

# Estadística Descriptiva Aplicada: Teórico y Práctico con Python

Julio 2025

## Resumen

Este documento explora los fundamentos de la estadística descriptiva, incluyendo medidas de tendencia central (media, mediana, moda), medidas de dispersión (desviación estándar, varianza, cuartiles, rango) y correlación entre variables. Se presentan conceptos teóricos con sus formulaciones matemáticas y ejemplos prácticos implementados en Python utilizando las bibliotecas Pandas y NumPy, basados en el libro *Python for Data Analysis* de Wes McKinney (Capítulo 5). Además, se analiza un dataset simulado para ilustrar la aplicación de estas técnicas en la identificación de características principales de los datos.

## 1. Introducción

La estadística descriptiva es una rama de la estadística que se centra en resumir y organizar datos para facilitar su interpretación. Según McKinney (2017), en *Python for Data Analysis* (Capítulo 5, p. 141), las estructuras de datos de Pandas, como Series y DataFrame, son ideales para realizar análisis descriptivos debido a su flexibilidad y capacidad para manejar grandes volúmenes de datos. Este documento aborda las medidas de tendencia central, dispersión y correlación, integrando teoría, fórmulas matemáticas y ejemplos prácticos en Python.

El objetivo es proporcionar una guía completa que combine fundamentos teóricos con aplicaciones prácticas, utilizando herramientas computacionales modernas como Pandas y NumPy. Los ejemplos incluyen cálculos de estadísticas descriptivas y visualizaciones, siguiendo las prácticas recomendadas por McKinney (Capítulo 5, pp. 177-181) para el análisis de datasets.

## 2. Medidas de Tendencia Central

Las medidas de tendencia central resumen un conjunto de datos en un valor representativo. Las más comunes son la media, la mediana y la moda.

### 2.1. Media

La media aritmética se calcula como el promedio de todos los valores en un dataset. Matemáticamente, para un conjunto de  $n$  observaciones  $x_1, x_2, \dots, x_n$ , la media se define

como:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

La media es sensible a valores extremos, lo que puede distorsionar su representatividad en datasets con outliers.

## 2.2. Mediana

La mediana es el valor que divide un dataset ordenado en dos partes iguales. Si  $n$  es impar, es el valor central; si  $n$  es par, es el promedio de los dos valores centrales. Es robusta frente a valores atípicos, lo que la hace útil en distribuciones sesgadas.

## 2.3. Moda

La moda es el valor que aparece con mayor frecuencia en el dataset. Puede haber múltiples modas (multimodal) o ninguna si todos los valores son únicos. Es especialmente útil para datos categóricos.

**Ejemplo Práctico:** Consideremos un dataset de ventas diarias (en miles): [10, 12, 12, 15, 20, 25, 25, 30, 50]. Usando Python:

```
1 import numpy as np
2 import pandas as pd
3 from scipy import stats
4
5 ventas = [10, 12, 12, 15, 20, 25, 25, 30, 50]
6 media = np.mean(ventas) # Resultado: 22.11
7 mediana = np.median(ventas) # Resultado: 20
8 moda = stats.mode(ventas)[0][0] # Resultado: 12
9 print(f"Media: {media:.2f}, Mediana: {mediana}, Moda: {moda}")
```

La media (22.11) está influenciada por el valor extremo (50), mientras que la mediana (20) ofrece una mejor representación del centro. La moda (12) indica el valor más frecuente.

## 3. Medidas de Dispersión

Las medidas de dispersión cuantifican la variabilidad de los datos, complementando las medidas de tendencia central.

### 3.1. Varianza

La varianza mide el promedio de las desviaciones al cuadrado respecto a la media. Para una muestra, se calcula como:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

El uso de  $n-1$  (grados de libertad) corrige el sesgo en muestras pequeñas, como se explica en McKinney (Capítulo 5, p. 177).

### 3.2. Desviación Estándar

La desviación estándar es la raíz cuadrada de la varianza:

$$s = \sqrt{s^2}$$

Proporciona una medida de dispersión en las mismas unidades que los datos, facilitando su interpretación.

### 3.3. Cuartiles

Los cuartiles dividen un dataset ordenado en cuatro partes iguales: Q1 (25 %), Q2 (50 %, mediana), Q3 (75 %). El rango intercuartílico (IQR) se calcula como  $IQR = Q3 - Q1$  y es útil para detectar outliers.

### 3.4. Rango

El rango es la diferencia entre el valor máximo y mínimo:

$$R = \text{máx}(x) - \text{mín}(x)$$

**Ejemplo Práctico:** Usando el mismo dataset de ventas:

```
1 varianza = np.var(ventas, ddof=1) # Resultado: 166.11
2 desv_std = np.std(ventas, ddof=1) # Resultado: 12.89
3 cuartiles = np.percentile(ventas, [25, 50, 75]) # Resultado:
  [12, 20, 25]
4 rango = np.ptp(ventas) # Resultado: 40
5 print(f"Varianza: {varianza:.2f}, Desv. Estándar:
  {desv_std:.2f}, Cuartiles: {cuartiles}, Rango: {rango}")
```

La alta desviación estándar (12.89) refleja la dispersión causada por el valor extremo (50). Los cuartiles muestran que el 50 % central de los datos está entre 12 y 25.

