderno ejercicio *Ejemplo\_72.ipynb* deben ser de la misma longitud. En caso de existir un índice, este debe ser de la misma longitud que los arrays y, en caso de no existir, se genera como

```
import numpy as np
import pandas as pd
from pandas import *
d = {'one': [1., 2., 3., 4.], 'two': [4., 3., 2., 1.]}
DataFrame(d)
```

índice la secuencia de números 0... longitud(array)-1.

	one	two
0	1.0	4.0
1	2.0	3.0
2	3.0	2.0
3	4.0	1.0

Figura 7.72. Diccionario de ndarrays.

### Lista de diccionarios

Lista de diccionarios (Figura 7.73.), cuaderno ejercicio *Ejemplo\_73.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
data2 = [{'a': 1, 'b': 2}, {'a': 5, 'b': 10, 'c': 20}]
DataFrame(data2)
```

	a	b	c
0	1	2	NaN
1	5	10	20.0

```
DataFrame(data2, index=['first', 'second'])
```

	a	b	c
first	1	2	NaN
second	5	10	20.0

Figura 7.73. Lista de diccionarios.

Diccionario de tuplas

Diccionario de tuplas (Figura 7.74.), cuaderno ejercicio *Ejemplo\_74.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *

DataFrame({
    ('a','b'): {('A','B'):1, ('A','B'):2},
    ('a','a'): {('A','C'):3, ('A','B'):4},
    ('a','c'): {('A','B'):5, ('A','C'):6},
    ('b','a'): {('A','C'):7, ('A','B'):8},
    ('b','b'): {('A','D'):9, ('A','B'):10}
})
```

		a			b	
		b	a	c	a	b
A	B	2.0	4.0	5.0	8.0	10.0
	C	NaN	3.0	6.0	7.0	NaN
	D	NaN	NaN	NaN	NaN	9.0

Figura 7.74. Diccionario de tuplas.

## Diccionario de series

(Figura 7.75), cuaderno ejercicio *Ejemplo\_75.ipynb*. El índice que resulta es la unión de los índices de las series. En caso de existir diccionarios anidados, primero se convierten en diccionarios. Además, si no se pasan columnas, se toman como tales la lista ordenada de las claves de los diccionarios.

```
import numpy as np
import pandas as pd
from pandas import *
d = {'one': Series([1.,2.,3.], index=['a','b','c']),
     'two': Series([1.,2.,3.,4.], index=['a','b','c','d'])
}
df = DataFrame(d)
df
```

	one	two
a	1.0	1.0
b	2.0	2.0
c	3.0	3.0
d	NaN	4.0

```
DataFrame(d, index=['d','b','a'])
```

	one	two
b	2.0	2.0
a	1.0	1.0
d	NaN	4.0

**Figura 7.75.** Diccionario de series.

Puede accederse a las etiquetas de las columnas y de las filas a través de los atributos índice y columnas. Téngase en cuenta que cuando un conjunto particular de columnas se pasa como argumento con el diccionario, entonces las columnas sobrescriben en las claves del diccionario.



### Array estructurado

Array estructurado (Figura 7.76.), cuaderno ejercicio *Ejemplo\_76.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
data = np.zeros((2,), dtype=[
    ('A', 'i4'), ('B', 'f4'), ('C', 'a10'),
])
data[:] = [(1, 2., 'Hello'), (2, 3., "World")]
DataFrame(data)
```

	A	B	C
0	1	2.0	b'Hello'
1	2	3.0	b'World'

Figura 7.76. Array estructurado

## 5.2.2. Operaciones sobre dataframes

### Manipulación de un dataframe

Un dataframe es como un diccionario de series indexado, por lo que se pueden usar las mismas operaciones utilizadas con los diccionarios (Figura 7.77.), cuaderno ejercicio *Ejemplo\_77.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
d = { 'one': Series([1., 2., 3.], index=['a', 'b', 'c']),
      'two': Series([1., 2., 3., 4.], index=['a', 'b', 'c', 'd'])
}
df = DataFrame(d)
df['one']
```

```
a    1.0
b    2.0
c    3.0
d    NaN
Name: one, dtype: float64
```

```
df['three'] = df['one'] * df['two']
```

```
df['flag'] = df['one'] > 2
```

```
df
```

	one	two	three	flag
a	1.0	1.0	1.0	False
b	2.0	2.0	4.0	False
c	3.0	3.0	9.0	True
d	NaN	4.0	NaN	False

Figura 7.77. Manipulación de un dataframe.

