

# Manual Avanzado de Pandas

## Procesamiento de Datos para Agroindustria 4.0

Curso de IA Aplicada al Agro

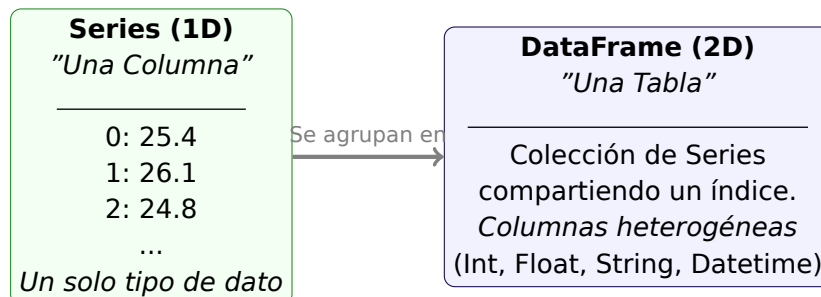
11 de enero de 2026

### Índice general

<b>1 Capítulo I: Ingesta de Datos (I/O)</b>	<b>2</b>
1.1 Lectura Robusta . . . . .	2
1.2 Inspección Inicial . . . . .	2
<b>2 Capítulo II: Indexación y Selección Quirúrgica</b>	<b>3</b>
2.1 Filtrado Booleano (Máscaras) . . . . .	3
<b>3 Capítulo III: Limpieza Profunda (Data Cleaning)</b>	<b>3</b>
3.1 Sanitización de Tipos . . . . .	3
3.2 Tratamiento de Nulos (Imputación) . . . . .	3
<b>4 Capítulo IV: Ingeniería de Características</b>	<b>4</b>
4.1 Operaciones Vectorizadas . . . . .	4
4.2 Lógica Personalizada (.apply) . . . . .	4
<b>5 Capítulo V: Series Temporales y Resampling</b>	<b>4</b>
<b>6 Capítulo VI: Fusión de Datos (Merge)</b>	<b>4</b>
<b>7 Capítulo VII: Ética y Documentación</b>	<b>6</b>
<b>8 Capítulo VIII: Taller de Ejercicios Prácticos</b>	<b>7</b>
8.1 Nivel 1: Laboratorio Guiado (Producción Lechera) . . . . .	7
8.2 Nivel 2: Ejercicios de Calistenia . . . . .	7
8.3 Nivel 3: El Desafío de Integración . . . . .	8

## Introducción: La Estructura de los Datos

Antes de limpiar datos, debemos entender cómo Pandas los organiza en la memoria RAM. A diferencia de Excel, donde todo es una celda visual, Pandas distingue rigurosamente entre tipos de estructuras.



## 1 Capítulo I: Ingesta de Datos (I/O)

El mundo real no vive solo en CSV. En agroindustria, lidiarás con Excel antiguos y JSON de APIs modernas.

### 1.1 Lectura Robusta

El comando `read_csv` tiene más de 50 parámetros. Aquí están los esenciales para evitar errores al cargar.

Listing 1: Carga Avanzada de Datos

```
1 import pandas as pd
2
3 # 1. Cargar CSV definiendo fechas y nulos personalizados
4 df_sensor = pd.read_csv(
5     'sensores_raw.csv',
6     sep=';', # A veces usan punto y coma
7     parse_dates=['fecha_hora'], # Convertir autom. a datetime
8     na_values=['error', '-', 'null'], # Tratar estos textos como NaN
9     dtype={'id_sensor': str} # Forzar ID como texto, no numero
10 )
11
12 # 2. Cargar Excel (requiere libreria 'openpyxl')
13 df_suelos = pd.read_excel(
14     'analisis_suelos.xlsx',
15     sheet_name='Lote_Norte',
16     header=1 # Usar la segunda fila como titulos
17 )
```

### 1.2 Inspección Inicial

Lo primero que haces al recibir un dataset:

```
1 print(df.shape) # (filas, columnas)
2 print(df.dtypes) # ¿Los numeros son float o object (texto)?
3 print(df.head(3)) # Muestra visual rapida
4 print(df['cultivo'].value_counts()) # Conteo de categorias unicas
```

## 2 Capítulo II: Indexación y Selección Quirúrgica

Acceder a los datos incorrectamente es la fuente #1 de *bugs* silenciosos.

