

Pandas para Procesamiento Industrial

Trazabilidad, Control de Calidad y Análisis de Producción
Agroindustria Alimenticia 4.0

Curso de IA Aplicada al Agro
Universidad del Valle

Enero 2026

Resumen

Este manual introduce Pandas desde la perspectiva de la ingeniería de datos aplicada a la industria alimenticia. A diferencia del enfoque tradicional basado en análisis exploratorio genérico, aquí abordamos problemas reales de trazabilidad de lotes, control estadístico de procesos (SPC), cumplimiento normativo (HACCP, FDA) y optimización de líneas de producción. Los estudiantes aprenderán a procesar datasets heterogéneos (fechas, categorías, mediciones numéricas) con eficiencia computacional y rigor científico.

Índice general

Prefacio: Del Campo a la Mesa

En la Semana 02 trabajaste con NumPy procesando matrices numéricas homogéneas (humedad del suelo en 365 días × 100 zonas). Ese enfoque funciona para sensores agrícolas, pero **la agroindustria moderna genera datos más complejos**:

- **Trazabilidad**: Cada lote de café procesado tiene ID (string), timestamp de entrada/-salida, temperatura de tostado (float), operario responsable (categoría), resultado QA (booleano).
- **Series temporales irregulares**: Sensores IoT envían datos cada 30 segundos, pero fallan aleatoriamente.
- **Relaciones entre tablas**: Para rastrear un recall de producto, debes cruzar 3 datasets: `lotes_producidos`, `pruebas_laboratorio`, `despachos_clientes`.

NumPy no está diseñado para esto. **Pandas sí.**

Contexto Industrial

Imagina una planta procesadora de alimentos que opera 24/7 en 3 turnos, con 5 líneas de producción y 1200 lotes/mes. Cada lote genera:

- 8 variables de proceso (temperatura, presión, pH, humedad, tiempo)
- 12 pruebas de laboratorio (microbiología, físico-químicas)
- Metadatos de trazabilidad (proveedor, lote de materia prima, destino)

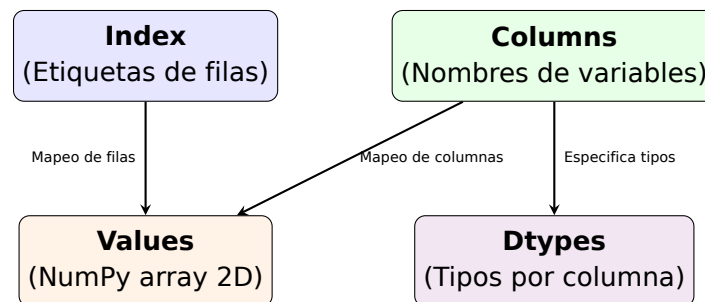
Total: 14,400 lotes/año × 20 variables = 288,000 datos/año.

No puedes analizar esto en Excel. Necesitas código profesional.

1 Capítulo I: Fundamentos — DataFrame como Base de Datos en Memoria

1.1 La Anatomía de un DataFrame

Un DataFrame es una **tabla en memoria RAM** con índice explícito y columnas etiquetadas. A diferencia de NumPy (donde accedes por posición), Pandas permite consultas tipo SQL.



Diferencia clave con NumPy:

- NumPy: `array[0, 3]` → Posición absoluta (fila 0, columna 3)
- Pandas: `df.loc["2026-01-01", "Temperatura"]` → Etiqueta semántica

1.2 Series vs DataFrame

Característica	Series	DataFrame
Dimensionalidad	1D (columna única)	2D (tabla)
Tipo de datos	Homogéneo (un dtype)	Heterogéneo (dtype por columna)
Index	Sí	Sí
Operaciones	Vectorizadas	Por columna/fila
Uso típico	Una medición	Dataset completo

Cuadro 1: Comparación Series-DataFrame

Listing 1: Crear Series y DataFrame desde código

```

1 import pandas as pd
2 import numpy as np
3
4 # Series: Una columna de temperaturas
5 temps = pd.Series([72.5, 73.1, 72.8, 74.0],
6                   index=['Lote_A', 'Lote_B', 'Lote_C', 'Lote_D'],
7                   name='Temperatura_Pasteurizacion')
8
9 print(temps['Lote_B']) # Acceso por etiqueta → 73.1
10
11 # DataFrame: Tabla completa de un turno
12 data = {
13     'id_lote': ['L001', 'L002', 'L003'],
14     'temp_C': [72.5, 73.1, 71.9],
15     'presion_bar': [2.8, 2.9, 2.7],
  
