

# Pandas

Semestre 02, 2025

## Introducción

---

Pandas es una librería externa para Python de código abierto. Para poder utilizarse es necesario instalarla primero ya que no está incluida en la instalación de Python.

💡 El nombre viene de PANel DAta.

Es útil para análisis y manipulación de datos.

## Instalación

---

Pandas se puede instalar usando el manejador de paquetes de Thonny.

## Prueba de instalación

---

Deben correr el siguiente comando ya sea desde una terminal, Thonny, un script de Python o un Jupyter Notebook. Si el comando no da error quiere decir que está correctamente instalado.

```
1  
2 import pandas as pd  
3
```

## CSV

---

CSV significa Comma-Separated Values o valores separados por comas. Es un formato de archivo de texto plano que se utiliza para almacenar datos en forma tabular, como si fuera una hoja de cálculo.

Es posible verlos desde Excel o cualquier editor de texto.

```
1
2  nombre,edad,ciudad
3
4  Ana,23,Guatemala
5
6  Luis,30,Quetzaltenango
7
8  Marta,27,Antigua
9
```

En este archivo tenemos una tabla de 3 columnas: Nombre, Edad y Ciudad.

Usamos CSV porque:

- Es fácil de crear y leer.
- Es compatible con Excel, Google Sheets, bases de datos y lenguajes de programación.
- Es ligero y no depende de ningún software específico.

## Data Frames

---

Son las estructuras de datos que usa Pandas. Constan de 3 componentes:

1. Columnas
2. Índices
3. Datos

Crear un data frame es muy parecido a crear un diccionario

```
1
2  df = pd.DataFrame({
3
4      "Carnet": [123, 456, 789, 987, 654],
5
6      "Hoja de trabajo 1": [80, 85, 75, 40, 100],
7
8      "Hoja de trabajo 2": [90, 90, 75, 61, 100],
9}
```

```
10 "Hoja de trabajo 3": [80, 80, 0, 50, 85],  
11  
12 })  
13
```

Imprimir un data muestra de forma ordenada la información.

```
1  
2 print(df)  
3
```

## Instrucciones básicas

---

La forma del data frame nos da información sobre la cantidad de información que este tiene.

```
1  
2 print(df.shape)  
3
```

Si se requiere estadística básica es posible utilizar la función describe.

```
1  
2 print(df.describe())  
3
```

Es posible sacar estadística de forma más específica:

Promedio:

```
1  
2 print(df.mean())  
3
```

Cantidad de registros

```
1  
2 print(df.count())  
3
```

## Máximo

```
1  
2 print(df.max())  
3
```

## Mínimo

```
1  
2 print(df.min())  
3
```

## Desviación estándar

```
1  
2 print(df.std())  
3
```

## Solo una columna

```
1  
2 print(df["Hoja de trabajo 1"])  
3
```

## Solo el promedio de una columna

```
1  
2 print(df["Hoja de trabajo 1"].mean())  
3
```

También es posible generar sub data frames con ciertas columnas

```
1  
2 df2 = df[["Hoja de trabajo 1", "Hoja de trabajo 3"]]  
3  
4  
5  
6 print(df2)  
7  
8  
9  
10 print(df2.shape)  
11
```

Si se requiere obtener un valor específico es posible hacerlo:

```
1
2     valor = df["Hoja de trabajo 1"][0]
3
4
5
6     print(valor)
7
```

## Búsqueda de información

---

Poder acceder a información específica es uno de los poderes más grandes de Pandas. Para lograr esto se puede utilizar LOC ya que permite acceder a un grupo de filas y columnas usando etiquetas.

```
1
2     print(df.loc[df["Hoja de trabajo 1"] >= 61])
3
```

También es posible filtrar la información en el dataframe para limitar el resultado.

```
1
2     distinguidos = df.loc[df["Hoja de trabajo 1"] >= 80]
3
4
5
6     print(distinguidos)
7
8
9
10    print("\n--\n")
11
12
13
14    print(distinguidos[["Carnet", "Hoja de trabajo 1"]])
15
```

## Agregar información

---

```
1
2     nuvas_notas = [100, 100, 100, 100, 100]
3
4
5
6     df["Hoja de trabajo 4"] = nuvas_notas
7
```

## Actualizar información

---

Cambiarle el nombre a una columna:

```
1
2     df = df.rename(columns={"Hoja de trabajo 1": "Tarea 1"})
3
4
5
6     df = df.rename(columns={
7
8         "Hoja de trabajo 2": "Tarea 2",
9
10        "Hoja de trabajo 3": "Tarea 3"
11
12    })
13
```

cambiar el valor:

```
1
2     df.loc[df["Carnet"] == 123, "Hoja de trabajo 2"] = 95
3
```

## Eliminar información

---

```
1
2     del df["Hoja de trabajo 4"]
3
```

# Leer archivos

---

Previo a leer archivos es una buena práctica usar `getcwd` (get current working directory). Esto nos da el directorio actual en donde está corriendo el script.

```
1  
2 import os  
3  
4  
5  
6 path = os.getcwd()  
7  
8  
9  
10 print(path)  
11
```

Esto permite ahorrar tiempo en definir la ruta completa del archivo que queremos leer.

```
1  
2 file = path + "/canciones.csv"  
3  
4  
5  
6 datos = pd.read_csv(file)  
7  
8  
9  
10 print(datos)  
11
```

En `datos` se creo un dataframe con la información del archivo CSV. Ahora ya se puede hacer análisis completo sobre los datos.

# Escribir archivos

---

Para escribir a un archivo se convierte un dataframe a un csv.

```
1
```

```
2 import os
3
4 import pandas as pd
5
6
7
8 path = os.getcwd()
9
10
11
12 canciones = [
13
14 {"titulo": "Cancion 1", "autor": "Autor 1", "duracion": 3.5},
15
16 {"titulo": "Cancion 2", "autor": "Autor 1", "duracion": 2.5},
17
18 {"titulo": "Cancion 3", "autor": "Autor 3", "duracion": 4.2},
19
20 {"titulo": "Cancion 4", "autor": "Autor 1", "duracion": 1.9},
21
22 ]
23
24
25
26 df = pd.DataFrame(canciones)
27
28
29
30 file = path + "/canciones.csv"
31
32
33
34 df.to_csv(file, index=False)
35
36
37
38 print("Archivo escrito exitosamente")
39
```