

ADSO-GA04-Actividad de Aprendizaje Seis Análisis de Datos

Eddie Alejandro Arenas Vital

Juan Pablo Oviedo Herrera

Jhon Jairo Guzmán Caballero

SENA

Ciencia De Datos

2828523

Luis Fernando Sánchez

4 de Abril de 2025

Análisis Exploratorio de Datos (EDA) - Conjunto de Datos Iris

1. Introducción

Este documento presenta un Análisis Exploratorio de Datos (EDA) utilizando el conjunto de datos Iris de la librería Seaborn. El objetivo es identificar patrones, valores atípicos y realizar visualizaciones básicas. Se usan las bibliotecas Pandas, Numpy, Matplotlib y Seaborn para el análisis y generación de gráficos.

2. Carga y descripción de los datos

Código para cargar y describir el conjunto de datos.

```
ActividadUno.py > ...
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.datasets import load_iris
6 from sklearn.preprocessing import LabelEncoder
7
8 # Carga el conjunto de datos Iris desde seaborn
9 data = sns.load_dataset('iris')
10
11 iris = load_iris()
12 df = pd.DataFrame(iris.data, columns=iris.feature_names)
13 df['species'] = iris.target_names[iris.target]
14
15 # Convierte las etiquetas de 'species' en valores numéricos con LabelEncoder
16 label_encoder = LabelEncoder()
17 df['species'] = label_encoder.fit_transform(df['species'])
18
19 # Ver las primeras filas del conjunto de datos
20 print(data.head())
21 print(data.describe())
```

3. Identificación de valores atípicos

Los valores atípicos pueden afectar el análisis de datos. Para identificarlos, utilizamos diagramas de caja, mejor conocidos como boxplots.

```

ActividadUno.py > ...
22
23 # Generar boxplots
24 plt.figure(figsize=(12, 8))
25 sns.boxplot(data=data[['sepal_length', 'sepal_width', 'petal_length', 'petal_width']])
26 plt.title("Boxplot de las características del Iris")
27 plt.savefig("boxplot_iris.png")
28 plt.show()

```

4. Distribución de los datos

Se analizan las distribuciones de las variables mediante histogramas.

```

ActividadUno.py > ...
30 # Histogramas de las características numéricas
31 data[['sepal_length', 'sepal_width', 'petal_length', 'petal_width']].hist(bins=15, figsize=(12, 8))
32 plt.suptitle("Histogramas de las características del Iris")
33 plt.savefig("histogramas_iris.png")
34 plt.show()

```

5. Relación entre características

Se analizan relaciones entre características con gráficos de dispersión.

```

Bienvenido 0.py > correlation_matrix
36 # Pairplot para ver la relación entre las características numéricas
37 sns.pairplot(data, hue='species')
38 plt.suptitle("Pairplot de las características del Iris", y=1.02)
39 plt.savefig("pairplot_iris.png")
40 plt.show()
41

```

6. Correlación entre características

Se genera un mapa de calor para visualizar la correlación entre variables.

```

ActividadUno.py > ...
42 # Matriz de correlación
43 correlation_matrix = df.corr()
44 plt.figure(figsize=(10, 8))
45 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', linewidths=0.5)
46 plt.title("Mapa de calor de la correlación entre características")
47 plt.savefig("heatmap_iris.png")
48 plt.show()

```

7. Comparación entre especies

Se comparan las longitudes del sépalos entre especies con un gráfico de caja.

```
ActividadUno.py > ...  
50 # Gráfico de caja para comparar la longitud del sépalos entre las especies  
51 plt.figure(figsize=(8, 6))  
52 sns.boxplot(x='species', y='sepal_length', data=data)  
53 plt.title("Comparación de la longitud del sépalos por especie")  
54 plt.savefig("boxplot_sepal_length.png")  
55 plt.show()  
56
```