```

```
16     'resultado_QA': ['Aprobado', 'Aprobado', 'Rechazado']
17 }
18
19 df = pd.DataFrame(data)
20 print(df.dtypes)
```

1.3 Carga de Datos Industriales

En la industria, los datos vienen de múltiples fuentes:

- **SCADA** (sistemas de control): CSV/Excel con timestamps
- **LIMS** (laboratorio): Resultados en archivos Excel con hojas múltiples
- **ERP** (SAP/Oracle): Exportaciones CSV con separadores raros
- **Sensores IoT**: JSON desde APIs REST

Listing 2: Carga robusta de datos industriales

```
1 import pandas as pd
2
3 # 1. CSV con problemas comunes
4 df_scada = pd.read_csv(
5     'datos_scada.csv',
6     sep=';', # Separador europeo
7     decimal=',', # Decimales con coma
8     encoding='latin1', # Codificación Windows
9     parse_dates=['timestamp'], # Convertir a datetime automáticamente
10    na_values=['error', 'offline', '-'], # Valores nulos personalizados
11    dtype={'id_lote': str} # Forzar ID como texto (evita 001 → 1)
12 )
13
14 # 2. Excel con múltiples hojas
15 df_lab = pd.read_excel(
16     'resultados_laboratorio.xlsx',
17     sheet_name='Microbiologia', # Hoja específica
18     header=2, # La fila 3 tiene los títulos
19     usecols='A:F' # Solo columnas A-F
20 )
21
22 # 3. JSON desde API de sensor IoT
23 import requests
24 response = requests.get('https://api.sensores.com/temperatura')
25 df_temp = pd.DataFrame(response.json()['data'])
```

2 Capítulo II: Indexación y Selección — El Fundamento de Todo

2.1 Los 3 Métodos de Acceso

⚠ Error #1 más común en Pandas

Confundir `.loc[]` (etiquetas) con `.iloc[]` (posiciones). Esto causa bugs silenciosos cuando el índice no es secuencial.

Método	Qué usa	Ejemplo industrial
<code>.loc[]</code>	Etiquetas (labels)	<code>df.loc["2026-01-15", "pH"]</code>
<code>.iloc[]</code>	Posición (enteros)	<code>df.iloc[0, 3]</code> (primera fila, cuarta columna)
<code>df[]</code>	Columnas (principalmente)	<code>df["Temperatura"]</code>

Cuadro 2: Métodos de indexación en Pandas

Listing 3: Ejemplos de indexación

```

1 import pandas as pd
2
3 # Dataset simulado: Control de calidad de leche
4 data = {
5     'id_lote': ['L001', 'L002', 'L003', 'L004'],
6     'fecha': ['2026-01-10', '2026-01-10', '2026-01-11', '2026-01-11'],
7     'temp_pasteurizacion': [72.5, 73.0, 71.8, 74.2],
8     'ph': [6.7, 6.6, 6.5, 6.9],
9     'resultado': ['Aprobado', 'Aprobado', 'Rechazado', 'Aprobado']
10 }
11
12 df = pd.DataFrame(data)
13
14 # 1. SELECCIÓN DE COLUMNAS
15 temps = df['temp_pasteurizacion'] # Retorna Series
16 subset = df[['id_lote', 'ph']]    # Retorna DataFrame (nota el [[ ]])
17
18 # 2. SELECCIÓN POR ETIQUETA (.loc)
19 # Sintaxis: df.loc[filas, columnas]
20 primera_fila = df.loc[0]          # Primera fila completa
21 ph_L002 = df.loc[1, 'ph']         # Celda específica: 6.6
22 rango = df.loc[0:2, 'temp_pasteurizacion'] # Filas 0-2, una columna
23
24 # 3. SELECCIÓN POR POSICIÓN (.iloc)
25 primera_celda = df.iloc[0, 0]     # 'L001'
26 subcuadro = df.iloc[0:2, 1:3]     # 2 filas x 2 columnas

```

2.2 Filtrado Booleano (Máscaras)

El poder real de Pandas está en las **consultas vectorizadas**. No uses bucles for — usa máscaras booleanas.