**Borrado de columnas**

Se pueden borrar columnas (Figura 7.78.), cuaderno ejercicio *Ejemplo\_78.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
d = { 'one': Series([1.,2.,3.], index=['a','b','c']),
      'two': Series([1.,2.,3.,4.], index=['a','b','c','d'])
}
df = DataFrame(d)
df['one']
```

```
a    1.0
b    2.0
c    3.0
d    NaN
Name: one, dtype: float64
```

```
df['three'] = df['one'] * df['two']
```

```
df['flag'] = df['one'] > 2
```

```
del df['two']
```

```
three = df.pop('three')
```

```
df
```

	one	flag
a	1.0	False
b	2.0	False
c	3.0	True
d	NaN	False

Figura 7.78. Borrado de columnas.

## Inserción de valores

Cuando se inserta un valor escalar, entonces se propaga a toda la columna (Figura 7.79.),cuaderno ejercicio *Ejemplo\_79.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
d = { 'one': Series([1.,2.,3.], index=['a','b','c']),
      'two': Series([1.,2.,3.,4.], index=['a','b','c','d'])
}
df = DataFrame(d)
df['three'] = df['one'] * df['two']
df['flag'] = df['one'] > 2
del df['two']
three = df.pop('three')
df['foo'] = 'bar'
df
```

	one	flag	foo
a	1.0	False	bar
b	2.0	False	bar
c	3.0	True	bar
d	NaN	False	bar

**Figura 7.79.** Propagación a toda la columna.

## Creación de nuevo índice

Cuando se inserta una serie que no tiene el mismo índice, se crea el índice para el dataframe (Figura 7.80.), cuaderno ejercicio *Ejemplo\_80.ipynb*:

```
import numpy as np
import pandas as pd
from pandas import *
d = { 'one': Series([1.,2.,3.], index=['a','b','c']),
      'two': Series([1.,2.,3.,4.], index=['a','b','c','d'])
}
df = DataFrame(d)
df['three'] = df['one'] * df['two']
df['flag'] = df['one'] > 2
del df['two']
three = df.pop('three')
df['foo'] = 'bar'
df['one_trunc'] = df['one'][:2]
df
```

	one	flag	foo	one_trunc
a	1.0	False	bar	1.0
b	2.0	False	bar	2.0
c	3.0	True	bar	NaN
d	NaN	False	bar	NaN

**Figura 7.80.** Inserción de una serie con distinto índice.

### Selección del punto de inserción

Las columnas, por defecto, se insertan al final, sin embargo, se puede elegir el lugar de inserción usando la función **insert** (Figura 7.81.), cuaderno ejercicio *Ejemplo\_81.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
d = { 'one': Series([1.,2.,3.], index=['a','b','c']),
      'two': Series([1.,2.,3.,4.], index=['a','b','c','d'])
}
df = DataFrame(d)
df['three'] = df['one'] * df['two']
df['flag'] = df['one'] > 2
del df['two']
three = df.pop('three')
df['foo'] = 'bar'
df['one_trunc'] = df['one'][:2]
df.insert(1, 'bar', df['one'])
df
```

	one	bar	flag	foo	one_trunc
a	1.0	1.0	False	bar	1.0
b	2.0	2.0	False	bar	2.0
c	3.0	3.0	True	bar	NaN
d	NaN	NaN	False	bar	NaN

Figura 7.81. Inserción de un dataframe.

### Indexación/Selección

#### Selección por columnas

Se puede seleccionar por columnas (Figura 7.82.), cuaderno ejercicio *Ejemplo\_82.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
d = { 'one': Series([1.,2.,3.], index=['a','b','c']),
      'two': Series([1.,2.,3.,4.], index=['a','b','c','d'])
}
df = DataFrame(d)
df['one']
```

a	1.0
b	2.0
c	3.0
d	NaN

Name: one, dtype: float64

Figura 7.82. Selección por columnas.

### Selección por etiqueta

Se puede seleccionar por etiqueta (Figura 7.83.), cuaderno ejercicio *Ejemplo\_83.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
d = { 'one': Series([1.,2.,3.], index=['a','b','c']),
      'two': Series([1.,2.,3.,4.], index=['a','b','c','d'])
}
df = DataFrame(d)
df.loc['b']
```

one	2.0
two	2.0

Name: b, dtype: float64

**Figura 7.83.** Selección por etiqueta.

### Selección por entero

Selección de fila por entero (Figura 7.84.), cuaderno ejercicio *Ejemplo\_84.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
d = { 'one': Series([1.,2.,3.], index=['a','b','c']),
      'two': Series([1.,2.,3.,4.], index=['a','b','c','d'])
}
df = DataFrame(d)
df.iloc[2]
```

