

Misión 2: Herramientas y Software para Análisis de Datos

Adriana Lucia González Ardila
Anthoni Lexandre Hernández Díaz
Carlos Alberto Robayo Melendez
José Arley León Méndez

Sector de Análisis y Justificación:

Para el desarrollo de la Misión 2: Herramientas y Software para Análisis de Datos se analizan los datos con información sobre la cantidad de defunciones ocurridas durante los años 2018 a 2023 en una clínica del Área Metropolitana de Bucaramanga, la cantidad de defunciones por sexo y la edad de las personas en el momento de su defunción. Esta información es muy importante para el país dado que a partir de ella se establecen estadísticas respecto a tasa de mortalidad como un indicador de salud pública que permite a las instituciones del estado conocer los problemas de salud que afectan a la población, gestionar los diversos servicios de salud y establecer políticas que mejoren la calidad de vida de los habitantes.

Programa en Python para el cargue de los datos

```
1 # Importar las librerías necesarias
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.metrics import mean_squared_error, r2_score
6 import numpy as np
7 from scipy import stats
8
9 # 1. Cargar los datos desde un archivo Excel
10 df = pd.read_excel(r'c:\Carlos A Robayo M\Analitica de Datos Bootcamp\MORTALIDAD_2018a2022.xlsx', sheet_name='MORTALIDAD')
11
12 # Convertir las columnas "Edad Fallecido" y "AÑO" a enteros
13 df['Edad Fallecido'] = df['Edad Fallecido'].astype('Int64')
14 df['AÑO'] = df['AÑO'].astype('Int64')
15
16 # Selección de las columnas que se desea analizar
17 dataframe = df[['Edad Fallecido', 'AÑO', 'Sexo']]
18
19 # Eliminar filas donde el año es 2017, 2018 o 2014
20 dataframe = dataframe[~dataframe['AÑO'].isin([2016, 2017, 2024])]
```

Calcular las siguientes variables estadísticas de tendencia central:

```

1 # Función para imprimir estadísticas descriptivas
2 def print_statistics(column, col_name):
3     print(f"\nEstadísticas para {col_name}:")
4     print(f"Media: {column.mean()}")
5     print(f"Mediana: {column.median()}")
6     # Calcular moda y manejar ambos casos (valor único o lista)
7     moda = stats.mode(column, nan_policy='omit')
8     try:
9         moda_valor = moda.mode[0] # Si devuelve un array
10    except IndexError:
11        moda_valor = moda.mode if moda.mode.size > 0 else "No disponible"
12    print(f"Moda: {moda_valor}")
13    print(f"Varianza: {column.var()}")
14    print(f"Desviación Estándar: {column.std()}")
15
16 # Calcular estadísticas para "Sexo" y "Año"
17 # Para "Sexo" y "Año", solo mostramos la moda, ya que media, mediana, varianza y desviación estándar no tienen
18 # sentido en una variable categórica.
19 def calcular_moda_categorica(column):
20     valores, conteos = np.unique(column.dropna(), return_counts=True)
21     moda_valor = valores[np.argmax(conteos)]
22     return moda_valor
23
24 print("\nEstadísticas para Sexo:")
25 moda_sexo_valor = calcular_moda_categorica(dataframe['Sexo'])
26 print(f"Moda: {moda_sexo_valor}")
27
28 print("\nEstadísticas para el Año:")
29 moda_año_valor = calcular_moda_categorica(dataframe['AÑO'])
30 print(f"Moda: {moda_año_valor}")
31
32 # Calcular estadísticas para "Edad Fallecido"
33 print_statistics(dataframe['Edad Fallecido'], 'Edad Fallecido')

```

La edad media de una persona en el momento de su defunción es de 72,31 años; la edad de mayor frecuencia de las personas en el momento de su defunción es de 79 años y la edad que ocupa el lugar central de todos los datos cuando éstos están ordenados de menor a mayor es de 75 años.

Estadísticas para Edad Fallecido:

Media: 72.31820848816741

Mediana: 75.0

Moda: 79

La moda es que el sexo de la persona fallecida sea MASCULINO.

Estadísticas para Sexo:

Moda: MASCULINO

La moda para el mayor numero de fallecidos es el año 2021.

Estadísticas para el Año:

Moda: 2021

Calcular las siguientes variables estadísticas de dispersión:

Varianza: La medida de dispersión que representa la variabilidad de la serie de datos de la edad de las personas en el momento de su defunción con respecto a su media es de 224.87.

Desviación Estándar: La desviación estándar de la serie de datos de la edad de las personas en el momento de su defunción es de 14.99. Esta medida nos permite medir la dispersión de los valores en un conjunto de datos.

