

Domando Datos Reales con Pandas

Limpieza, Análisis y Series Temporales

Curso de IA Aplicada al Agro

11 de enero de 2026

Índice general

1 Capítulo I: El DataFrame	3
1.1 Creación y Carga de Datos	3
1.2 Selección de Datos: loc vs iloc	3
2 Capítulo II: Limpieza de Datos (Data Cleaning)	4
2.1 Detectar Datos Faltantes	4
2.2 Estrategias de Corrección	4
3 Capítulo III: Filtrado y Estadística	5
3.1 Agrupación (GroupBy)	5

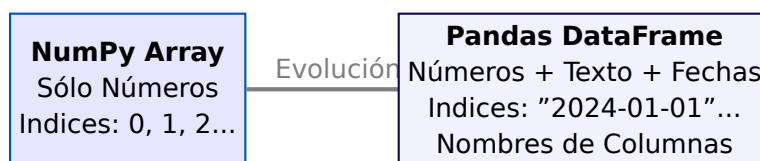
Introducción: ¿Por qué NumPy no es suficiente?

En la Semana 02 aprendimos que NumPy es rapidísimo para matrices numéricas. Pero el mundo real no son solo números; son etiquetas, fechas y datos sucios.

Imagina una tabla de Excel de una estación meteorológica. Tiene:

- *Fechas* (Time Series).
- *Texto* (Nombres de sensores, ubicación).
- *Huecos* (Celdas vacías porque se fue la luz).

Para esto existe **Pandas**. Es una librería construida sobre NumPy que añade etiquetas y manejo de errores.



1 Capítulo I: El DataFrame

El objeto principal de Pandas es el 'DataFrame'. Piensa en él como una hoja de cálculo programable.

1.1 Creación y Carga de Datos

Rara vez creamos DataFrames a mano; usualmente los cargamos desde archivos.

Listing 1: Cargar un CSV de Cosecha

```
1 import pandas as pd
2
3 # Leer archivo separado por comas
4 # index_col=0 usa la primera columna como índice (ID)
5 df = pd.read_csv('cosecha_2024.csv', index_col=0)
6
7 # Ver las primeras 5 filas
8 print(df.head())
9
10 # Ver informacion tecnica (tipos de datos, memoria)
11 print(df.info())
```

1.2 Selección de Datos: loc vs iloc

Esta es la fuente de confusión #1.

- `.loc[]`: Busca por Etiqueta (Nombre).
- `.iloc[]`: Busca por Índice Numérico (Posición).

```
1 # Seleccionar columna 'Humedad' de la fila 'Sensor_A'
2 dato = df.loc['Sensor_A', 'Humedad']
3
4 # Seleccionar fila 0, columna 2
5 dato_raw = df.iloc[0, 2]
```

2 Capítulo II: Limpieza de Datos (Data Cleaning)

En el agro, los datos son ruidosos. Un sensor de humedad puede desconectarse y enviar valores vacíos o basura.

El Peligro del NaN

En computación científica, **NaN** significa *Not a Number*. Si intentas sumar $5 + \text{NaN}$, el resultado es NaN. Un solo dato faltante puede arruinar todo el cálculo de promedio anual.

2.1 Detectar Datos Faltantes

Pandas nos permite ver dónde están los huecos.

```
1 # Crear un reporte de datos faltantes por columna
2 print(df.isna().sum())
3
4 # Salida ejemplo:
5 # Temperatura      0
6 # Humedad          15 <-- 15 lecturas perdidas!
7 # pH               2
```

2.2 Estrategias de Corrección

Como ingenieros, debemos decidir qué hacer con esos huecos.

1. Eliminar (Drop): Si hay pocos datos malos, bórralos. 2. Imputar (Fill): Rellenar con un promedio o el valor anterior.

Listing 2: Sanitización de Sensores

```
1 # ESTRATEGIA 1: Borrar filas con huecos
2 df_limpio = df.dropna()
3
4 # ESTRATEGIA 2: Rellenar (Ideal para series de tiempo)
5 # 'ffill' (Forward Fill) toma el ultimo valor valido y lo repite.
6 # Si a las 10:00 hizo 25C y a las 11:00 es NaN, asume 25C.
7 df_relleno = df.fillna(method='ffill')
8
9 # ESTRATEGIA 3: Rellenar con promedio
10 promedio = df['Humedad'].mean()
11 df['Humedad'] = df['Humedad'].fillna(promedio)
```

3 Capítulo III: Filtrado y Estadística

Podemos hacer preguntas complejas a los datos sin usar bucles.

Caso de Uso: Alerta de Hongos

Los hongos proliferan con humedad alta (> 80%) y temperaturas moderadas (> 20°C). Filtremos esos días.

Listing 3: Filtrado Multicriterio

```
1 # Filtro booleano (Igual que NumPy, pero con nombres)
2 riesgo_hongos = df[ (df['Humedad'] > 80) & (df['Temperatura'] > 20) ]
3
4 print("Días con riesgo de hongos:")
5 print(riesgo_hongos)
```

3.1 Agrupación (GroupBy)

El poder de las tablas dinámicas de Excel, pero en código.

```
1 # Supongamos que tenemos columna 'Zona' (Norte, Sur, Este)
2 # Queremos el promedio de rendimiento por zona.
3
4 resumen = df.groupby('Zona')['Rendimiento'].mean()
5 print(resumen)
```

Reto Semanal: El Dataset Sucio

Se te entregará un archivo `clima_corrupto.csv`.

1. Cárgalo en un DataFrame.
2. Encuentra cuántas filas tienen NaN.
3. Rellena las temperaturas faltantes usando interpolación lineal (`'.interpolate()'`).
4. Calcula la temperatura promedio mensual.
5. Exporta el resultado a `clima_limpio.csv`.

```
1 # Pista para exportar
2 df.to_csv('clima_limpio.csv')
```