## 4. Correlación entre Variables

La correlación mide la relación lineal entre dos variables. El coeficiente de correlación de Pearson se define como:

$$r = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2}}$$

donde  $-1 \leq r \leq 1$ . Un valor cercano a 1 indica una correlación positiva fuerte, cercano a -1 una correlación negativa, y cercano a 0 una relación débil o nula.

**Ejemplo Práctico:** Analicemos la relación entre horas de estudio y calificaciones:

```
1 df = pd.DataFrame({
2     'horas_estudio': [2, 4, 6, 8, 10],
3     'calificaciones': [60, 70, 75, 85, 90]
4 })
5 correlacion = df['horas_estudio'].corr(df['calificaciones']) #
  Resultado: 0.97
6 print(f"Correlación: {correlacion:.2f}")
```

El valor de 0.97 indica una fuerte correlación positiva: a más horas de estudio, mejores calificaciones.

## 5. Análisis Descriptivo con Pandas y NumPy

El análisis descriptivo combina medidas de tendencia central, dispersión y correlación para caracterizar un dataset. Según McKinney (Capítulo 5, pp. 177181), el método `describe()` de Pandas proporciona un resumen estadístico que incluye conteo, media, desviación estándar, cuartiles y valores extremos.

**Ejemplo Completo:** Consideremos un dataset simulado de ventas y gastos publicitarios:

```
1 df = pd.DataFrame({
2     'ventas': [10, 12, 12, 15, 20, 25, 25, 30, 50],
3     'publicidad': [5, 6, 6, 8, 10, 12, 12, 15, 20]
4 })
5 # Resumen estadístico
6 resumen = df.describe()
7 print(resumen)
8 # Matriz de correlación
9 matriz_corr = df.corr()
10 print(matriz_corr)
11 # Visualización
12 import seaborn as sns
13 import matplotlib.pyplot as plt
14 sns.heatmap(matriz_corr, annot=True, cmap='coolwarm')
15 plt.title('Matriz de Correlación')
16 plt.show()
```

El método `describe()` genera: - Conteo: 9 observaciones. - Media de ventas: 22.11, publicidad: 10.44. - Desviación estándar: 12.89 (ventas), 4.85 (publicidad). - Cuartiles: Q1=12, Q2=20, Q3=25 (ventas); Q1=6, Q2=10, Q3=12 (publicidad). La matriz de correlación muestra una correlación positiva fuerte (e.g., 0.95) entre ventas y publicidad, indicando que mayores gastos en publicidad se asocian con mayores ventas.

## 6. Visualización de Datos

La visualización es clave para interpretar estadísticas descriptivas. McKinney (Capítulo 9, pp. 271303) destaca el uso de bibliotecas como Matplotlib y Seaborn para crear gráficos informativos. Por ejemplo, un diagrama de caja (boxplot) puede mostrar la distribución y detectar outliers:

```
1 sns.boxplot(data=df['ventas'])
2 plt.title('Diagrama de Caja de Ventas')
3 plt.show()
```

Este gráfico revela el valor extremo (50) como un outlier, confirmando la influencia en la media y la desviación estándar.

## 7. Aplicaciones Prácticas

El análisis descriptivo es fundamental en la preparación de datos para minería de datos y aprendizaje automático. Por ejemplo, identificar outliers puede guiar la limpieza de datos, mientras que la correlación ayuda a seleccionar variables relevantes para modelos predictivos. McKinney (Capítulo 5, p. 179) enfatiza que el método `describe()` es una herramienta inicial para explorar datasets antes de análisis más avanzados.

**Ejemplo Avanzado:** Analicemos un dataset más complejo con datos faltantes:

```
1 df = pd.DataFrame({
2     'ventas': [10, 12, None, 15, 20, 25, 25, 30, 50],
3     'publicidad': [5, 6, 6, 8, None, 12, 12, 15, 20],
4     'region': ['A', 'B', 'B', 'A', 'C', 'C', 'A', 'B', 'C']
5 })
6 # Manejo de datos faltantes
7 df.fillna(df.mean(), inplace=True)
8 # Estadísticas por región
9 stats_por_region = df.groupby('region').describe()
10 print(stats_por_region)
11 # Correlación
12 print(df.corr())
```

Este código maneja valores faltantes imputando la media y calcula estadísticas descriptivas por región, mostrando cómo Pandas facilita el análisis de subgrupos.

## 8. Conclusión

La estadística descriptiva proporciona una base sólida para entender datasets, identificar patrones y preparar datos para análisis avanzados. Las bibliotecas Pandas y NumPy, como se detalla en McKinney (Capítulo 5), permiten realizar cálculos eficientes de medidas de tendencia central, dispersión y correlación, además de generar visualizaciones informativas. Los ejemplos prácticos presentados demuestran cómo aplicar estas técnicas a datasets reales o simulados, destacando la importancia de interpretar resultados en contexto, especialmente en presencia de valores extremos o datos faltantes.

## 9. Referencias

- McKinney, W. (2017). *Python for Data Analysis*, Capítulo 5, pp. 141181. O'Reilly Media.
- Hernández Orallo, J. (2004). *Introducción a la Minería de Datos*, Capítulo 2.