

# Sesión 2: Estadística descriptiva con Matplotlib (2 horas)

## 2.1 Introducción a la estadística descriptiva

### 2.1.1 Medidas de tendencia central: media, mediana, moda

La **media** es el promedio aritmético de los valores en un conjunto de datos. La **mediana** es el valor central que separa los datos en dos partes iguales. La **moda** es el valor que ocurre con mayor frecuencia en el conjunto de datos.

```
1 import numpy as np
2
3 data = np.array([4, 5, 6, 4, 6, 7, 8, 4, 9, 10])
4
5 media = np.mean(data)
6 mediana = np.median(data)
7 moda = np.bincount(data).argmax()
8
9 print("Media:", media)
10 print("Mediana:", mediana)
11 print("Moda:", moda)
```

### 2.1.2 Medidas de dispersión: rango, varianza, desviación estándar

El **rango** es la diferencia entre el valor máximo y el valor mínimo del conjunto de datos. La **varianza** es el promedio de las diferencias al cuadrado entre los valores y la media. La **desviación estándar** es la raíz cuadrada de la varianza.

```
1 rango = np.max(data) - np.min(data)
2 varianza = np.var(data)
3 desviacion_estandar = np.std(data)
4
5 print("Rango:", rango)
6 print("Varianza:", varianza)
7 print("Desviación estándar:", desviacion_estandar)
```

### 2.1.3 Cuartiles y diagrama de caja y bigotes (boxplot)

Los **cuartiles** son valores que dividen un conjunto de datos en cuatro partes iguales. El primer cuartil (Q1) es el valor que separa el 25% inferior de los datos, el segundo cuartil (Q2) es la mediana, y el tercer cuartil (Q3) es el valor que separa el 75% inferior de los datos.

```
1 Q1, Q2, Q3 = np.percentile(data, [25, 50, 75])
2
3 print("Primer cuartil (Q1):", Q1)
4 print("Segundo cuartil (Q2 - Mediana):", Q2)
5 print("Tercer cuartil (Q3):", Q3)
```

## 2.2 Introducción a Matplotlib

### 2.2.1 Instalación y configuración

```
1 pip install matplotlib
```

```
1 import matplotlib.pyplot as plt
2 plt.style.use("ggplot")
```

### 2.2.2 Creación de gráficos básicos: histogramas, barras, pastel

```
1 # Histograma
2 data_hist = np.random.randn(1000)
3 plt.hist(data_hist, bins=30)
4 plt.title("Histograma")
5 plt.xlabel("Valores")
6 plt.ylabel("Frecuencia")
7 plt.show()
8
9 # Gráfico de barras
10 categorias = ["A", "B", "C", "D", "E"]
11 valores = [3, 7, 2, 5, 8]
12 plt.bar(categorias, valores)
13 plt.title("Gráfico de barras")
14 plt.xlabel("Categorías")
15 plt.ylabel("Valores")
16 plt.show()
17
18 # Gráfico de pastel
19 labels = ["Python", "Java", "C++", "JavaScript", "Ruby"]
20 sizes = [45, 30, 15, 8, 2]
21 colors = ["blue", "green", "red", "yellow", "purple"]
22 explode = (0.1, 0, 0, 0, 0)
23
24 plt.pie(sizes, labels=labels, colors=colors, explode=explode,
25         autopct="%1.1f%%", shadow=True, startangle=90)
26 plt.title("Lenguajes de programación")
27 plt.axis("equal")
28 plt.show()
```

Con este bloque de código, hemos creado un gráfico de pastel que muestra la proporción de distintos lenguajes de programación. Las variables `labels` y `sizes` representan las etiquetas y los tamaños de las porciones del gráfico de pastel, respectivamente. La variable `colors` contiene los colores de las porciones y `explode` indica qué porción debe sobresalir del gráfico. La función `plt.pie()` crea el gráfico de pastel y la función `plt.show()` lo muestra en pantalla.

## 2.3 Visualización de estadísticas descriptivas con Matplotlib

### 2.3.1 Histogramas y gráficos de densidad

Los **histogramas** y **gráficos de densidad** son útiles para visualizar la distribución de los datos. Utilice `plt.hist` para crear histogramas y `plt.plot` con un objeto de densidad de Kernel para crear gráficos de densidad.

```
1 import seaborn as sns
2
3 # Histograma
4 data_hist = np.random.randn(1000)
5 plt.hist(data_hist, bins=30, density=True, alpha=0.6, color='g')
6 plt.title("Histograma")
7 plt.xlabel("Valores")
8 plt.ylabel("Densidad")
9 plt.show()
10
11 # Gráfico de densidad
12 sns.kdeplot(data_hist, color="r")
13 plt.title("Gráfico de densidad")
14 plt.xlabel("Valores")
15 plt.ylabel("Densidad")
16 plt.show()
```

### 2.3.2 Gráficos de caja y bigotes (boxplot)

Utilice `plt.boxplot` para crear gráficos de caja y bigotes que muestren la dispersión de los datos y los posibles valores atípicos.

```
1 data_box = [np.random.normal(0, std, 100) for std in range(1, 4)]
2
3 plt.boxplot(data_box, vert=True, patch_artist=True)
4 plt.title("Gráfico de caja y bigotes")
5 plt.xlabel("Categorías")
6 plt.ylabel("Valores")
7 plt.show()
```

### 2.3.3 Gráficos de dispersión y correlación

Los **gráficos de dispersión** son útiles para visualizar la relación entre dos variables numéricas. Los gráficos de correlación muestran la fuerza y la dirección de una relación lineal entre dos variables.

```
1 x = np.random.randn(100)
2 y = 2 * x + np.random.randn(100)
3
4 plt.scatter(x, y, marker="o", color="b")
5 plt.title("Gráfico de dispersión")
6 plt.xlabel("Variable X")
7 plt.ylabel("Variable Y")
8 plt.show()
9
10 correlacion = np.corrcoef(x, y)[0, 1]
11 print("Correlación:", correlacion)
```

## 2.4 Ejercicios prácticos con Python y Matplotlib

---

1. Utilice un conjunto de datos de su elección y calcule las medidas de tendencia central y de dispersión. Muestre los resultados en un gráfico de barras.
2. Cree un gráfico de pastel que muestre la proporción de diferentes categorías en su conjunto de datos. Por ejemplo, si su conjunto de datos contiene información sobre las ventas de productos, muestre la proporción de ventas de cada producto.
3. Dado un conjunto de datos con dos variables numéricas, cree un gráfico de dispersión para visualizar la relación entre las dos variables. Calcule la correlación entre las variables y muestre el coeficiente de correlación en el gráfico.