one	3.0
two	3.0

Name: c, dtype: float64

**Figura 7.84.** Selección por entero.

### Selección por rangos

Selección por rangos (Figura 7.85.), cuaderno ejercicio *Ejemplo\_85.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
d = { 'one': Series([1.,2.,3.], index=['a','b','c']),
      'two': Series([1.,2.,3.,4.], index=['a','b','c','d'])
}
df = DataFrame(d)
df.iloc[5:10]
```

one two

Figura 7.85. Selección por rangos.

### Alineación y aritmética

#### Alineación

Cuando se alinea un dataframe, se realiza sobre las columnas y el índice, de manera que se hace la unión de las etiquetas de las columnas y de las filas (Figura 7.86.), cuaderno ejercicio *Ejemplo\_86.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
randn = np.random.rand
df = DataFrame( randn(10,4), columns=['A','B','C','D'] )
df2 = DataFrame( randn(7, 3), columns=['A','B','C'] )
df+df2
```

	A	B	C	D
0	0.957629	1.420427	1.208260	NaN
1	0.049203	1.433719	1.549568	NaN
2	1.270981	0.644690	1.884569	NaN
3	0.323450	0.338493	1.301843	NaN
4	1.408984	0.811337	1.375654	NaN
5	1.555031	0.975821	0.722530	NaN
6	1.348029	1.079629	0.404529	NaN
7	NaN	NaN	NaN	NaN
8	NaN	NaN	NaN	NaN
9	NaN	NaN	NaN	NaN

Figura 7.86. Alineamiento.

### Alineación de operación con dataframe y series

Cuando se operan con dataframe y series, se alinea el índice de las series sobre las columnas del dataframe (Figura 7.87.), cuaderno ejercicio *Ejemplo\_87.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
randn = np.random.rand
df = DataFrame( randn(10,4), columns=['A','B','C','D'] )
df=df.iloc[0]
```

	A	B	C	D
0	0.000000	0.000000	0.000000	0.000000
1	-0.484808	0.082457	0.359350	0.352122
2	-0.539635	-0.545852	0.838071	0.433856
3	-0.855391	-0.426257	0.687573	0.393009
4	-0.041807	-0.751054	-0.034094	0.573807
5	-0.523541	-0.376184	0.106900	0.699860
6	-0.141314	-0.636719	0.412637	-0.032672
7	-0.081670	0.048396	0.357219	0.480282
8	-0.925606	-0.621518	0.835190	-0.009714
9	-0.512766	-0.423692	0.258417	0.480712

**Figura 7.87.** Alineamiento de dataframe y serie.



## Operaciones con escalares

Se pueden hacer operaciones con escalares (Figura 7.88.), cuaderno ejercicio *Ejemplo\_88.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
randn = np.random.rand
df = DataFrame( randn(10,4), columns=['A','B','C','D'] )
df*5+2
```

	A	B	C	D
0	2.322141	5.327164	2.880365	4.138989
1	5.756474	5.063875	6.789093	6.764717
2	5.683400	2.664712	2.031166	6.859473
3	4.144601	4.184243	6.805077	2.725346
4	4.444096	5.666809	6.295436	5.824817
5	6.992702	3.663723	4.154809	5.610130
6	5.531661	5.667566	6.283478	4.461741
7	6.485985	3.629512	2.656262	6.062589
8	6.738205	3.170034	6.605348	2.714790
9	2.993485	3.824260	4.099124	2.432600

Figura 7.88. Operaciones con escalares.

## Operaciones lógicas

Se pueden hacer operaciones lógicas (Figura 7.89.), cuaderno ejercicio *Ejemplo\_89.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
randn = np.random.rand
df1 = DataFrame( {'a':[1,0,1], 'b':[0,1,1]}, dtype=bool )
df2 = DataFrame( {'a':[0,1,1], 'b':[1,1,0]}, dtype=bool )
df1 & df2
```

	a	b
0	False	False
1	False	True
2	True	False

Figura 7.89. Operaciones lógicas.

## Transposición

Para transponer, se utiliza el atributo T (Figura 7.90.), cuaderno ejercicio *Ejemplo\_90.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
randn = np.random.rand
df = DataFrame( randn(10,4), columns=['A', 'B', 'C', 'D'] )
df[:5].T
```

	0	1	2	3	4
A	0.147590	0.752690	0.069236	0.521695	0.754726
B	0.746538	0.401150	0.328227	0.981864	0.825076
C	0.276812	0.522900	0.114420	0.716107	0.358262
D	0.098769	0.473946	0.733130	0.674665	0.759104

**Figura 7.90.** Transposición.

## Visualización

Hay varias formas de visualizar la información de un dataframe:

