



Bootcamp Data Engineering 2

#Week 2

Ejercicio 1: **CSV**

Archivo ventas_home.csv

```
ventas_home.csv X
data > ventas_home.csv > data
1 producto,categoria,precio,cantidad
2 Manzana ,fruta,100,5.0
3 Banana,Fruta,80,
4 Pera,frutas,90,10.0
5 Notebook,Electronica,50000,2.0
6 Auriculares,electrónica,1500,4.0
7 Remera,indumentaria,2000,3.0
8 Pantalón,Indumentaria,3500,1.0
9 Campera, indumentaria ,8000,2.0
10 Tablet,electrónica,30000,1.0
11 Teclado,Electronica,2500,
12 Camisaco,indumentaria,2000,3.0
13 camisaco,indumentaria,3000,1
14 Jeans,indumentaria,4000,
15 Tablet,electrónica,30000,15
16
17
```

Archivo ventas_home_limpio.csv

```
ventas_home_limpio.csv X
outputs > ventas_home_limpio.csv > data
1 producto, categoria, precio, cantidad, ingresos
2 Manzana, Fruta, 100.0, 5, 500.0
3 Banana, Fruta, 80.0, 0, 0.0
4 Pera, Fruta, 90.0, 10, 900.0
5 Notebook, Electrónica, 50000.0, 2, 100000.0
6 Auriculares, Electrónica, 1500.0, 4, 6000.0
7 Remera, Indumentaria, 2000.0, 3, 6000.0
8 Pantalón, Indumentaria, 3500.0, 1, 3500.0
9 Campera, Indumentaria, 8000.0, 2, 16000.0
10 Tablet, Electrónica, 30000.0, 1, 30000.0
11 Teclado, Electrónica, 2500.0, 0, 0.0
12 Camisaco, Indumentaria, 2000.0, 3, 6000.0
13 camisaco, Indumentaria, 3000.0, 1, 3000.0
14 Jeans, Indumentaria, 4000.0, 0, 0.0
15 Tablet, Electrónica, 30000.0, 15, 450000.0
16
```

Archivo resumen_ventas.csv

```
resumen_ventas.csv X
outputs > resumen_ventas.csv > data
1 categoria, ingresos_total, productos_distintos
2 Electrónica, 586000.0, 4
3 Indumentaria, 34500.0, 6
4 Fruta, 1400.0, 3
5
```