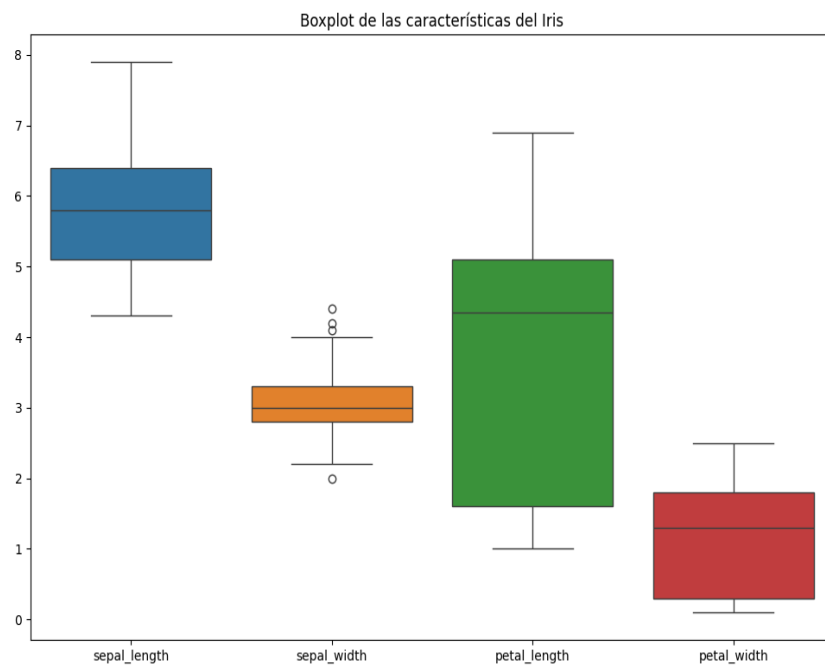
8. Conclusiones

El análisis del conjunto de datos "Iris" reveló diferencias significativas entre las especies en características como la longitud y el ancho del sépalos y pétalos, con Iris Setosa mostrando dimensiones más pequeñas. La visualización de datos a través de histogramas y scatter plots facilitó la identificación de patrones y correlaciones, anotando una alta correlación entre longitud del pétalos y longitud del sépalos. Se identificaron valores atípicos que requieren atención antes de implementar modelos predictivos.

9. Evidencias y Entrega

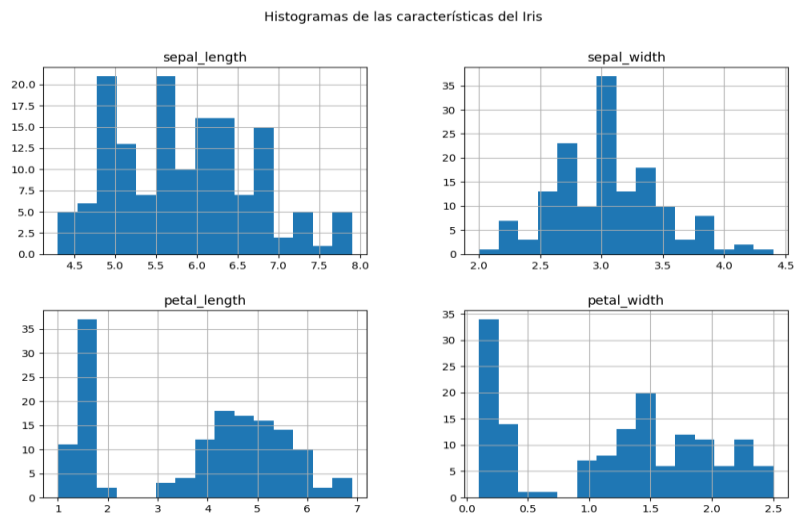
- Boxplot de las características del Iris

Este gráfico muestra la dispersión y posibles valores atípicos de las variables del conjunto de datos.



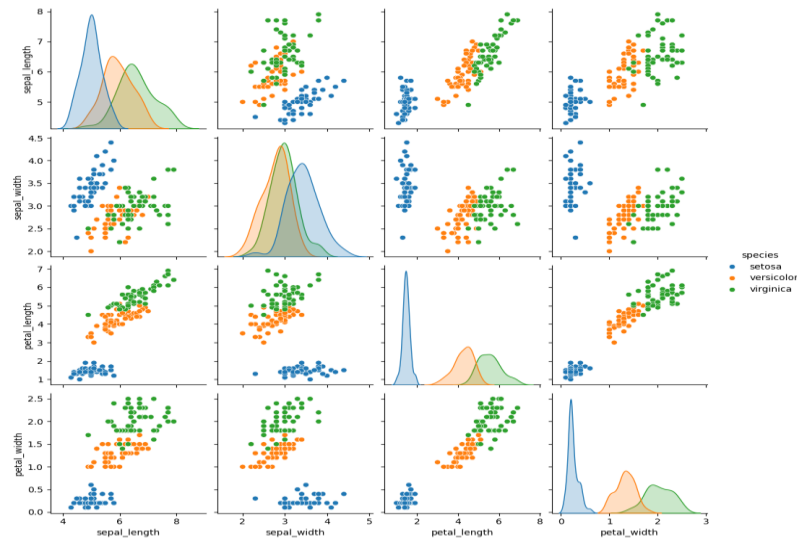
- Histogramas de las características del Iris