Estadísticas para Edad Fallecido:

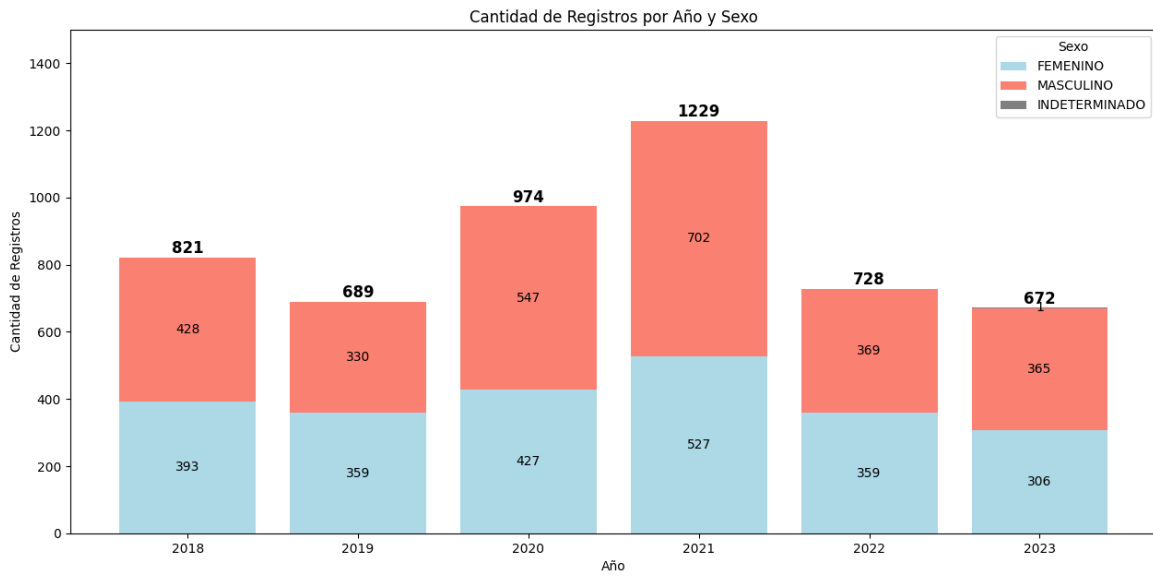
Varianza: 224.8730975723301

Desviación Estándar: 14.995769322456587

Realizar las siguientes graficas en Python

- **Gráficos de barras**

```
1 # Agrupar por "AÑO" y "Sexo", y contar las ocurrencias
2 barras = dataframe.groupby(['AÑO', 'Sexo']).size().unstack(fill_value=0)
3
4 # Agregar una columna "Total" con la suma de los registros por año
5 barras['Total'] = barras.sum(axis=1)
6
7 #Data de la lista a usar con año sexo y total
8 Data = barras.reset_index()
9
10 # Configuración del gráfico de barras
11 fig, ax = plt.subplots(figsize=(12, 8))
12
13 # Crear las barras apiladas
14 ax.bar(Data['AÑO'], Data['FEMENINO'], label='FEMENINO', color='lightblue')
15 ax.bar(Data['AÑO'], Data['MASCULINO'], bottom=Data['FEMENINO'], label='MASCULINO', color='salmon')
16
17 # Agregar la categoría 'INDETERMINADO' solo si existe
18 if 'INDETERMINADO' in Data.columns:
19     ax.bar(Data['AÑO'], Data['INDETERMINADO'],
20           bottom=Data['FEMENINO'] + Data['MASCULINO'],
21           label='INDETERMINADO', color='gray')
22
23 # Añadir etiquetas de valor para FEMENINO, MASCULINO e INDETERMINADO
24 for i, year in enumerate(Data['AÑO']):
25     fem_value = Data.loc[i, 'FEMENINO']
26     masc_value = Data.loc[i, 'MASCULINO']
27     total_value = Data.loc[i, 'Total']
28
29     # Etiqueta para FEMENINO
30     ax.text(year, fem_value / 2, str(fem_value), ha='center', va='center', color='black', fontsize=10)
31
32     # Etiqueta para MASCULINO, apilada sobre FEMENINO
33     ax.text(year, fem_value + (masc_value / 2), str(masc_value), ha='center', va='center', color='black', fontsize=10)
34
35     # Etiqueta para INDETERMINADO si existe
36     if 'INDETERMINADO' in Data.columns:
37         indet_value = Data.loc[i, 'INDETERMINADO']
38         if indet_value > 0:
39             ax.text(year, fem_value + masc_value + (indet_value / 2), str(indet_value), ha='center', va='center', color='black', fontsize=10)
40
41     # Etiqueta del total encima de la barra
42     ax.text(year, total_value + 2, str(total_value), ha='center', va='bottom', color='black', fontsize=12, fontweight='bold')
43
44 # Etiquetas y titulo
45 ax.set_xlabel('Año')
46 ax.set_ylabel('Cantidad de Registros')
47 ax.set_title('Cantidad de Registros por Año y Sexo')
48 ax.legend(title='Sexo')
49
50 # Configurar limites del eje y
51 ax.set_ylim(0, 1500)
52
53 # Mostrar el gráfico de barras
54 plt.tight_layout()
55 plt.show()
```