Listing 4: Filtrado avanzado para control de calidad

```

1 import pandas as pd
2
3 # Cargar datos de producción
4 df = pd.read_csv('produccion_cafe_enero.csv')

```

```
5
6 # 1. CONSULTA SIMPLE: Lotes rechazados
7 rechazados = df[df['resultado'] == 'Rechazado']
8
9 # 2. CONSULTAS COMPUESTAS: Temperatura fuera de spec Y presión baja
10 # Rango de pasteurización: 72-76°C, Presión mínima: 2.5 bar
11 problemas_criticos = df[
12     ((df['temp_C'] < 72) | (df['temp_C'] > 76)) &
13     (df['presion_bar'] < 2.5)
14 ]
15
16 # 3. FILTRO POR LISTA (isin): Solo líneas L1 y L3
17 lineas_foco = df[df['linea'].isin(['L1', 'L3'])]
18
19 # 4. FILTRO POR STRING (contiene): Lotes de turno nocturno
20 nocturnos = df[df['id_lote'].str.contains('NOCHE')]
21
22 # 5. QUERY (sintaxis SQL-like)
23 # Nota: Solo funciona si nombres de columnas no tienen espacios
24 criticos = df.query('temp_C > 76 and resultado == "Rechazado"')
```

▣ Complejidad Computacional de Máscaras

Una máscara booleana `df['temp'] > 72` tiene complejidad $O(n)$ donde n es el número de filas. Internamente:

1. Pandas delega la comparación a NumPy (código C optimizado)
2. Se crea un array booleano en memoria del mismo tamaño que la columna
3. El filtrado `df[mask]` usa fancy indexing de NumPy

Para un DataFrame de 1M filas, esto toma ~ 10 ms. Un bucle for equivalente tomaría ~ 2 segundos (200x más lento).

3 Capítulo III: Limpieza de Datos — Fail Fast en Producción

3.1 El Problema de los Tipos Incorrectos

⚠ Tipo object

Si una columna numérica aparece como dtype: object, significa que Pandas la leyó como texto. No podrás hacer operaciones matemáticas hasta convertirla.

Causas comunes:

- Un solo valor con texto ("Error", "N/A", "-") contamina toda la columna
- Formato de número europeo: "3,14" en lugar de "3.14"
- Espacios en blanco: " 25.5 " no se convierte automáticamente

Listing 5: Diagnóstico y corrección de tipos

```
1 import pandas as pd
2
3 df = pd.read_csv('sensores_planta.csv')
4
5 # 1. DIAGNÓSTICO
6 print(df.dtypes)
7 print(df.info()) # Muestra tipos y valores no-nulos
8
9 # Ejemplo de salida problemática:
10 # temperatura    object ←⚠ Debería ser float64
11 # presion        object ←⚠ Debería ser float64
12
13 # 2. INSPECCIÓN MANUAL
14 print(df['temperatura'].unique()) # Ver valores únicos
15 # Output: ['25.5', '26.1', 'Error', '24.8', ...] ← "Error" causa el problema
16
17 # 3. CONVERSIÓN FORZADA (errores → NaN)
18 df['temperatura'] = pd.to_numeric(df['temperatura'], errors='coerce')
19 df['presion'] = pd.to_numeric(df['presion'], errors='coerce')
20
21 # 4. VERIFICACIÓN
22 print(df.dtypes)
23 # temperatura    float64 ✓ Corregido
24 # presion        float64 ✓ Corregido
25
26 print(df['temperatura'].isna().sum()) # Contar cuántos NaN se generaron
```

3.2 Tratamiento de Valores Faltantes

En la industria alimenticia, **un dato faltante puede significar un fallo crítico**. No siempre es correcto rellenar con el promedio.

