

Responsible AI (AI354): Lab Assignment 2

Linear Regression Analysis

Devesh Singh Chauhan (I23MA002)

January 18, 2026

1 Question 1: Impact of Shuffling

1.1 Methodology

The goal of this experiment was to see if the way we organize our data before feeding it into the model changes the results.

First, we took the salary dataset and split it into two parts: 70% was used to train the model, and the remaining 30% was hidden away to test it later.

To test the impact of shuffling, we didn't just split the data once. We repeated the process five separate times. In each "run," we shuffled the data in a different random order before splitting it. This meant that different data points ended up in the test set each time. For every run, we calculated the Mean Squared Error (MSE) to see how much the error rate fluctuated.

1.2 Results

As shown in the graph below, the error rate was not constant. It changed significantly depending on which specific data points happened to fall into the test set.

This happens because our dataset is quite small (only 30 rows). In a small dataset, a single "unusual" data point (like someone with low experience but a very high salary) can skew the results heavily if it ends up in the test set versus the training set.

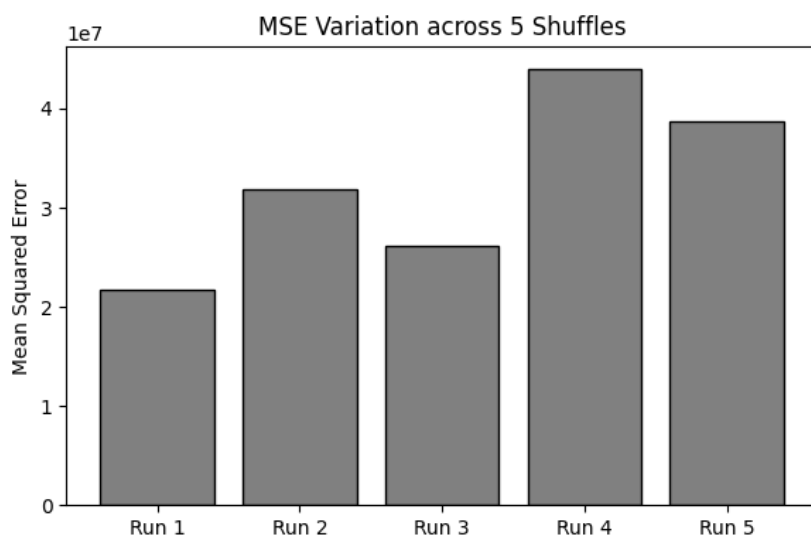


Figure 1: Fluctuation of MSE across 5 different random shuffles.

2 Question 2: Impact of Normalization

2.1 Methodology

In this second experiment, we wanted to see if the "scale" or "units" of the numbers mattered to the Linear Regression model. For example, does the model get confused because "Salary" (e.g., 40,000) is a much larger number than "Years of Experience" (e.g., 2.5)?

We used an 80:20 split for this task. We then trained three separate models:

1. **Raw Data:** We fed the numbers in exactly as they appeared in the CSV file.
2. **Min-Max Scaling:** We mathematically squeezed the "Years of Experience" so all values fell strictly between 0 and 1.
3. **Standard Normalization:** We adjusted the values so they had an average (mean) of 0 and a standard deviation of 1.

2.2 Results

After testing all three models, we found that the Mean Squared Error was identical for all of them ($MSE \approx 49.8 \times 10^6$).

This proves that for a standard Linear Regression model, normalization is not necessary for accuracy. The math behind the model automatically adjusts the "slope" (m) to compensate for the size of the input numbers. If we shrink the inputs, the model simply makes the slope steeper to arrive at the same prediction.

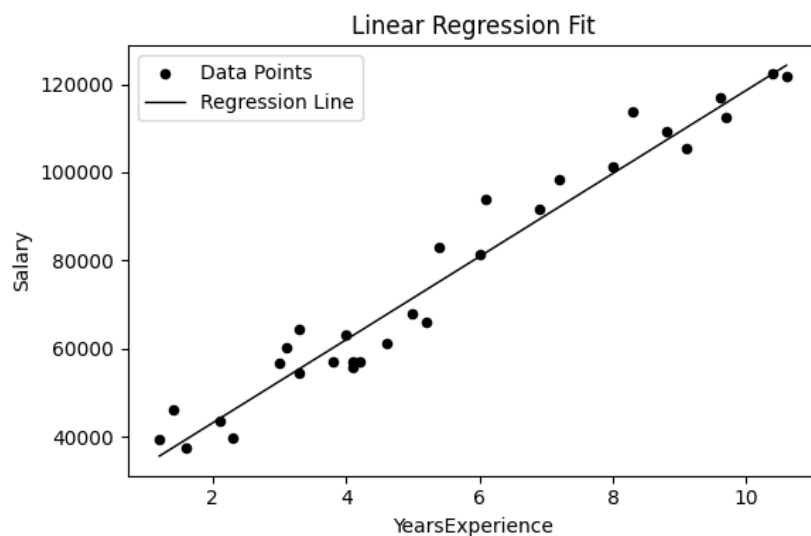


Figure 2: The Linear Regression fit used for the calculation.

3 Appendix: Python Code

3.1 Code for Question 1 (Shuffling)

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import LinearRegression
4 from sklearn.metrics import mean_squared_error
5
6 df = pd.read_csv('Salary_dataset.csv')
7 X = df[['YearsExperience']]
8 y = df['Salary']
9
10 # We loop through 5 different "seeds" to simulate 5 different shuffles
11 seeds = [10, 20, 30, 40, 50]
12 for seed in seeds:
13     X_train, X_test, y_train, y_test = train_test_split(
14         X, y, test_size=0.3, random_state=seed, shuffle=True
15     )
16     model = LinearRegression()
17     model.fit(X_train, y_train)
18     mse = mean_squared_error(y_test, model.predict(X_test))
19     print(f"Seed_{seed}_MSE: {mse}")
```

3.2 Code for Question 2 (Normalization)

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import LinearRegression
4 from sklearn.metrics import mean_squared_error
5 from sklearn.preprocessing import MinMaxScaler, StandardScaler
6
7 df = pd.read_csv('Salary_dataset.csv')
8 X = df[['YearsExperience']]
9 y = df['Salary']
10
11 # Fixed split so we can compare apples to apples
12 X_train, X_test, y_train, y_test = train_test_split(
13     X, y, test_size=0.2, random_state=42
14 )
15
16 # 1. No Normalization
17 model = LinearRegression()
18 model.fit(X_train, y_train)
19 print(f"MSE_None: {mean_squared_error(y_test, model.predict(X_test))}")
20
21 # 2. Min-Max Scaling
22 scaler = MinMaxScaler()
23 X_tr_mm = scaler.fit_transform(X_train)
24 X_te_mm = scaler.transform(X_test)
25 model.fit(X_tr_mm, y_train)
26 print(f"MSE_MinMax: {mean_squared_error(y_test, model.predict(X_te_mm))}")
27
28 # 3. Standard Normalization
29 scaler = StandardScaler()
30 X_tr_std = scaler.fit_transform(X_train)
31 X_te_std = scaler.transform(X_test)
32 model.fit(X_tr_std, y_train)
33 print(f"MSE_Std: {mean_squared_error(y_test, model.predict(X_te_std))}")
```