Método	¿Qué usa?	Ejemplo
<code>.loc[]</code>	<b>Etiquetas</b> (Nombres)	<code>df.loc["2024-01-01", "Humedad"]</code>
<code>.iloc[]</code>	<b>Posición</b> (Enteros)	<code>df.iloc[0, 3]</code>
<code>df[ ]</code>	<b>Columnas</b> (Principalmente)	<code>df["Temperatura"]</code>

### 2.1 Filtrado Booleano (Máscaras)

No uses `if` dentro de un `for`. Usa máscaras vectorizadas.

Listing 2: Consultas Complejas

```
1 # 1. Crear la mascara (devuelve Series de True/False)
2 mask_calor = df['Temperatura'] > 30
3 mask_seco = df['Humedad'] < 40
4
5 # 2. Aplicar filtro
6 alerta_estres = df[mask_calor & mask_seco]
7
8 # 3. Filtrar por lista de valores permitidos
9 lotes_interes = ['Lote_A1', 'Lote_B5']
10 df_foco = df[df['id_lote'].isin(lotes_interes)]
```

## 3 Capítulo III: Limpieza Profunda (Data Cleaning)

### Tipos de Datos Incorrectos

Si una columna de números tiene un solo texto (ej: "Sin dato"), Pandas convertirá **toda** la columna a `object` (texto). No podrás hacer sumas ni promedios.

### 3.1 Sanitización de Tipos

```
1 # Convertir a numerico, forzando errores a NaN
2 # Si hay un texto "Error", se vuelve NaN en lugar de romper el codigo
3 df['ph_suelo'] = pd.to_numeric(df['ph_suelo'], errors='coerce')
4
5 # Eliminar duplicados (ej: sensor envio el dato 2 veces)
6 df = df.drop_duplicates(subset=['fecha', 'id_sensor'])
```

### 3.2 Tratamiento de Nulos (Imputación)

- `fillna(0)`: Peligroso en agro (0 de lluvia es válido, 0 de pH es ácido mortal).
- `ffill()`: Rellenar con el último valor conocido (ideal sensores).
- `interpolate()`: Trazar una línea entre puntos.

```
1 # Interpolacion lineal limitada a huecos pequeños
2 df['Humedad'] = df['Humedad'].interpolate(method='time', limit=2)
```

## 4 Capítulo IV: Ingeniería de Características

La IA no aprende sola; a veces debemos darle "pistas" creando nuevas columnas.

### 4.1 Operaciones Vectorizadas

```
1 # Crear columna nueva basada en calculo matematico
2 # Diferencia de temperatura (Dia - Noche)
3 df['amplitud_termica'] = df['temp_max'] - df['temp_min']
```

### 4.2 Lógica Personalizada (.apply)

Cuando la matemática simple no basta, aplicamos una función a cada fila. **\*\*Nota:\*\*** Es más lento que vectorizar, úsalo con cuidado.

```
1 def clasificar_riesgo(fila):
2     if fila['Humedad'] > 90 and fila['Temperatura'] > 25:
3         return "ALTO (Hongos)"
4     elif fila['Viento'] > 50:
5         return "ALTO (Daño Fisico)"
6     else:
7         return "BAJO"
8
9 # axis=1 significa "pasar la fila completa a la funcion"
10 df['tipo_riesgo'] = df.apply(clasificar_riesgo, axis=1)
```

## 5 Capítulo V: Series Temporales y Resampling

El agro se mueve por ciclos (diarios, mensuales, estacionales).