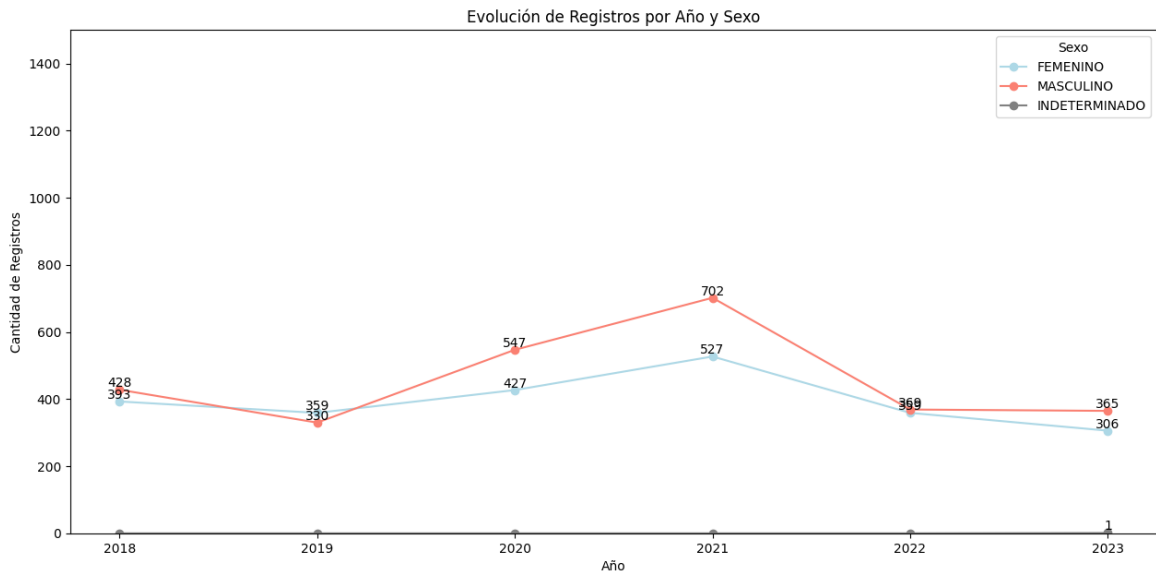


- **Gráficos de líneas**

```

1 # Gráfico de líneas
2 fig, ax = plt.subplots(figsize=(12, 8))
3
4 # Crear las líneas para cada categoría de sexo
5 ax.plot(Data['AÑO'], Data['FEMENINO'], marker='o', color='lightblue', label='FEMENINO')
6 ax.plot(Data['AÑO'], Data['MASCULINO'], marker='o', color='salmon', label='MASCULINO')
7
8 # Agregar la categoría 'INDETERMINADO' solo si existe
9 if 'INDETERMINADO' in Data.columns:
10     ax.plot(Data['AÑO'], Data['INDETERMINADO'], marker='o', color='gray', label='INDETERMINADO')
11
12 # Etiquetas de valor en cada punto
13 for i, year in enumerate(Data['AÑO']):
14     ax.text(year, Data.loc[i, 'FEMENINO'], str(Data.loc[i, 'FEMENINO']), ha='center', va='bottom', color='black')
15     ax.text(year, Data.loc[i, 'MASCULINO'], str(Data.loc[i, 'MASCULINO']), ha='center', va='bottom', color='black')
16
17     if 'INDETERMINADO' in Data.columns and Data.loc[i, 'INDETERMINADO'] > 0:
18         ax.text(year, Data.loc[i, 'INDETERMINADO'], str(Data.loc[i, 'INDETERMINADO']), ha='center', va='bottom', color='black')
19
20 # Configurar límites del eje y y etiquetas
21 ax.set_ylim(0, 1500)
22 ax.set_xlabel('Año')
23 ax.set_ylabel('Cantidad de Registros')
24 ax.set_title('Evolución de Registros por Año y Sexo')
25 ax.legend(title='Sexo')
26
27 # Mostrar el gráfico de líneas
28 plt.tight_layout()
29 plt.show()