Listing 6: Imputación contextual para sensores

```
1 import pandas as pd
2
3 df = pd.read_csv('temperatura_camara_fria.csv', parse_dates=['timestamp'])
4 df = df.set_index('timestamp')
5
```

Método	Cuándo usarlo	Riesgo
fillna(0)	Contadores (eventos)	0 puede ser válido en agro
ffill()	Series temporales (sensores)	Oculto fallos prolongados
interpolate()	Datos continuos (temperatura)	Inventa datos inexistentes
dropna()	QA crítico	Pierdes información

Cuadro 3: Estrategias de imputación de datos faltantes

```

6 # CASO 1: Interpolación limitada (máximo 2 valores consecutivos)
7 # Si faltan >2 valores, algo falló y no deberíamos inventar datos
8 df['temp'] = df['temp'].interpolate(method='time', limit=2)
9
10 # CASO 2: Forward fill con límite temporal
11 # Rellenar con el último valor conocido, pero solo por 10 minutos
12 df['humedad'] = df['humedad'].fillna(method='ffill', limit=20) # 20 registros = 10 min
13
14 # CASO 3: Marcar como fallo en lugar de imputar
15 df['sensor_falla'] = df['temp'].isna() # Columna booleana de alertas
16
17 # CASO 4: Eliminar filas con datos críticos faltantes
18 df_limpio = df.dropna(subset=['ph', 'acidez']) # Solo si faltan variables críticas

```

3.3 Detección de Outliers

Listing 7: Detección estadística de anomalías

```

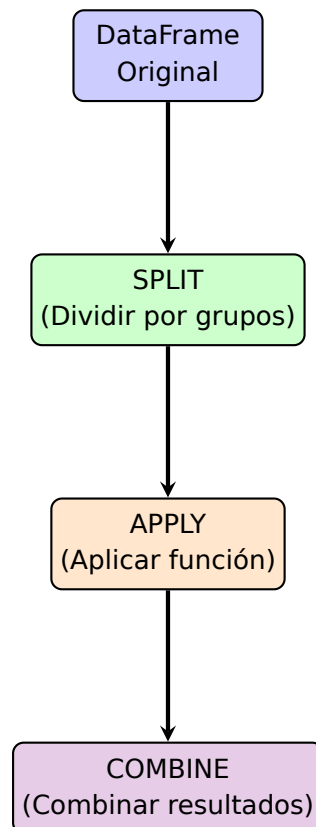
1 import pandas as pd
2 import numpy as np
3
4 df = pd.read_csv('temperatura_pasteurizacion.csv')
5
6 # MÉTODO 1: Rango intercuartílico (IQR) – Robusto a valores extremos
7 Q1 = df['temp'].quantile(0.25)
8 Q3 = df['temp'].quantile(0.75)
9 IQR = Q3 - Q1
10
11 limite_inferior = Q1 - 1.5 * IQR
12 limite_superior = Q3 + 1.5 * IQR
13
14 outliers = df[(df['temp'] < limite_inferior) | (df['temp'] > limite_superior)]
15 print(f"Detectados {len(outliers)} outliers")
16
17 # MÉTODO 2: Z-score (asume distribución normal)
18 mean = df['temp'].mean()
19 std = df['temp'].std()
20 df['z_score'] = (df['temp'] - mean) / std
21
22 # Outliers: |z| > 3 (regla de 3 sigmas)
23 outliers_zscore = df[np.abs(df['z_score']) > 3]
24
25 # MÉTODO 3: Límites físicos (conocimiento del dominio)
26 # La temperatura de pasteurización NUNCA puede ser > 100°C
27 errores_sensor = df[df['temp'] > 100]
28 df.loc[df['temp'] > 100, 'temp'] = np.nan # Marcar como faltante

```


4 Capítulo IV: GroupBy — El Motor de Agregación

4.1 El Paradigma Split-Apply-Combine

`.groupby()` es la operación más importante en Pandas. Implementa el patrón *split-apply-combine*:



Listing 8: Análisis de productividad por línea

```
1 import pandas as pd
2
3 # Dataset: 2000 lotes de café procesados en enero
4 df = pd.read_csv('produccion_cafe_enero.csv', parse_dates=['timestamp_inicio', 'timestamp_fin'])
5
6 # Calcular duración de cada lote
7 df['duracion_min'] = (df['timestamp_fin'] - df['timestamp_inicio']).dt.total_seconds() / 60
8
9 # AGREGACIÓN 1: Productividad por línea
10 productividad = df.groupby('linea').agg({
11     'kg_procesados': 'sum',          # Total de kilos
12     'duracion_min': 'mean',         # Duración promedio
13     'id_lote': 'count'              # Cantidad de lotes
14 })
15
16 print(productividad)
17 # Output:
18 #      kg_procesados  duracion_min  id_lote
19 # linea
20 # L1              45000          87.2     650
21 # L2              38000          92.1     520
22 # L3              42000          89.5     600
23
```

```

24 # AGREGACIÓN 2: Rechazos por turno
25 rechazos = df.groupby(['turno', 'resultado']).size().unstack(fill_value=0)
26 print(rechazos)
27
28 # AGREGACIÓN 3: Múltiples estadísticas
29 stats = df.groupby('linea')['duracion_min'].agg(['mean', 'std', 'min', 'max'])