Representación de la distribución de las variables numéricas mediante histogramas.



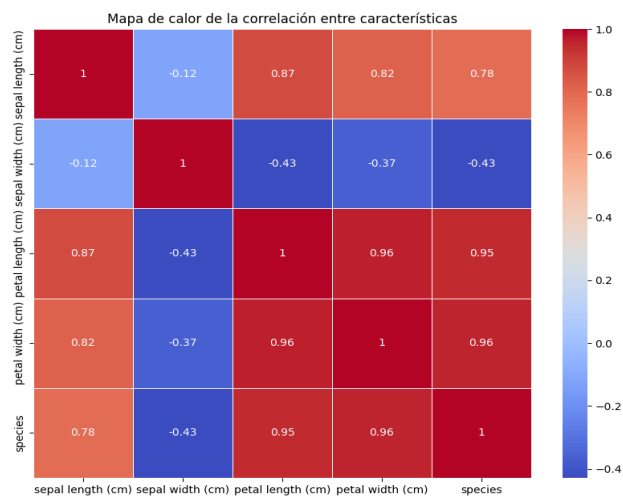
- Pairplot de las características del Iris

Visualización de las relaciones entre las características del conjunto de datos, agrupadas por especie.



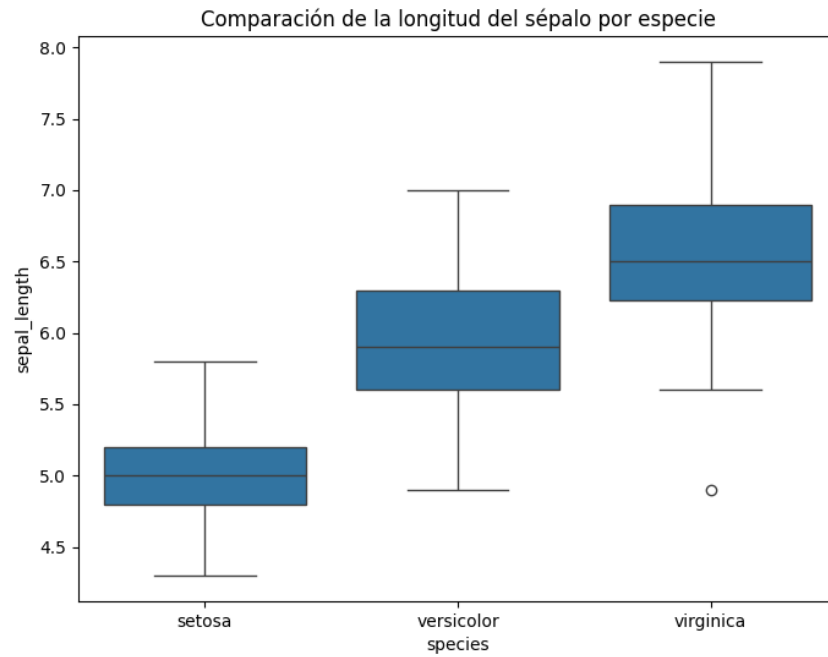
- Mapa de calor de la correlación entre características

Matriz de correlación entre las variables del conjunto de datos para analizar relaciones.



- Comparación de la longitud del sépalo por especie

Boxplot comparando la longitud del sépalo entre las tres especies de Iris.