### Resampling vs Rolling

- **Resample:** Cambiar la frecuencia (de Hora → Día). Como un "Zoom Out".
- **Rolling:** Ventana deslizante. Suaviza el gráfico manteniendo la frecuencia.

Listing 3: Análisis Temporal

```
1 # Asegurar que el indice es Datetime
2 df = df.set_index('fecha')
3
4 # 1. RESAMPLE: Promedio Mensual ('M')
5 clima_mensual = df.resample('M').mean()
6
7 # 2. ROLLING: Media movil de 7 dias (Suavizar picos)
8 # Elimina el ruido de dias atipicos
9 df['temp_suavizada'] = df['Temperatura'].rolling(window=7).mean()
```

## 6 Capítulo VI: Fusión de Datos (Merge)

Rara vez tienes toda la info en una tabla.

- Tabla A: Cosecha (kilos, fecha, id\_lote)

- Tabla B: Suelo (id\_lote, tipo\_tierra, ph)

```
1 # Unir tablas usando 'id_lote' como llave
2 # how='left' mantiene todas las cosechas, y pega info de suelo si existe
3 df_completo = pd.merge(
4     df_cosecha,
5     df_suelo,
6     on='id_lote',
7     how='left'
8 )
```

## 7 Capítulo VII: Ética y Documentación

La limpieza de datos es un proceso subjetivo que altera la realidad observada.

### Checklist de Integridad de Datos

1. **Trazabilidad:** ¿Guardo el dataset raw original sin tocar?
2. **Transparencia:** ¿Documento cuántas filas eliminé y por qué?
3. **Sesgo de Imputación:** Al rellenar nulos con el promedio, ¿estoy ocultando un fallo sistemático de un sensor en una zona pobre?

## Reto Profesional: Auditoría Climática

Se te entrega estacion\_falla.csv (10,000 filas). Contiene datos horarios de 2023.

### Problemas conocidos:

1. El sensor de lluvia se "traba" y repite el valor exacto por horas.
2. Hay temperaturas de 200°C (error de voltaje).
3. Faltan días completos en Octubre.

### Tu Misión:

```
1 # 1. Cargar parseando fechas
2 df = pd.read_csv('estacion_falla.csv', parse_dates=['timestamp'], index_col='timestamp')
3
4 # 2. Filtrar Outliers (Temp > 50 es imposible)
5 df.loc[df['temp'] > 50, 'temp'] = pd.NA
6
7 # 3. Detectar "valores congelados" (Diff = 0 por mucho tiempo)
8 # shift(1) mueve la columna un paso abajo para comparar con la anterior
9 df['cambio_lluvia'] = df['lluvia'] - df['lluvia'].shift(1)
10
11 # 4. Resamplear a diario sumando lluvia y promediando temp
12 df_diario = df.resample('D').agg({
13     'lluvia': 'sum',
14     'temp': 'mean',
15     'humedad': 'max'
16 })
17
18 # 5. Visualizar comparando Raw vs Limpio
19 import matplotlib.pyplot as plt
20 plt.figure(figsize=(12,6))
21 df['temp'].plot(alpha=0.3, label='Raw (Horario)', color='gray')
22 df_diario['temp'].plot(linewidth=2, label='Limpio (Diario)', color='blue')
23 plt.legend()
24 plt.title("Auditoria de Temperatura 2023")
25 plt.savefig("reporte_calidad.pdf")
```

## 8 Capítulo VIII: Taller de Ejercicios Prácticos

La única forma de aprender Pandas es escribiendo código. En este capítulo, resolveremos problemas reales de la agroindustria.

### 8.1 Nivel 1: Laboratorio Guiado (Producción Lechera)

En este ejercicio, crearemos un pequeño dataset manualmente para entender el flujo completo sin depender de archivos externos.

#### Objetivo

Analizar la calidad de leche de 3 vacas durante una semana y detectar anomalías en el pH (Rango normal: 6.6 - 6.8).