Scripts:

```
cargar_csv.py X
cargar_csv.py > ...
1 # cargar_csv.py
2 from pathlib import Path
3 import pandas as pd
4
5 ROOT = Path(__file__).resolve().parent
6 CSV = ROOT / "data" / "ventas_home.csv"
7
8 if not CSV.exists():
9     raise FileNotFoundError(f"No encuentro el archivo: {CSV}")
10
11 # UTF-8 primero; si falla, probamos latin-1 (útil por tildes)
12 try:
13     df = pd.read_csv(CSV, encoding="utf-8")
14 except UnicodeDecodeError:
15     df = pd.read_csv(CSV, encoding="latin-1")
16
17 print("✅ Cargado OK")
18 print("Ruta:", CSV)
19 print("Shape:", df.shape) # (filas, columnas)
20 print("\nPrimeras filas:")
21 print(df.head(10).to_string(index=False))
22 print("\nInfo del DataFrame:")
23 print(df.info())
24
```



```
limpiar.py x
limpiar.py > ...
22         .str.lower()
23         .str.replace(r"\s+", "_", regex=True)
24     )
25
26     # Chequeo mínimo
27     required = {"producto", "categoria", "precio", "cantidad"}
28     faltan = required - set(df.columns)
29     if faltan:
30         raise KeyError(f"Faltan columnas en el CSV: {sorted(faltan)}. Encabezados: {list(df.columns)}")
31
32     # 2.1 Quitar espacios extra en producto
33     df["producto"] = (
34         df["producto"].astype(str)
35         .str.replace(r"\s+", " ", regex=True)
36         .str.strip()
37     )
38
39     # 2.2 Uniformar categoria (Fruta, Electrónica, Indumentaria)
40     def sin_tildes(s: str) -> str:
41         s = str(s)
42         return "".join(c for c in unicodedata.normalize("NFKD", s) if not unicodedata.combining(c))
43
44     cat_norm = (df["categoria"].astype(str)
45                 .str.strip().str.lower()
46                 .apply(sin_tildes))
47
48     mapa = {
49         "fruta": "Fruta",
50         "frutas": "Fruta",
51         "electronica": "Electrónica",
52         "electronic": "Electrónica",
53         "indumentaria": "Indumentaria",
54     }
55     df["categoria"] = cat_norm.apply(lambda x: mapa.get(x, "Otros"))
56
57     # 2.3 Reemplazar precios en texto por números
58     def to_number(x):
59         if pd.isna(x):
60             return pd.NA
61         s = str(x).strip()
62         s = s.replace("$", "").replace("USD", "").replace("ARS", "")
63         s = re.sub(r"[\s](?=\d{3})(?:\D|$)", "", s) # quita miles
64         s = s.replace(",", ".") # coma -> punto
65         try:
66             return float(s)
67         except Exception:
68             return pd.NA
```

```
69
70     # ← línea equivalente a tu 49, usando apply para evitar warnings del linter
71     df["precio"] = df["precio"].apply(to_number)
72
73     # 2.4 Rellenar nulos en cantidad con 0 (asegurar int) ← tu línea 50
74     df["cantidad"] = pd.to_numeric(df["cantidad"], errors="coerce").fillna(0).astype(int)
75
76     # Guardar dataset limpio SIN ingresos (solo Parte 2)
77     df.to_csv(CSV_OUT, index=False, encoding="utf-8")
78     print(f"Archivo limpio (sin ingresos) guardado en: {CSV_OUT}")
79
```

```
columna_ingresos.py X
columna_ingresos.py > ...
1 # columna_ingresos.py
2 from pathlib import Path
3 import pandas as pd
4
5 ROOT = Path(__file__).resolve().parent
6 IN_OUT = ROOT / "outputs" / "ventas_home_limpio.csv" # viene del Paso 2 y se actualiza acá
7
8 if not IN_OUT.exists():
9     raise FileNotFoundError(f"No encuentro el archivo limpio del Paso 2: {IN_OUT}")
10
11 # Cargar robusto
12 try:
13     df = pd.read_csv(IN_OUT, encoding="utf-8")
14 except UnicodeDecodeError:
15     df = pd.read_csv(IN_OUT, encoding="latin-1")
16
17 # Asegurar tipos (por si el Paso 2 no forzó numéricos)
18 df["precio"] = pd.to_numeric(df["precio"], errors="coerce").fillna(0.0)
19 df["cantidad"] = pd.to_numeric(df["cantidad"], errors="coerce").fillna(0).astype(int)
20
21 # PASO 3: crear ingresos = precio * cantidad
22 df["ingresos"] = df["precio"] * df["cantidad"]
23
24 # Guardar (sobrescribe el limpio con la nueva columna)
25 df.to_csv(IN_OUT, index=False, encoding="utf-8")
26
27 print("✅ Paso 3 OK - columna 'ingresos' creada")
28 print(f"Archivo actualizado: {IN_OUT}")
29 print(df.head(10).to_string(index=False))
30
```

```
agrupar_resumen.py X
agrupar_resumen.py > ...
1 # agrupar_resumen.py
2 from pathlib import Path
3 import pandas as pd
4
5 ROOT = Path(__file__).resolve().parent
6 INP = ROOT / "outputs" / "ventas_home_limpio.csv" # viene de los pasos 2 y 3
7 OUT = ROOT / "outputs" / "resumen_ventas.csv"
8
9 if not INP.exists():
10     raise FileNotFoundError(f"Falta el archivo limpio con 'ingresos': {INP}")
11
12 # Cargar
13 try:
14     df = pd.read_csv(INP, encoding="utf-8")
15 except UnicodeDecodeError:
16     df = pd.read_csv(INP, encoding="latin-1")
17
18 # Si por algún motivo no existe 'ingresos', lo calculo en memoria (sin modificar el archivo)
19 if "ingresos" not in df.columns:
20     df["precio"] = pd.to_numeric(df.get("precio"), errors="coerce").fillna(0.0)
21     df["cantidad"] = pd.to_numeric(df.get("cantidad"), errors="coerce").fillna(0).astype(int)
22     df["ingresos"] = df["precio"] * df["cantidad"]
23
24 # Paso 4: agrupar por categoría
25 resumen = (
26     df.groupby("categoria", dropna=False)
27     .agg(
28         ingresos_total=("ingresos", "sum"),
29         productos_distintos=("producto", "nunique")
30     )
31     .reset_index()
32     .sort_values("ingresos_total", ascending=False)
33 )
34
35 # Paso 5: guardar
36 OUT.parent.mkdir(parents=True, exist_ok=True)
37 resumen.to_csv(OUT, index=False, encoding="utf-8")
38 print(f"✅ Resumen generado: {OUT}")
39 print(resumen)
40
```