```

4.2 GroupBy con Transformaciones

A veces no quieres reducir el DataFrame, sino **agregar columnas calculadas por grupo**.

Listing 9: Normalización por grupo

```

1 import pandas as pd
2
3 df = pd.read_csv('lotes_produccion.csv')
4
5 # CASO 1: Calcular % de productividad de cada lote respecto a su línea
6 df['kg_promedio_linea'] = df.groupby('linea')['kg_procesados'].transform('mean')
7 df['performance_relativo'] = (df['kg_procesados'] / df['kg_promedio_linea']) * 100
8
9 # CASO 2: Ranking dentro de cada turno
10 df['ranking_turno'] = df.groupby('turno')['kg_procesados'].rank(ascending=False)
11
12 # CASO 3: Detectar lotes atípicos (> 2 std de su grupo)
13 df['media_linea'] = df.groupby('linea')['duracion_min'].transform('mean')
14 df['std_linea'] = df.groupby('linea')['duracion_min'].transform('std')
15 df['es_atipico'] = (df['duracion_min'] - df['media_linea']).abs() > (2 * df['std_linea'])

```

Complejidad de GroupBy

Internamente, `.groupby()` usa un algoritmo de hashing para agrupar filas:

1. Calcula hash de cada valor en la columna de agrupación: $O(n)$
2. Ordena los índices por hash: $O(n \log n)$
3. Aplica función a cada grupo: $O(n)$

Complejidad total: $O(n \log n)$. Para 1M filas, esto toma ~ 100 ms en un CPU moderno.

Comparación: Un bucle manual con diccionarios tomaría ~ 5 segundos (50x más lento).

5 Capítulo V: Series Temporales — El Corazón de la Industria

5.1 Datetime como Index

En la industria, **el tiempo es el índice natural** de los datos. Convertir el DataFrame a índice temporal desbloquea operaciones avanzadas.

Listing 10: Configurar índice temporal

```
1 import pandas as pd
2
3 # Cargar datos de sensor con timestamps
4 df = pd.read_csv('temperatura_camara.csv')
5
6 # PASO 1: Convertir columna a datetime
7 df['timestamp'] = pd.to_datetime(df['timestamp'])
8
9 # PASO 2: Establecer como índice
10 df = df.set_index('timestamp')
11
12 # PASO 3: Ordenar por tiempo (¡importante!)
13 df = df.sort_index()
14
15 # Ahora puedes hacer selección por rangos de fecha:
16 enero = df['2026-01-01':'2026-01-31']
17 primera_semana = df['2026-01-01':'2026-01-07']
```

5.2 Resampling — Cambiar la Frecuencia

Resampling vs Rolling

- **Resample:** Cambia la frecuencia temporal. Ejemplo: datos cada 30 seg → promedio diario.
- **Rolling:** Ventana deslizante. Mantiene la frecuencia original pero suaviza con promedios móviles.

Listing 11: Resampling para reportes diarios

```
1 import pandas as pd
2
3 # Datos de temperatura cada 30 segundos
4 df = pd.read_csv('temp_pasteurizacion.csv', parse_dates=['timestamp'], index_col='timestamp')
5
6 # RESAMPLE 1: Promedio diario
7 temp_diaria = df['temperatura'].resample('D').mean()
8
9 # RESAMPLE 2: Máximo por hora
10 temp_horaria_max = df['temperatura'].resample('H').max()
11
12 # RESAMPLE 3: Múltiples agregaciones
13 stats_diarias = df.resample('D').agg({
14     'temperatura': ['mean', 'min', 'max', 'std'],
15     'presion': 'mean'
16 })
17
18 # RESAMPLE 4: Contar eventos por turno (8 horas)
19 eventos_turno = df.resample('8H').count()
```