```



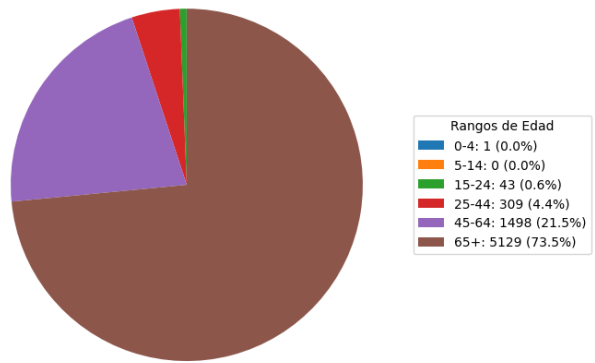
- **Gráficos de pastel**

```

1 # Grafico de Torta
2 # 1. Reemplazar valores en blanco en 'Edad Fallecido' con 0
3 dataframe['Edad Fallecido'] = dataframe['Edad Fallecido'].fillna(0)
4
5 # Definir los rangos de edad con un máximo de 120
6 bins = [0, 4, 14, 24, 44, 64, 120]
7 labels = ['0-4', '5-14', '15-24', '25-44', '45-64', '65+']
8 df['Rango de Edad'] = pd.cut(df['Edad Fallecido'], bins=bins, labels=labels, right=True)
9
10 # Contar la cantidad de ocurrencias en cada rango de edad
11 age_group_counts = df['Rango de Edad'].value_counts().sort_index()
12
13 # Crear el gráfico de torta sin etiquetas directamente en la torta
14 fig, ax = plt.subplots(figsize=(12, 8))
15 wedges, _ = ax.pie(age_group_counts, labels=['']*len(age_group_counts), startangle=90)
16
17 # Añadir el porcentaje en la leyenda
18 legend_labels = [f'{label}: {count} ({count / age_group_counts.sum():.1%})' for label, count in zip(age_group_counts.index, age_group_counts)]
19 plt.legend(wedges, legend_labels, title="Rangos de Edad", bbox_to_anchor=(1, 0.5), loc="center left")
20
21 # Título y ajuste de layout
22 plt.title('Distribución de Edad de Fallecimiento')
23 plt.tight_layout()
24 plt.show()

