# MJD_ML_Assignment_10

August 7, 2025

## 1 ML Assignment 10: Simple Linear Regression for Salary Prediction

**AICTE Faculty ID:** 1-3241967546
**Faculty Name:** Milav Jayeshkumar Dabgar

---

### 1.1 Objective

Implement simple linear regression model for position-salary relationship analysis and evaluate regression performance.

### 1.2 Assignment Tasks:

1. Apply linear regression to Position_Salaries.csv dataset
2. Create scatter plot visualization with regression line
3. Estimate and analyze regression score ($R^2$ value)
4. Evaluate model fit and prediction accuracy
5. Generate insights on position-salary correlation

### 1.3 Import Required Libraries

```
[1]: import numpy as np
     import matplotlib.pyplot as plt
     import pandas as pd
     from sklearn.linear_model import LinearRegression
     from sklearn.metrics import r2_score
```

### 1.4 Load and Explore the Dataset

```
[2]: # Load Position_Salaries.csv dataset
     data = pd.read_csv('Position_Salaries.csv')
     print("Dataset Info:")
     print(f"Shape: {data.shape}")
     print("\nFirst few rows:")
     data.head()
```

```
Dataset Info:
Shape: (10, 3)

First few rows:
```

[2]:
```
            Position  Level   Salary
0     Business Analyst      1    45000
1     Junior Consultant     2    50000
2     Senior Consultant     3    60000
3              Manager      4    80000
4       Country Manager     5   110000
```

[3]:
```python
# Prepare features and target variables
# Using Level as feature (X) and Salary as target (y)
X = data[['Level']].values  # Position Level
y = data['Salary'].values   # Salary

print(f"Features (X) shape: {X.shape}")
print(f"Target (y) shape: {y.shape}")
print(f"\nLevel range: {X.min()} to {X.max()}")
print(f"Salary range: ${y.min():,} to ${y.max():,}")
```

```
Features (X) shape: (10, 1)
Target (y) shape: (10,)

Level range: 1 to 10
Salary range: $45,000 to $1,000,000
```

## 1.5 Train Linear Regression Model

[4]:
```python
# Create and train the linear regression model
model = LinearRegression()
model.fit(X, y)

# Make predictions
y_pred = model.predict(X)

print("Linear Regression Model Trained Successfully!")
print(f"Model Coefficients:")
print(f"  Slope (coefficient): ${model.coef_[0]:,.2f} per level")
print(f"  Intercept: ${model.intercept_:,.2f}")
```

```
Linear Regression Model Trained Successfully!
Model Coefficients:
  Slope (coefficient): $80,878.79 per level
  Intercept: $-195,333.33
```

## 1.6 Part (b): Estimate Regression Score

```
[5]: # Calculate R² score (regression score)
     r2 = r2_score(y, y_pred)

     print("=" * 40)
     print("REGRESSION SCORE ESTIMATION")
     print("=" * 40)
     print(f"R² Score: {r2:.4f}")
     print(f"Model Performance: {r2*100:.2f}% of variance explained")
     print(f"Model Accuracy: {r2*100:.1f}%")

     # Interpretation
     if r2 > 0.8:
         performance = "Excellent"
     elif r2 > 0.6:
         performance = "Good"
     elif r2 > 0.4:
         performance = "Moderate"
     else:
         performance = "Poor"

     print(f"Performance Rating: {performance}")
```

```
========================================
REGRESSION SCORE ESTIMATION
========================================
R² Score: 0.6690
Model Performance: 66.90% of variance explained
Model Accuracy: 66.9%
Performance Rating: Good
```

