

Data Scientist Senior Toolkit

Roadmap, Curación y Feature Engineering en Pandas

Curso: Fundamentos en Ciencia de Datos(Maestría) | **Periodo:** 2026-1

Docente: Jorge Iván Padilla-Buriticá | Universidad EAFIT

Filosofía de la Excelencia en Datos

Un Consultor Senior entiende que el algoritmo es solo el 10 % del éxito. El 90 % restante reside en la calidad del dato. Documentar cada transformación no es burocracia, es **trazabilidad**. Si un dashboard no es auditável, su recomendación no tiene valor de negocio.

1 Roadmap: Pasos previos a la Visualización

Antes de construir cualquier gráfico o dashboard, el científico de datos debe resolver el *Lineage* y la *Integridad*.

1. **Comprendión del Negocio:** ¿Qué KPI estamos moviendo? (Ej: Churn rate, Rentabilidad por SKU).
2. **Perfilamiento (Profiling):** Auditoría de tipos de datos, nulos y duplicados.
3. **Limpieza (Cleaning):** Estandarización de texto, manejo de nulos y outliers.
4. **Ingeniería (Feature Engineering):** Creación de variables que capturen la esencia del negocio.

2 Pandas Mastery: 20 Funciones Fundamentales

2.1 Auditoría e Identificación

1. isna() & sum()

```
1 # Identifica nulos por columna. Paso 1 de cualquier EDA.  
2 df.isna().sum()
```

2. info()

```
1 # Muestra tipos de datos y uso de memoria. Vital para optimizar tipos.  
2 df.info()
```

3. describe()

```
1 # Estadísticos rápidos. Ayuda a detectar ceros donde no debería haber.  
2 df.describe()
```

4. nunique()

```
1 # Cuenta valores únicos. Sirve para detectar variables constantes.  
2 df.nunique()
```

5. value_counts()

```
1 # Frecuencia de categorías. Detecta desbalanceo de clases.  
2 df['categoria'].value_counts(normalize=True)
```

2.2 Limpieza y Corrección

ALERTA TÉCNICA: Efectos del Cero

En campos como 'Presión Arterial' o 'Costo Unitario', el valor 0 es un nulo técnico. No lo dejes así; cámbialo a NaN antes de imputar.

6. replace()

```
1 import numpy as np  
2 df['variable'] = df['variable'].replace(0, np.nan)
```

7. dropna()

```
1 # Elimina filas con nulos críticos (ej. en la variable objetivo)  
2 df.dropna(subset=['target_y'], inplace=True)
```

8.fillna()

```
1 # Imputación rápida. Se recomienda usar la mediana para evitar  
# outliers.  
2 df['edad'].fillna(df['edad'].median(), inplace=True)
```

9. str.strip() & lower()

```
1 # Limpieza de strings para evitar duplicados por espacios o mayúsculas  
2 df['nombre'] = df['nombre'].str.strip().str.lower()
```

10. to_datetime()

```
1 # Conversión de fechas con manejo de errores (coerce).  
2 df['fecha'] = pd.to_datetime(df['fecha'], errors='coerce')
```

2.3 Transformación y Feature Engineering

11. astype()

```
1 # Casting de tipos para ahorrar memoria (ej: float64 a float32).
2 df['id'] = df['id'].astype(str)
```

12. get_dummies()

```
1 # One-Hot Encoding: Convierne categorias en variables numericas (0/1).
2 df_enc = pd.get_dummies(df, columns=['ciudad'], drop_first=True)
```

13. apply()

```
1 # Aplicar funciones personalizadas fila a fila o columna a columna.
2 df['rango_edad'] = df['edad'].apply(lambda x: 'Adulto' if x >= 18 else 'Menor')
```

14. clip()

```
1 # Acota outliers (Winsorizacion) sin eliminar la muestra.
2 df['salario'] = df['salario'].clip(lower=1000, upper=50000)
```

15. map()

```
1 # Sustitucion de valores basada en diccionarios (limpieza de categorias).
2 mapping = {'Med': 'Medellin', 'MDE': 'Medellin'}
3 df['ciudad'] = df['ciudad'].map(mapping).fillna(df['ciudad'])
```

16. groupby().transform()

```
1 # Normalizacion por grupos (ej: restar la media del departamento).
2 df['rel_precio'] = df['precio'] - df.groupby('cat')['precio'].transform('mean')
```

17. drop_duplicates()

```
1 # Elimina filas identicas garantizando integridad.
2 df.drop_duplicates(subset=['id_transaccion'], keep='first', inplace=True)
```

18. melt()

```
1 # Pasa datos de formato "Ancho" a "Largo" (Tidy Data).
2 df_long = df.melt(id_vars=['ID'], value_vars=['2024', '2025'])
```

19. merge()

```

1 # Une tablas. El 'left' garantiza que no pierdas ventas sin SKU
   asociado.
2 df_final = df_ventas.merge(df_productos, on='SKU', how='left')

```

20. nlargest()

```

1 # Obtiene los N registros mas altos sin ordenar todo el dataset.
2 df.nlargest(10, 'ventas')

```

3 Guía Técnica de Visualización Estadística

Seleccionar el gráfico incorrecto puede llevar a interpretaciones erróneas.

Tipo de Variable	Gráfico Sugerido	Propósito Estadístico
Numérica Continua	Histograma / KDE	Ver la forma de la distribución (Normal vs Sesgada).
Numérica Continua	Boxplot	Identificar Outliers y el Rango Intercuartílico (IQR).
Categórica	Gráfico de Barras	Comparar frecuencias o magnitudes entre grupos.
Num vs Num	Scatter Plot	Analizar correlación y densidad de puntos.
Num vs Tiempo	Gráfico de Líneas	Identificar tendencias, ciclos y estacionalidad.
Categórica vs Num	Violin Plot	Ver la distribución numérica dentro de cada categoría.

ALERTA TÉCNICA: Sobrealuste Visual

Evite los gráficos de torta (Pie charts) cuando tenga más de 3 categorías. El ojo humano tiene dificultades para comparar áreas; prefiera siempre gráficos de barras horizontales.

4 Resumen de Imputación Senior

Método	Escenario	Efecto Secuencial
Media	Distribución Normal	Mantiene el promedio, reduce la varianza.
Mediana	Distribución con Outliers	Robusta ante valores extremos.
Moda	Variables Categóricas	Mantiene la consistencia de clases.
Forward Fill	Series de Tiempo	Asume que el valor actual es igual al anterior.

"Un dashboard es la punta del iceberg; la ingeniería de datos es lo que lo mantiene a flote."