```

Distribución de Edad de Fallecimiento

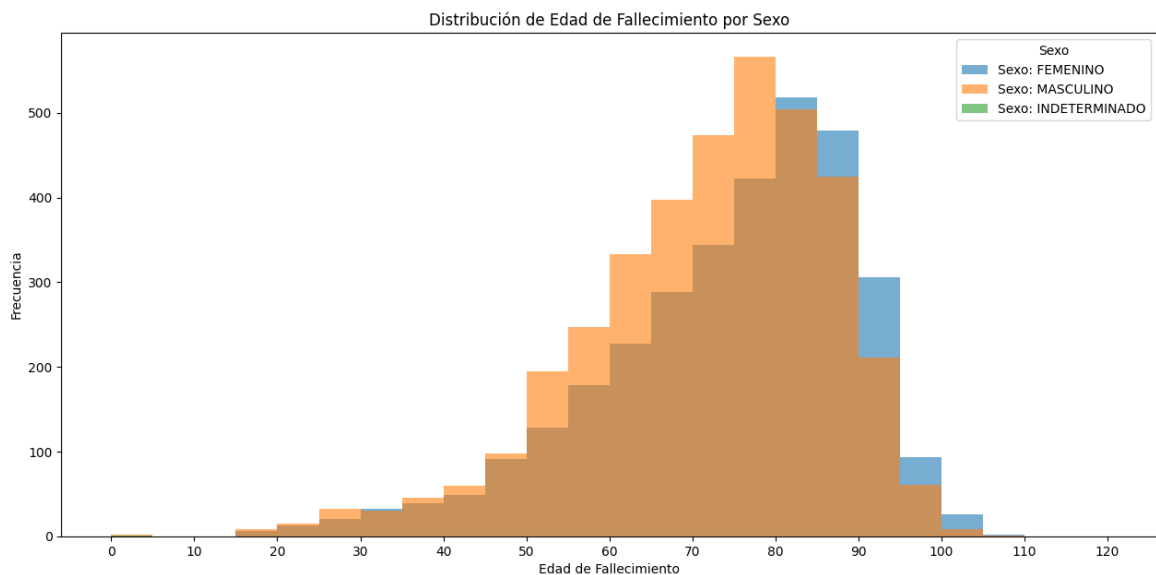


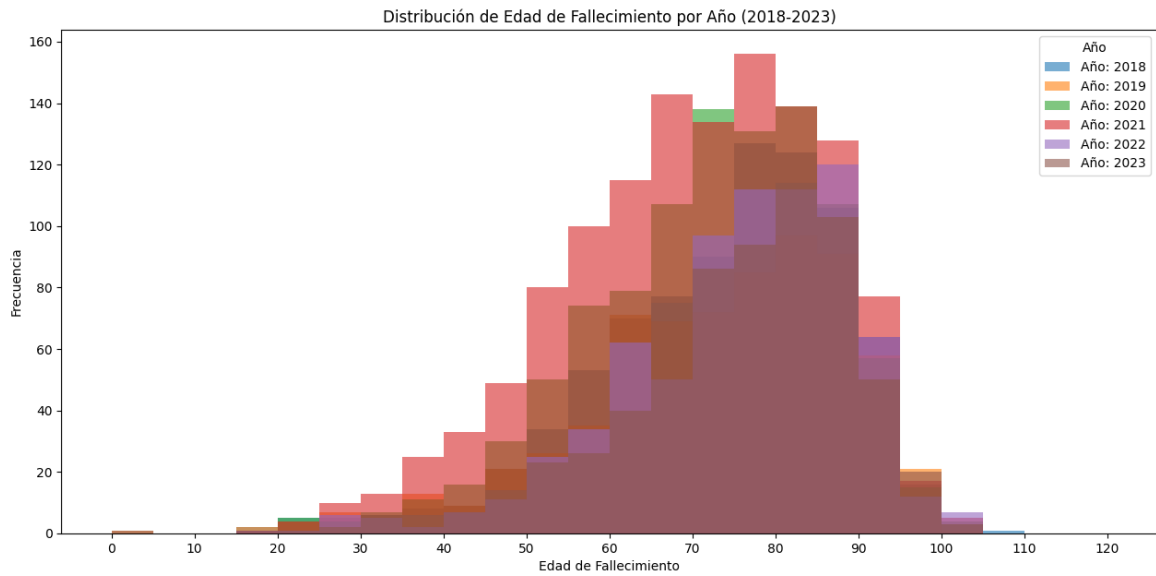
- **Histograma**

```

1  # Histograma
2  # Asegurar que 'Edad Fallecido' es numérico y eliminar valores nulos
3  df['Edad Fallecido'] = pd.to_numeric(df['Edad Fallecido'], errors='coerce').fillna(0).astype(int)
4
5  # Crear el histograma de edad, separado por sexo
6  plt.figure(figsize=(12, 8))
7  for sexo in df['Sexo'].unique():
8      subset = df[df['Sexo'] == sexo]
9      plt.hist(subset['Edad Fallecido'], bins=range(0, 121, 5), alpha=0.6, label=f'Sexo: {sexo}')
10
11 # Configuración del gráfico
12 plt.xlabel('Edad de Fallecimiento')
13 plt.ylabel('Frecuencia')
14 plt.title('Distribución de Edad de Fallecimiento por Sexo')
15 plt.legend(title='Sexo')
16 plt.xticks(range(0, 121, 10))
17 plt.tight_layout()
18 plt.show()
19
20 # Histograma de edad y año
21 # Filtrar los datos para solo incluir los años 2018 a 2023
22 df_filtered = df[(df['AÑO'] >= 2018) & (df['AÑO'] <= 2023)]
23
24 # Crear el histograma de edad, separado por año (solo de 2018 a 2023)
25 plt.figure(figsize=(12, 8))
26 for year in sorted(df_filtered['AÑO'].unique()):
27     subset = df_filtered[df_filtered['AÑO'] == year]
28     plt.hist(subset['Edad Fallecido'], bins=range(0, 121, 5), alpha=0.6, label=f'Año: {year}')
29
30 # Configuración del gráfico
31 plt.xlabel('Edad de Fallecimiento')
32 plt.ylabel('Frecuencia')
33 plt.title('Distribución de Edad de Fallecimiento por Año (2018-2023)')
34 plt.legend(title='Año')
35 plt.xticks(range(0, 121, 10))
36 plt.tight_layout()
37 plt.show()

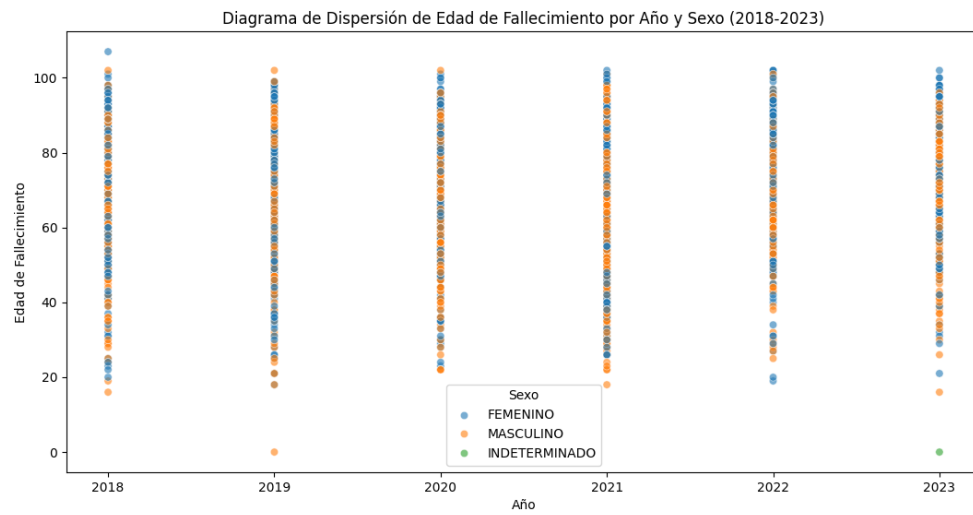
```