### A partir del nombre

Todos los datos a partir del nombre del dataframe (Figura 7.91.), cuaderno ejercicio *Ejemplo\_91.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
randn = np.random.rand
df = DataFrame( randn(10,4), columns=['A', 'B', 'C', 'D'] )
df
```

	A	B	C	D
0	0.879912	0.754664	0.450893	0.548262
1	0.668827	0.016458	0.182358	0.061669
2	0.626958	0.423782	0.552083	0.869737
3	0.923197	0.531068	0.804448	0.987810
4	0.216840	0.254141	0.101713	0.054160
5	0.301915	0.617831	0.483386	0.660352
6	0.641980	0.237544	0.286541	0.680817
7	0.431369	0.572311	0.450401	0.731337
8	0.055978	0.860867	0.160450	0.012160
9	0.346314	0.077852	0.430650	0.481859

**Figura 7.91.** Datos de un dataframe.

## Resumen

Un resumen de la información contenida, usando el método `info()` (Figura 7.92.), cuaderno ejercicio *Ejemplo\_92.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
randn = np.random.rand
df = DataFrame( randn(10,4), columns=['A','B','C','D'] )
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10 entries, 0 to 9
Data columns (total 4 columns):
A      10 non-null float64
B      10 non-null float64
C      10 non-null float64
D      10 non-null float64
dtypes: float64(4)
memory usage: 400.0 bytes
```

**Figura 7.92.** Datos de un dataframe usando `info()`.

## Cadena

Como una cadena usando el método `to_string()` (Figura 7.93.), cuaderno ejercicio *Ejemplo\_93.ipynb*.

```
import numpy as np
import pandas as pd
from pandas import *
randn = np.random.rand
df = DataFrame( randn(10,4), columns=['A','B','C','D'] )
print(df.to_string())
```

	A	B	C	D
0	0.034685	0.914251	0.946942	0.083720
1	0.972669	0.142089	0.788396	0.013883
2	0.325790	0.850953	0.380769	0.602521
3	0.039785	0.360459	0.619699	0.168697
4	0.416723	0.932684	0.873190	0.249169
5	0.855723	0.608578	0.062883	0.336972
6	0.836523	0.497623	0.205874	0.121346
7	0.904203	0.011100	0.847687	0.518168
8	0.352006	0.465943	0.377309	0.191254
9	0.280667	0.284980	0.545421	0.602264

usando `to_string()`.

**Figura 7.93.** Datos de un dataframe

## VI. Resumen



PyCharm

En esta unidad, se han presentado las principales librerías científicas de Python, las cuales se utilizan para realizar análisis de datos.

La librería NumPy facilita la manipulación de los arrays de datos como si fueran tipos de datos básicos, lo que posibilita realizar operaciones con los mismos sin tener que acudir a estructuras de datos más complejas.

La librería Matplotlib permite la representación gráfica de puntos, pudiéndolos mostrar en diferentes tipos de gráficos. Asimismo, es posible utilizar diferentes tipos de gráficos para representar los datos. Los gráficos se pueden decorar con distintos elementos, como son el título del gráfico o los valores de los ejes.

Por último, la librería Pandas permite manipular los datos en general de una manera simple, lo que facilita operaciones propias del análisis de datos tales como la limpieza o normalización de datos, o bien la selección de determinados datos.

Como información complementaria se recomienda visualizar la siguiente clase magistral: *Introducción a DataFrames con Python*, y descargar y ejecutar la mencionada clase magistral con los archivos que se facilitan.



Para ello, **en este enlace** se puede descargar una carpeta comprimida con los archivos que se necesitarán:

- LMA\_INTRO\_DATAFRAMES\_PYTHON.ipynb también disponible en html
- LMA\_INTRO\_DATAFRAMES\_PYTHON.html
- Iris.csv
- Iris\_merge.csv



Respecto a la clase, es posible verla en el siguiente enlace:

- Juan Manuel Moreno. *LMA BIG DATA: Introducción a DataFrames con Python*. IMF Business School, YouTube, 29/01/2019. [En línea] URL disponible en: <https://www.youtube.com/watch?v=CZcS8PaLiwY>

## VII. Apéndice

En las unidades "Fundamentos de programación en Python", "Fundamentos de tecnologías de internet" y "Fundamentos de tratamiento de datos para Data Science" de este módulo se estudiaron las funcionalidades generales del lenguaje de programación Python con el fin de, posteriormente, orientarse a las librerías científicas para análisis de datos. En este apartado, se incluyen una serie de ejercicios complementarios en los que convergen diferentes técnicas como manejo de archivos, creación de funciones y desarrollo de diferentes estructuras de datos. Se recomienda descargar y ejecutar los cuadernos Jupyter Notebook, entre los que se incluye:

- Variables numéricas y de texto, además operadores booleanos.
- Estructuras de control y ciclos.
- Tuplas, listas y diccionarios.
- Funciones, JSON y CSV.