Ejercicio 2: Diccionarios y JSON

Archivo JSON

```
{} stock.json X
outputs > {} stock.json > ...
1 {
2   "Manzana": 50,
3   "Banana": 40,
4   "Naranja": 25
5 }
```

Scripts:

```
stock_json.py X
stock_json.py > ...
1 # stock_json.py
2 from pathlib import Path
3 import json
4
5 ROOT = Path(__file__).resolve().parent
6 OUT_DIR = ROOT / "outputs"
7 OUT_DIR.mkdir(parents=True, exist_ok=True)
8 JSON_PATH = OUT_DIR / "stock.json"
9
10 # 1) Diccionario inicial
11 stock = {
12     "Manzana": 50,
13     "Banana": 30,
14     "Pera": 20
15 }
16
17 # 2) Agregar producto nuevo
18 stock["Naranja"] = 25
19
20 # 3) Actualizar cantidad existente (ej: +10 bananas)
21 stock["Banana"] = stock.get("Banana", 0) + 10 # 40
22
23 # 4) Eliminar un producto (ej: Pera)
24 stock.pop("Pera", None)
25
26 # 5) Guardar a JSON (pretty)
27 with open(JSON_PATH, "w", encoding="utf-8") as f:
28     json.dump(stock, f, ensure_ascii=False, indent=2)
29
30 # 6) Leer y mostrar contenido
31 with open(JSON_PATH, "r", encoding="utf-8") as f:
32     data = json.load(f)
33
34 print("✅ stock.json generado en:", JSON_PATH)
35 print(json.dumps(data, ensure_ascii=False, indent=2))
36
```

Ejercicio 3: **APIs públicas** *(opcional)*

JSON final

```
... {} todos_completados.json X
outputs > {} todos_completados.json > ...
1  [
2    {
3      "userId": 1,
4      "id": 4,
5      "title": "et porro tempora",
6      "completed": true
7    },
8    {
9      "userId": 1,
10     "id": 8,
11     "title": "quo adipisci enim quam ut ab",
12     "completed": true
13   },
14   {
15     "userId": 1,
16     "id": 10,
17     "title": "illo est ratione doloremque quia maiores aut",
18     "completed": true
19   },
20   {
21     "userId": 1,
22     "id": 11,
23     "title": "vero rerum temporibus dolor",
24     "completed": true
25   },
26   {
27     "userId": 1,
28     "id": 12,
29     "title": "ipsa repellendus fugit nisi",
30     "completed": true
31   },
32   {
33     "userId": 1,
34     "id": 14,
35     "title": "repellendus sunt dolores architecto voluptatum",
36     "completed": true
37   },
38   {
39     "userId": 1,
40     "id": 15,
41     "title": "ab voluptatum amet voluptas",
42     "completed": true
43   },
44   {
45     "userId": 1,
46     "id": 16,
47     "title": "accusamus eos facilis sint et aut voluptatem",
48     "completed": true
49   }
50 ]
```

Scripts:

```
todos_api.py X
todos_api.py > ...
1  # todos_api.py
2  from pathlib import Path
3  import json
4  import requests
5
6  ROOT = Path(__file__).resolve().parent
7  OUT_DIR = ROOT / "outputs"
8  OUT_DIR.mkdir(parents=True, exist_ok=True)
9  OUT_PATH = OUT_DIR / "todos_completados.json"
10
11  URL = "https://jsonplaceholder.typicode.com/todos" # endpoint público
12
13  # 1) Hacer la solicitud HTTP GET
14  resp = requests.get(URL, timeout=30) # timeout por si no responde
15  resp.raise_for_status() # si no es 200 OK, lanza error
16
17  # 2) Convertir la respuesta JSON (texto) a objetos Python (lista/dict)
18  todos = resp.json() # ahora es una lista de dicts
19
20  # 3) Filtrar solo los completados
21  completados = [t for t in todos if t.get("completed") is True]
22
23  # 4) Guardar el resultado en disco como JSON bonito
24  with open(OUT_PATH, "w", encoding="utf-8") as f:
25      json.dump(completados, f, ensure_ascii=False, indent=2)
26
27  print(f"OK → completados: {len(completados)}")
28  print(f"Guardado en: {OUT_PATH}")
29
```