

```
1 # Import necessary libraries
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.model_selection import train_test_split
6 from sklearn.linear_model import LinearRegression
7 from sklearn.metrics import mean_squared_error, r2_score
8
9 # Load the dataset (replace with the correct path if needed)
10 file_path = "/content/Salary_dataset.csv"
11 df = pd.read_csv(file_path)
12
13 # Step 1: EDA
14 # Remove unnecessary column
15 df_cleaned = df.drop(columns=["Unnamed: 0"])
16
17 # Display basic info and statistics
18 print("Dataset Overview:\n", df_cleaned.head())
19 print("\nSummary Statistics:\n", df_cleaned.describe())
20
21 # Scatter plot to visualize the relationship
22 plt.figure(figsize=(8, 5))
23 sns.scatterplot(data=df_cleaned, x="YearsExperience", y="Salary", color="blue", s=80)
24 plt.title("Years of Experience vs Salary", fontsize=16)
25 plt.xlabel("Years of Experience", fontsize=12)
26 plt.ylabel("Salary (in ₹)", fontsize=12)
27 plt.grid(True)
28 plt.show()
29
30 # Correlation calculation
31 correlation = df_cleaned.corr().loc["YearsExperience", "Salary"]
32 print(f"\nCorrelation between YearsExperience and Salary: {correlation:.2f}")
33
34 # Step 2: Prepare Data for Linear Regression
35 X = df_cleaned[["YearsExperience"]] # Independent variable
36 y = df_cleaned["Salary"] # Dependent variable
37
38 # Split the dataset into training and testing sets (80% train, 20% test)
39 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
40
41 # Step 3: Train the Linear Regression Model
42 model = LinearRegression()
43 model.fit(X_train, y_train)
44
45 # Step 4: Make Predictions
46 y_pred = model.predict(X_test)
47
48 # Step 5: Evaluate the Model
49 mse = mean_squared_error(y_test, y_pred)
50 r2 = r2_score(y_test, y_pred)
51 print(f"\nMean Squared Error (MSE): {mse:.2f}")
52 print(f"\nR-squared (R2 Score): {r2:.2f}")
53
54 # Step 6: Visualize the Regression Line
55 plt.figure(figsize=(8, 5))
56 sns.scatterplot(x=X_train["YearsExperience"], y=y_train, color="blue", label="Training Data")
57 sns.scatterplot(x=X_test["YearsExperience"], y=y_test, color="red", label="Test Data", marker="x")
58 plt.plot(X_test, y_pred, color="green", linewidth=2, label="Regression Line")
59 plt.title("Simple Linear Regression: Years of Experience vs Salary", fontsize=16)
60 plt.xlabel("Years of Experience", fontsize=12)
61 plt.ylabel("Salary (in ₹)", fontsize=12)
62 plt.legend()
63 plt.grid(True)
64 plt.show()
65
```



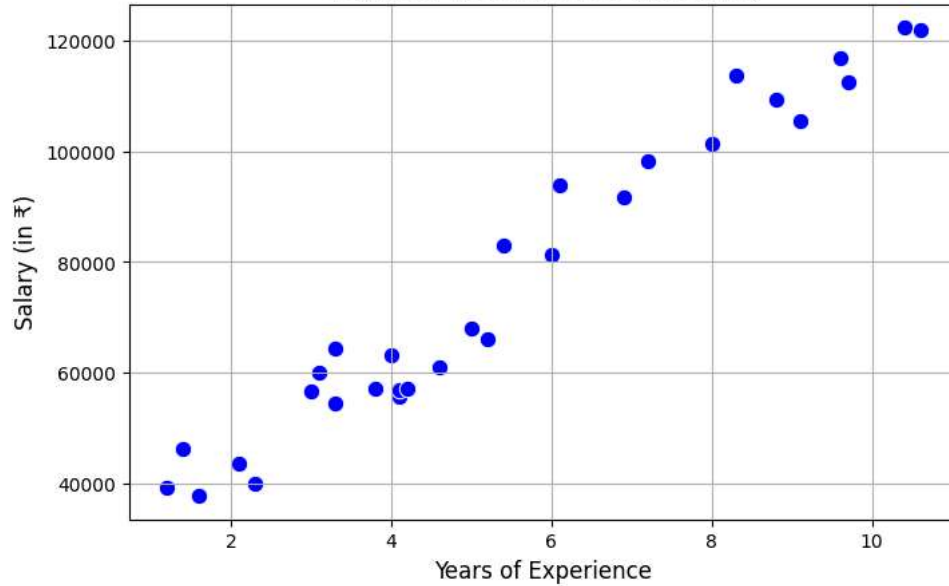
Dataset Overview:

	YearsExperience	Salary
0	1.2	39344.0
1	1.4	46206.0
2	1.6	37732.0
3	2.1	43526.0
4	2.3	39892.0

Summary Statistics:

	YearsExperience	Salary
count	30.000000	30.000000
mean	5.413333	76004.000000
std	2.837888	27414.429785
min	1.200000	37732.000000
25%	3.300000	56721.750000
50%	4.800000	65238.000000
75%	7.800000	100545.750000
max	10.600000	122392.000000

Years of Experience vs Salary



Correlation between YearsExperience and Salary: 0.98

Mean Squared Error (MSE): 49830096.86

R-squared (R2 Score): 0.90

Simple Linear Regression: Years of Experience vs Salary