Actividad 2: Implementación de Modelos de Machine Learning

1. Introducción

En esta actividad, se implementan diversos modelos de Machine Learning para la clasificación y regresión, utilizando el conjunto de datos Iris. Se exploraron modelos como la Regresión Logística, Árboles de Decisión, SVM, KNN, Redes Neuronales, Random Forest y Regresión Lineal.

2. Objetivos

- Implementar diferentes modelos de Machine Learning.
- Evaluar su desempeño en la clasificación de especies de flores.
- Realizar una regresión para predecir una variable continua del conjunto de datos.

3. Desarrollo

3.1. Carga de Datos y Preprocesamiento

Se utilizó la librería seaborn para cargar el conjunto de datos Iris. Posteriormente, se seleccionaron las variables independientes (sepal_length, sepal_width, petal_width) y la variable objetivo (petal_length). Además, se dividió el conjunto de datos en entrenamiento y prueba usando train_test_split.

3.2. Implementación del Modelo de Regresión Lineal

Se utilizó LinearRegression de scikit-learn para entrenar un modelo predictivo de petal_length. Por otra parte se calcularon los coeficientes y la intersección del modelo, realizando predicciones sobre el conjunto de prueba.

Se evaluó el desempeño mediante el error cuadrático medio (MSE) y el coeficiente de determinación (R^2) además de generar una gráfica comparativa entre los valores reales y las predicciones.

3.3. Evaluación de Resultados

Los valores obtenidos de MSE y R^2 indicaron el nivel de ajuste del modelo. Además, la gráfica de dispersión mostró la relación entre valores reales y predichos.

4. Conclusiones

- Se implementó exitosamente un modelo de regresión lineal para predecir petal_length.
- Se observó la importancia de seleccionar adecuadamente las variables predictoras.
- Se identificó la precisión del modelo y su aplicabilidad en problemas de predicción continua.

5. Anexos

Código utilizado para la implementación.

```
Actividad2 > ActividadDos.py > ...
1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  from sklearn.ensemble import RandomForestClassifier
6  from sklearn.metrics import mean_squared_error, r2_score, accuracy_score
7  from sklearn.model_selection import train_test_split
8  from sklearn.linear_model import LogisticRegression, LinearRegression
9  from sklearn.datasets import load_iris
10 from sklearn.neighbors import KNeighborsClassifier
11 from sklearn.neural_network import MLPClassifier
12 from sklearn.tree import DecisionTreeClassifier
13 from sklearn.svm import SVC
14
15 # Cargar dataset de seaborn (usar un solo dataset)
16 data = sns.load_dataset('iris')
17
18 # Seleccionamos las características y variable objetivo
19 x = data[['sepal_length', 'sepal_width', 'petal_width']] # Variables independientes
20 y = data['petal_length'] # Variable objetivo
21
22 # Dividir en conjunto de entrenamiento y prueba
23 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=42)
24
25 # Verificar tamaños de los conjuntos
26 print(f"Tamaño del conjunto de entrenamiento: {x_train.shape}, {y_train.shape}")
27 print(f"Tamaño del conjunto de prueba: {x_test.shape}, {y_test.shape}")
28
29 # Crear y entrenar el modelo de regresión lineal
30 model = LinearRegression()
31 model.fit(x_train, y_train)
32
33 # Ver coeficientes e intercepto
34 print(f"Coefficientes: {model.coef}")
35 print(f"Intersección (intercepto): {model.intercept}")
36
37 # Realizar predicciones
38 y_pred = model.predict(x_test)
39
40 # Mostrar predicciones
41 predictions_df = pd.DataFrame({'Real': y_test.values, 'Predicción': y_pred})
42 print(predictions_df.head())
43
44 # Calcular MSE y R²
45 mse = mean_squared_error(y_test, y_pred)
46 r2 = r2_score(y_test, y_pred)
47 print(f"Error cuadrático medio (MSE): {mse}")
48 print(f"Coefficiente de determinación R²: {r2}")
49
50 # Graficar resultados
51 plt.figure(figsize=(8, 6))
52 plt.scatter(y_test, y_pred, color='blue', label='Predicciones')
53 plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='red', linewidth=2, label="Línea de referencia")
54 plt.xlabel("Valores reales")
55 plt.ylabel("Predicciones")
56 plt.title("Valores reales vs Predicciones")
57 plt.legend()
58 plt.show()
```

Gráfica generada