5.3 Rolling Windows — Suavizar Ruido

Listing 12: Ventanas móviles para control de procesos

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 df = pd.read_csv('temperatura_real_time.csv', parse_dates=['timestamp'], index_col='timestamp')
5
6 # Media móvil de 10 minutos (window=20 si los datos son cada 30 seg)
7 df['temp_suavizada'] = df['temperatura'].rolling(window=20).mean()
8
9 # Desviación estándar móvil (detectar variabilidad)
10 df['temp_std_movil'] = df['temperatura'].rolling(window=20).std()
11
12 # Detectar derivas: si la std móvil supera 2°C, el proceso está inestable
13 df['proceso_inestable'] = df['temp_std_movil'] > 2.0
14
15 # Visualización
16 plt.figure(figsize=(12, 6))
17 plt.plot(df.index, df['temperatura'], alpha=0.3, label='Datos crudos')
18 plt.plot(df.index, df['temp_suavizada'], linewidth=2, label='Media móvil 10 min')
19 plt.legend()
20 plt.title('Control de Temperatura - Pasteurización')
21 plt.savefig('control_temperatura.png', dpi=150)
```

6 Capítulo VI: Merge y Trazabilidad — Conectar las Piezas

6.1 El Problema de la Trazabilidad

En agroindustria alimenticia, la trazabilidad es **un requisito legal** (HACCP, ISO 22000, FDA). Debes poder responder:

- ¿Qué clientes recibieron lotes de un proveedor contaminado?
- ¿Qué lotes fueron procesados por un operario específico en una fecha?
- ¿Qué materia prima se usó en un lote con defecto?

Esto requiere **cruzar múltiples tablas**.

6.2 Tipos de Merge

Tipo	Comportamiento
inner	Solo filas con match en ambas tablas (intersección)
left	Todas las filas de la tabla izquierda + matches de la derecha
right	Todas las filas de la tabla derecha + matches de la izquierda
outer	Todas las filas de ambas tablas (unión)

Cuadro 4: Tipos de merge en Pandas

Listing 13: Caso HACCP: Rastreo de lote contaminado

```

1 import pandas as pd
2
3 # PASO 1: Cargar las 3 tablas
4 lotes = pd.read_csv('lotes_producidos.csv')
5 # Columnas: id_lote, fecha_produccion, kg, linea
6
7 pruebas = pd.read_csv('pruebas_laboratorio.csv')
8 # Columnas: id_lote, ph, acidez, resultado
9
10 despachos = pd.read_csv('despachos.csv')
11 # Columnas: id_lote, cliente, fecha_envio, destino
12
13 # PASO 2: Primera unión (lotes + pruebas)
14 lotes_con_qa = pd.merge(lotes, pruebas, on='id_lote', how='left')
15
16 # PASO 3: Segunda unión (agregar despachos)
17 trazabilidad_completa = pd.merge(lotes_con_qa, despachos, on='id_lote', how='left')
18
19 # PASO 4: Identificar lotes problemáticos (pH < 4.3)
20 lotes_problematicos = trazabilidad_completa[trazabilidad_completa['ph'] < 4.3]
21
22 # PASO 5: Listar clientes afectados
23 clientes_afectados = lotes_problematicos[['id_lote', 'cliente', 'destino', 'fecha_envio']]
24 print(clientes_afectados)
25

```

```
26 # PASO 6: Exportar para reporte de recall
27 clientes_afectados.to_csv('recall_lista_clientes.csv', index=False)
```

⚠ Error Común: Claves con tipos diferentes

Si `lotes['id_lote']` es string y `pruebas['id_lote']` es int, el merge fallará silenciosamente (0 matches).

Solución: Siempre verificar tipos antes de merge:

```
1 print(lotes['id_lote'].dtype)
2 print(pruebas['id_lote'].dtype)
3
4 # Si son diferentes, convertir:
5 lotes['id_lote'] = lotes['id_lote'].astype(str)
6 pruebas['id_lote'] = pruebas['id_lote'].astype(str)
```

7 Capítulo VII: Ingeniería de Características

7.1 Creación de Columnas Derivadas

Listing 14: Variables calculadas para análisis