- **Diagramas de dispersión**

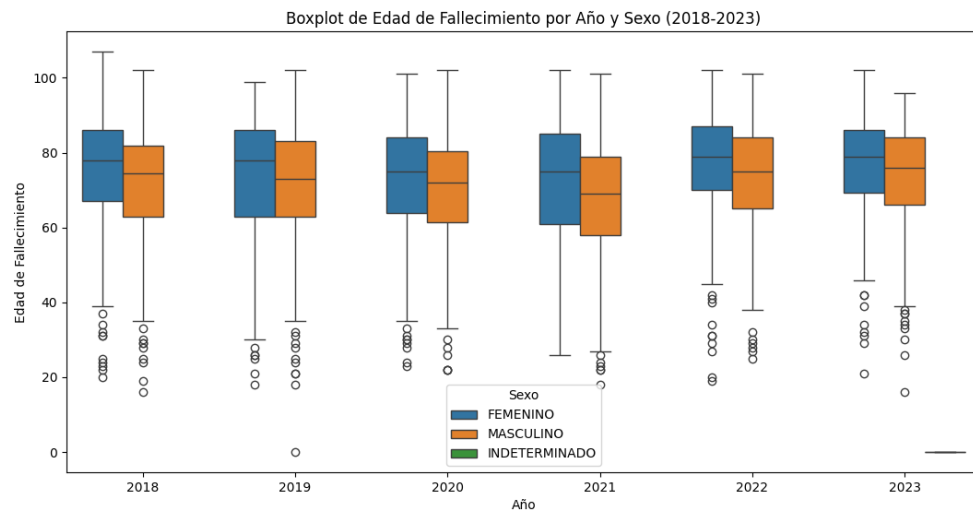
```
1 # Grafico de dispersion
2 # Filtrar los datos para los años 2018 a 2023
3 df_filtered = df[(df['AÑO'] >= 2018) & (df['AÑO'] <= 2023)]
4
5 # Crear el diagrama de dispersión
6 plt.figure(figsize=(12, 8))
7 sns.scatterplot(data=df_filtered, x='AÑO', y='Edad Fallecido', hue='Sexo', alpha=0.6)
8 plt.xlabel('Año')
9 plt.ylabel('Edad de Fallecimiento')
10 plt.title('Diagrama de Dispersión de Edad de Fallecimiento por Año y Sexo (2018-2023)')
11 plt.legend(title='Sexo')
12 plt.show()
```

- **Boxplots**



```
1 # Grafico de boxpot
2 plt.figure(figsize=(12, 8))
3 sns.boxplot(data=df_filtered, x='AÑO', y='Edad Fallecido', hue='Sexo')
4 plt.xlabel('Año')
5 plt.ylabel('Edad de Fallecimiento')
6 plt.title('Boxplot de Edad de Fallecimiento por Año y Sexo (2018-2023)')
7 plt.legend(title='Sexo')
8 plt.show()
```