- DataFrame y Pyplot.



Se pueden descargar todos los *notebooks* con los ejercicios complementarios en este enlace:

- [Ejercicios\\_apéndice](#)

## VIII. Caso práctico

### Se pide

**a**

Crear una serie denominada "Sueldos\_hombres" que contenga los datos [2500,1800,1900,2000,2100] y cuyos índices sean ["Euskadi", "Murcia", "Madrid", "Barcelona", "Zaragoza"], y otro denominada "Sueldos\_mujeres" que contenga los datos [2300,1600,1980,1900,2150] y cuyos índices sean ["Euskadi", "Murcia", "Madrid", "Barcelona", "Córdoba"]

**b**

Usando las series anteriores obtén otra serie que contenga la diferencia de sueldos entre mujeres y hombres.

**c**

Usando la serie anterior que contiene las diferencias de sueldos, obtén las ciudades donde la diferencia de sueldos es mayor de 100 euros.

### Solución

1

```
#Solucion A
import pandas as pd
sueldoshom = pd.Series( [2500,1800,1900,2000,2100],
                        index = ["Euskadi", "Murcia", "Madrid", "Barcelona", "Zaragoza"]
                        )
print(sueldoshom)
sueldosmuj = pd.Series( [2300,1600,1980,1900,2150],
                        index = ["Euskadi", "Murcia", "Madrid", "Barcelona", "Córdoba"]
                        )
print(sueldosmuj)
```

```
Euskadi      2500
Murcia       1800
Madrid       1900
Barcelona    2000
Zaragoza     2100
dtype: int64
Euskadi      2300
Murcia       1600
Madrid       1980
Barcelona    1900
Córdoba      2150
dtype: int64
```

2

```
#Solucion B
difsueldos = sueldoshom - sueldosmuj
difsueldos
```

```
Barcelona    100.0
Córdoba      NaN
Euskadi      200.0
Madrid       -80.0
Murcia       200.0
Zaragoza     NaN
dtype: float64
```

3

```
#Solución C
mascara = (difsueldos > 100)
difsueldos[mascara]
```

```
Euskadi      200.0
Murcia       200.0
dtype: float64
```

## Recursos

### Bibliografía

- **Pandas: powerful Python data analysis toolkit:**

*Pandas: powerful Python data analysis toolkit.* 2017. [En línea] URL disponible en <http://pandas.pydata.org/pandas-docs/stable/index.html>

- **Python Data Analysis Library—pandas: Python Data Analysis Library:**

*Python Data Analysis Library—pandas: Python Data Analysis Library.* [En línea] URL disponible en <https://pandas.pydata.org/>

- **Python for Data Analysis :**

McKinney, W. *Python for Data Analysis*. EE. UU.: O'Reilly, 2012.

- **Tentative Numpy Tutorial. :**

*Tentative Numpy Tutorial.* [En línea] URL disponible en [http://wiki.scipy.org/Tentative\\_NumPy\\_Tutorial](http://wiki.scipy.org/Tentative_NumPy_Tutorial)

- **Matplotlib :**

[En línea] URL disponible en <http://matplotlib.org/>

### Glosario.

- **Axe:** Es un objeto que representa un área donde se puede dibujar.
- **Dataframe:** Es una estructura que contiene una colección ordenada de columnas, cada una de las cuales puede tener valores de diferentes tipos. Está formada por datos y, opcionalmente, un índice (etiquetas de las filas) y un conjunto de columnas (etiquetas de las columnas). En caso de no existir un índice, se genera a partir de los datos.
- **Matplotlib:** Es una librería de Python para realizar gráficos. Se caracteriza porque es fácil de usar, flexible y se puede configurar de múltiples maneras.
- **Ndarray:** Son arrays multidimensionales donde todos sus elementos son del mismo tipo y están indexados por una tupla de números positivos.
- **Numpy:** El módulo NumPy (Numerical Python) es una extensión de Python que proporciona funciones y rutinas matemáticas para la manipulación de arrays y matrices de datos numéricos de una forma eficiente.
- **Pandas:** Es una librería construida sobre Numpy que ofrece estructuras de datos de alto nivel que facilitan el análisis de datos desde Python.
- **Series:** Una serie es un objeto, como un array, que está formado por un array de datos y un array de etiquetas denominado índice.
- **Subgráfico:** Son gráficas formadas por un conjunto de gráficos.