```
1 import pandas as pd
2
3 df = pd.read_csv('produccion_diaria.csv', parse_dates=['fecha'])
4
5 # 1. Duración de proceso (timedelta a minutos)
6 df['duracion_min'] = (df['hora_fin'] - df['hora_inicio']).dt.total_seconds() / 60
7
8 # 2. Rendimiento (kg/hora)
9 df['rendimiento'] = df['kg_producidos'] / (df['duracion_min'] / 60)
10
11 # 3. Categorización de turnos
12 def clasificar_turno(hora):
13     if 6 <= hora < 14:
14         return 'Mañana'
15     elif 14 <= hora < 22:
16         return 'Tarde'
17     else:
18         return 'Noche'
19
20 df['turno'] = df['hora_inicio'].dt.hour.apply(clasificar_turno)
21
22 # 4. Día de la semana (útil para detectar patrones)
23 df['dia_semana'] = df['fecha'].dt.day_name()
24
25 # 5. Semana del año (agrupación temporal)
26 df['semana'] = df['fecha'].dt.isocalendar().week
```

7.2 Discretización (Binning)

Listing 15: Clasificar variables continuas en categorías

```
1 import pandas as pd
2
3 df = pd.read_csv('analisis_ph.csv')
4
5 # Clasificar pH en categorías
6 bins = [0, 6.5, 7.0, 14]
7 labels = ['Ácido', 'Neutro', 'Alcalino']
8 df['categoria_ph'] = pd.cut(df['ph'], bins=bins, labels=labels)
9
10 # Clasificar temperatura de pasteurización en zonas
11 bins_temp = [0, 71, 74, 77, 100]
12 labels_temp = ['Bajo Spec', 'Óptimo Bajo', 'Óptimo Alto', 'Sobre Spec']
13 df['zona_temp'] = pd.cut(df['temp'], bins=bins_temp, labels=labels_temp)
```

8 Capítulo VIII: Ética y Calidad de Datos

⚖ Responsabilidad en la Limpieza de Datos

Cada decisión de limpieza altera la realidad registrada. En la industria alimenticia, esto tiene implicaciones legales y de salud pública.

Principios éticos:

1. **Trazabilidad:** Guardar dataset original sin modificar (raw/)
2. **Documentación:** Registrar cada transformación en un log
3. **Transparencia:** Reportar cuántas filas se eliminaron y por qué
4. **Sesgo de imputación:** No ocultar fallos sistemáticos rellenando con promedios

8.1 Pipeline de Limpieza Documentado

Listing 16: Pipeline con logging

```
1 import pandas as pd
2 import logging
3
4 # Configurar logging
5 logging.basicConfig(filename='limpieza.log', level=logging.INFO)
6
7 def limpiar_dataset(path_entrada, path_salida):
8     # Cargar datos crudos
9     df = pd.read_csv(path_entrada)
10    filas_originales = len(df)
11    logging.info(f"Dataset cargado: {filas_originales} filas")
12
13    # 1. Eliminar duplicados
14    df = df.drop_duplicates()
15    duplicados = filas_originales - len(df)
16    logging.info(f"Duplicados eliminados: {duplicados}")
17
18    # 2. Convertir tipos
19    df['temperatura'] = pd.to_numeric(df['temperatura'], errors='coerce')
20    nulos_generados = df['temperatura'].isna().sum()
21    logging.info(f"Valores no numéricos convertidos a NaN: {nulos_generados}")
22
23    # 3. Eliminar outliers
24    Q1 = df['temperatura'].quantile(0.25)
25    Q3 = df['temperatura'].quantile(0.75)
26    IQR = Q3 - Q1
27    df_limpio = df[
28        (df['temperatura'] >= Q1 - 1.5*IQR) &
29        (df['temperatura'] <= Q3 + 1.5*IQR)
30    ]
31    outliers = len(df) - len(df_limpio)
32    logging.info(f"Outliers eliminados: {outliers}")
33
34    # Guardar dataset limpio
35    df_limpio.to_csv(path_salida, index=False)
36    logging.info(f"Dataset final: {len(df_limpio)} filas guardadas en {path_salida}")
37
38    return df_limpio
39
```