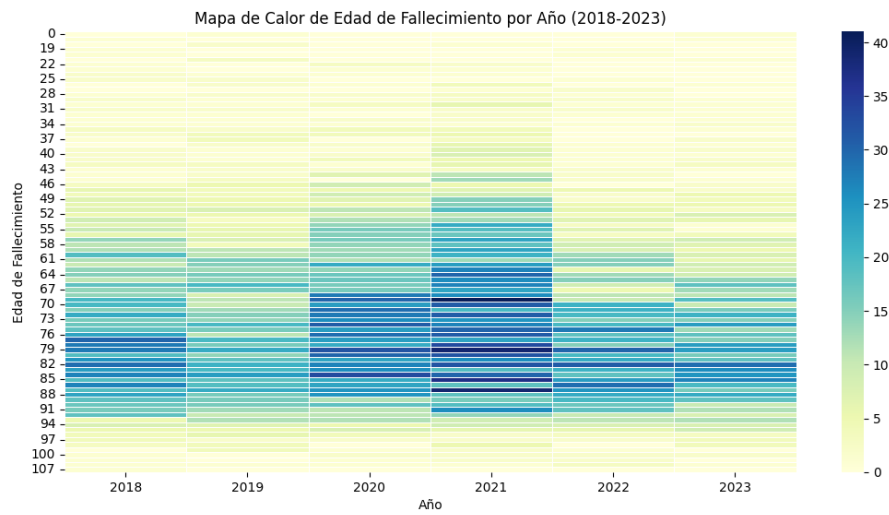


- **Mapas de calor**

```

1 # Grafico de calor
2 # Crear una tabla de frecuencias entre edad y año
3 heatmap_data = pd.crosstab(df_filtered['Edad Fallecido'], df_filtered['AÑO'])
4
5 # Crear el mapa de calor
6 plt.figure(figsize=(12, 8))
7 sns.heatmap(heatmap_data, cmap='YlGnBu', linewidths=0.5)
8 plt.xlabel('Año')
9 plt.ylabel('Edad de Fallecimiento')
10 plt.title('Mapa de Calor de Edad de Fallecimiento por Año (2018-2023)')
11 plt.show()

```



8. Realizar un modelo de regresión lineal

- Cargar lo datos en un dataframe

```

1 # Cargar los datos en un DataFrame
2 # Filtrar los datos para los años 2018 a 2023
3 df_filtered = df[(df['AÑO'] >= 2018) & (df['AÑO'] <= 2023)]
4
5 # Seleccionar las columnas para el modelo
6 df_model = df_filtered[['AÑO', 'Edad Fallecido']]

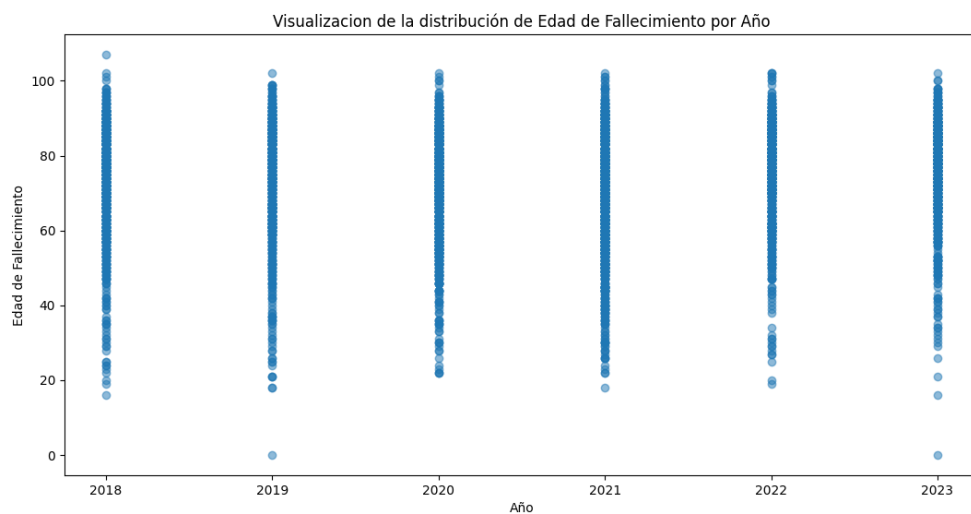
```

- Visualizar los datos

```

1 # Visualizar los datos
2 plt.figure(figsize=(12, 8))
3 plt.scatter(df_model['AÑO'], df_model['Edad Fallecido'], alpha=0.5)
4 plt.xlabel('Año')
5 plt.ylabel('Edad de Fallecimiento')
6 plt.title('Visualizacion de la distribución de Edad de Fallecimiento por Año')
7 plt.show()