Listing 4: Simulación y Análisis

```
1 import pandas as pd
2 import numpy as np
3
4 # 1. CREAR DATOS SIMULADOS
5 data = {
6     'id_vaca': ['V01', 'V02', 'V01', 'V03', 'V02', 'V03'],
7     'fecha': ['2024-01-01', '2024-01-01', '2024-01-02',
8              '2024-01-01', '2024-01-02', '2024-01-02'],
9     'litros': [20.5, 18.2, 21.0, 19.5, 18.0, 15.0],
10    'ph': [6.7, 6.6, 6.2, 6.8, 6.9, 7.1] # Notar valores acidos/alcalinos
11 }
12
13 df_leche = pd.DataFrame(data)
14
15 # 2. LIMPIEZA: FILTRAR LECHE ACIDA (pH < 6.5)
16 # Esto se descarta para consumo humano
17 leche_mala = df_leche[df_leche['ph'] < 6.5]
18 print("--- ALERTA DE ACIDEZ ---")
19 print(leche_mala)
20
21 # 3. ANALISIS: TOTAL DE LITROS POR VACA (SOLO LECHE BUENA)
22 df_buena = df_leche[df_leche['ph'] >= 6.5]
23 reporte = df_buena.groupby('id_vaca')['litros'].sum()
24
25 print("\n--- REPORTE DE PAGO ---")
26 print(reporte)
```

### 8.2 Nivel 2: Ejercicios de Calistenia

Resuelve estos problemas cortos. Asume que tienes un DataFrame llamado `df_clima`.

- **Ejercicio A (Selección):** Crea un nuevo DataFrame llamado `df_lluvia` que solo contenga las columnas "fecha" y "precipitacion".
- **Ejercicio B (Filtro Simple):** Filtra los días donde la "temperatura" fue mayor a 35 grados.
- **Ejercicio C (Corrección de Texto):** La columna "zona" tiene valores " Norte" (con espacio extra). Usa `df['zona'].str.strip()` para arreglarlo.

- **Ejercicio D (Creación de Columnas):** Convierte la temperatura de Celsius a Fahrenheit en una nueva columna:  $F = C \times 1.8 + 32$ .

### 8.3 Nivel 3: El Desafío de Integración

#### Escenario Real

Tienes dos fuentes de datos no sincronizadas:

1. `ventas.csv`: Registro de sacos de café vendidos (con fechas).
2. `precio_internacional.csv`: Precio del café en bolsa (diario).

**Problema:** El archivo de ventas no tiene el precio. Debes cruzar los datos para saber cuánto dinero real ingresó.

#### Pasos para la solución:

```
1 # PISTA PARA LA SOLUCION
2
3 # 1. Cargar ambos DataFrames
4 ventas = pd.read_csv('ventas.csv')
5 precios = pd.read_csv('precio_internacional.csv')
6
7 # 2. Asegurar que las fechas sean tipo 'datetime' en AMBOS
8 ventas['fecha'] = pd.to_datetime(ventas['fecha'])
9 precios['fecha'] = pd.to_datetime(precios['fecha'])
10
11 # 3. Realizar el MERGE (Cruce)
12 # Usamos 'left' para no perder ventas si falta algun precio
13 df_final = pd.merge(ventas, precios, on='fecha', how='left')
14
15 # 4. Calcular ingreso total
16 # Si faltan precios (NaN), rellenar con el promedio antes de calcular
17 precio_promedio = precios['valor'].mean()
18 df_final['valor'] = df_final['valor'].fillna(precio_promedio)
19
20 df_final['ingreso_total'] = df_final['kilos_vendidos'] * df_final['valor']
```

#### ¿Por qué fallan los Merges?

El error más común es que la columna clave (fecha) tenga formatos distintos.

- Tabla A: "2024-01-01" (String)
- Tabla B: 2024-01-01 (Timestamp)

**Regla de Oro:** Siempre normaliza tipos con `.info()` antes de unir tablas.