```
40 | # Ejecutar
41 | df_limpio = limpiar_dataset('data/raw/sensores.csv', 'data/processed/sensores_clean.csv')
```

9 Capítulo IX: Talleres Prácticos

9.1 Taller 1: Análisis de Línea de Producción

□ Caso: Planta de Procesamiento de Café

Tienes 2000 lotes procesados en enero en 3 líneas (L1, L2, L3). Debes analizar productividad, identificar cuellos de botella y generar reporte ejecutivo.

Dataset: produccion_cafe_enero.csv

Tareas:

1. Calcular duración promedio por línea
2. Identificar el turno más lento
3. Detectar lotes con duración > 2 desviaciones estándar
4. Generar tabla resumen con productividad (kg/hora)

9.2 Taller 2: Control de Calidad Temporal

□ Caso: Pasteurización de Leche

Sensor registra temperatura cada 30 segundos durante una semana. Debes detectar derivas, generar alertas y producir gráficos de control.

Dataset: temperatura_pasteurizacion_semana.csv

Tareas:

1. Resamplear a promedios de 10 minutos
2. Calcular desviación estándar móvil (ventana 20 lecturas)
3. Detectar periodos con temperatura fuera de $[72-76^{\circ}\text{C}]$ por >5 minutos
4. Visualizar con matplotlib

9.3 Taller 3: Trazabilidad y Recall

□ Caso HACCP: Lote Contaminado

Se detectó contaminación microbiológica en el lote L20260115_042. Debes rastrear qué clientes lo recibieron y generar lista para recall.

Datasets:

- lotes_producidos.csv
- pruebas_microbiologia.csv
- despachos_clientes.csv

Tareas:

1. Merge de las 3 tablas por `id_lote`
2. Filtrar lotes con resultado "Contaminado"
3. Generar CSV con: cliente, dirección, fecha_envío, kg_afectados
4. Crear reporte LaTeX con tabla de afectados

10 Capítulo X: Visualización Profesional con Matplotlib

10.1 Gráficos de Control Estadístico (SPC)

Los gráficos SPC son fundamentales en industria alimenticia para detectar derivas de proceso.

Listing 17: Gráfico de control

```
1 import matplotlib.pyplot as plt
2 import pandas as pd
3
4 df = pd.read_csv('ph_lotes.csv')
5 media = df['ph'].mean()
6 std = df['ph'].std()
7 UCL = media + 3*std
8 LCL = media - 3*std
9
10 fig, ax = plt.subplots(figsize=(14, 6))
11 ax.scatter(df.index, df['ph'], alpha=0.7)
12 ax.axhline(y=media, color='green', linewidth=2, label='Media')
13 ax.axhline(y=UCL, color='red', linestyle='--', label='UCL')
14 ax.axhline(y=LCL, color='red', linestyle='--', label='LCL')
15 ax.legend()
16 plt.savefig('spc.png', dpi=300)
```

11 Capítulo XI: Estadística con SciPy

11.1 Pruebas de Hipótesis

Listing 18: t-test entre turnos

```
1 from scipy import stats
2 import pandas as pd
3
4 df = pd.read_csv('defectos.csv')
5 diurno = df[df['turno'] == 'Diurno']['defectos']
6 nocturno = df[df['turno'] == 'Nocturno']['defectos']
7
8 t_stat, p_value = stats.ttest_ind(diurno, nocturno)
9
10 if p_value < 0.05:
11     print("Diferencia significativa entre turnos")
```

12 Capítulo XII: SQL para Ciencia de Datos

12.1 Conexión y Queries

Listing 19: Pandas + SQL

```
1 import pandas as pd
2 from sqlalchemy import create_engine
3
4 engine = create_engine('postgresql://user:pass@localhost/planta')
5
6 df = pd.read_sql_query(
7     "SELECT linea, AVG(kg) FROM lotes GROUP BY linea",
8     engine
9 )
```

Referencias y Recursos

Bibliografía Científica

1. McKinney, W. (2017). *Python for Data Analysis*. 2nd Edition. O'Reilly Media.
2. VanderPlas, J. (2016). *Python Data Science Handbook*. O'Reilly Media.
3. Pandas Development Team (2024). *Pandas Documentation*. <https://pandas.pydata.org/docs/>

Estándares Industriales

- ISO 22000:2018 — Food Safety Management Systems
- FDA 21 CFR Part 11 — Electronic Records and Signatures
- Codex Alimentarius — HACCP Principles

Datasets de Práctica

- Kaggle: Food Safety Inspections
- UCI Machine Learning Repository: Wine Quality Dataset
- Open Food Facts: Global food products database