```



- **Dividir los datos en conjuntos de entrenamiento y prueba**

```

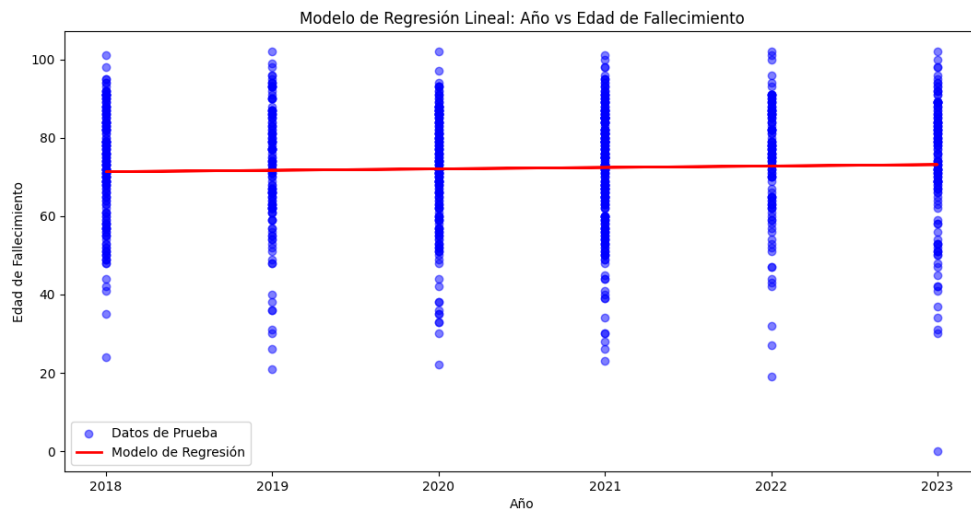
1 # Dividir los datos en conjuntos de entrenamiento y prueba
2 from sklearn.model_selection import train_test_split
3
4 X = df_model[['AÑO']]
5 y = df_model['Edad Fallecido']
6
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

```

- **Crear y entrenar el modelo de regresión lineal**



```
1 # Crear y entrenar el modelo de regresión lineal
2 from sklearn.linear_model import LinearRegression
3
4 model = LinearRegression()
5 model.fit(X_train, y_train)
```



- **Hacer predicciones**



```
1 # Hacer predicciones
2 y_pred = model.predict(X_test)
3
4 # Hacer la predicción para el año 2024 usando un DataFrame para conservar el nombre de la característica
5 edad_2024 = model.predict(pd.DataFrame([[2024]], columns=['AÑO']))
6 print(f"Predicción de la edad de fallecimiento para el año 2024: {edad_2024[0]:.2f}")
7
8 # Hacer la predicción para el año 2025 usando un DataFrame para conservar el nombre de la característica
9 edad_2025 = model.predict(pd.DataFrame([[2025]], columns=['AÑO']))
10 print(f"Predicción de la edad de fallecimiento para el año 2025: {edad_2025[0]:.2f}")
```

La edad de una persona en el momento de su defunción para el año 2024 será de 73,51 años y para el año 2025 será de 73,87 años.

Predicción de la edad de fallecimiento para el año 2024: 73.51

Predicción de la edad de fallecimiento para el año 2025: 73.87

- **Evaluar el modelo**

```
1 # Evaluar el modelo
2 # Calcular el Mean Squared Error (MSE)
3 mse = mean_squared_error(y_test, y_pred)
4
5 # Calcular el Root Mean Squared Error (RMSE)
6 rmse = np.sqrt(mse)
7
8 # Calcular el R2 Score
9 r2 = r2_score(y_test, y_pred)
10
11 print(f"\n\nMean Squared Error (MSE): {mse}")
12 print(f"Root Mean Squared Error (RMSE): {rmse}")
13 print(f"R2 Score: {r2}\n\n")
14
15 # Visualizar el modelo y las predicciones
16 plt.figure(figsize=(12, 8))
17 plt.scatter(X_test, y_test, color='blue', alpha=0.5, label='Datos de Prueba')
18 plt.plot(X_test, y_pred, color='red', linewidth=2, label='Modelo de Regresión')
19 plt.xlabel('Año')
20 plt.ylabel('Edad de Fallecimiento')
21 plt.title('Modelo de Regresión Lineal: Año vs Edad de Fallecimiento')
22 plt.legend()
23 plt.show()
```

- **Calcular el Mean Squared Error**

Cuanto menor sea el MSE, mejor será la precisión predictiva del modelo.

Mean Squared Error (MSE): 229.12684675665454

- **Calcular el Root Mean Squared Error**

Un valor RMSE más bajo indica que el modelo es más preciso y se ajusta bien a los datos.

Root Mean Squared Error (RMSE): 15.136936505008356

- **Calcular el R2 Score**

Es una medida estadística utilizada para determinar la proporción de varianza en una variable dependiente que puede predecirse o explicarse mediante una variable independiente. Un valor R-cuadrado de 0 significa que el modelo explica o predice el 0% de la relación entre las variables dependientes e independientes.

R2 Score: -0.0005196781800487216