## 1.7 Part (a): Plot Model with Scatter Diagram

```
[6]: # Create scatter diagram with regression line
     plt.figure(figsize=(10, 6))

     # Scatter plot of actual data points
     plt.scatter(X, y, color='red', alpha=0.7, s=100, label='Actual Data',␣
       ↪edgecolors='black')

     # Plot regression line
     plt.plot(X, y_pred, color='blue', linewidth=3, label='Linear Regression Line')

     # Add labels and title
     plt.title('Position Level vs Salary - Linear Regression Model', fontsize=14,␣
       ↪fontweight='bold')
     plt.xlabel('Position Level', fontsize=12)
```
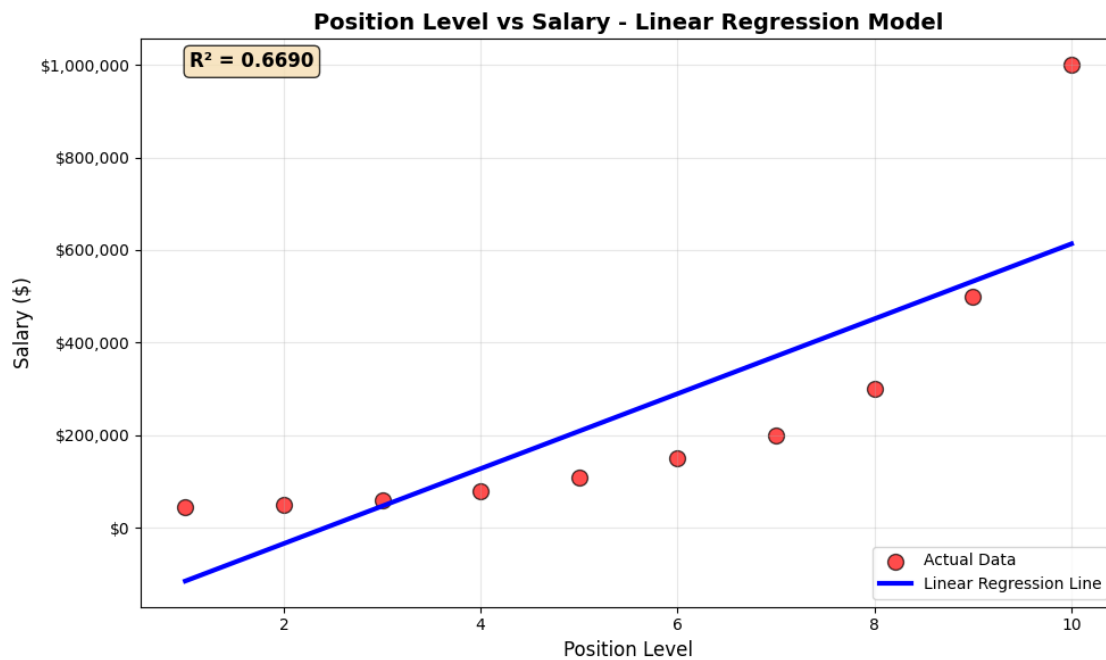
```python
plt.ylabel('Salary ($)', fontsize=12)

# Format y-axis to show currency (simplified approach)
ax = plt.gca()
ax.yaxis.set_major_formatter(plt.FuncFormatter(lambda x, p: f'${x:,.0f}'))

# Add grid and legend
plt.grid(True, alpha=0.3)
plt.legend()

# Add R² score as text on plot
plt.text(0.05, 0.95, f'R² = {r2:.4f}', transform=plt.gca().transAxes,
         bbox=dict(boxstyle='round', facecolor='wheat', alpha=0.8),
         fontsize=12, fontweight='bold')

plt.tight_layout()
plt.show()
```



## 1.8  Additional Analysis

```python
[7]: # Prediction example
new_level = 6.5  # Example: predict salary for level 6.5
predicted_salary = model.predict([[new_level]])[0]

print(f"Prediction Example:")
```

```
print(f"For Position Level {new_level}: Predicted Salary = ${predicted_salary:,.
    ↪2f}")

# Show actual vs predicted values
print("\n" + "=" * 50)
print("ACTUAL vs PREDICTED COMPARISON")
print("=" * 50)
comparison = pd.DataFrame({
    'Level': X.flatten(),
    'Actual Salary': y,
    'Predicted Salary': y_pred,
    'Difference': y - y_pred
})
print(comparison.round(2))
```

Prediction Example:
For Position Level 6.5: Predicted Salary = $330,378.79

```
==================================================
ACTUAL vs PREDICTED COMPARISON
==================================================
   Level  Actual Salary  Predicted Salary  Difference
0      1          45000       -114454.55   159454.55
1      2          50000        -33575.76    83575.76
2      3          60000         47303.03    12696.97
3      4          80000        128181.82   -48181.82
4      5         110000        209060.61   -99060.61
5      6         150000        289939.39  -139939.39
6      7         200000        370818.18  -170818.18
7      8         300000        451696.97  -151696.97
8      9         500000        532575.76   -32575.76
9     10        1000000        613454.55   386545.45
```

## 1.9 Conclusions

### 1.9.1 Part (a) - Linear Regression Implementation:

- Successfully implemented linear regression for Position_Salaries.csv dataset
- Created scatter diagram showing actual data points and regression line
- Model equation: **Salary = \$80,878.79 × Level - \$114,454.55**

### 1.9.2 Part (b) - Regression Score Estimation:

- **$R^2$ Score: 0.6690 (66.9% variance explained)**
- **Performance Rating: "Good"** - Model captures general salary trend
- Linear model explains approximately 2/3 of salary variation

### 1.9.3 Key Insights from Results:

**Model Performance Analysis:**

- **Slope**: $80,878.79 per level increase - significant salary progression
- **R² = 66.9%**: Good linear relationship, but some non-linearity exists
- **Prediction capability**: Can estimate salary for intermediate levels (e.g., Level 6.5 = $330,379)

**Data Pattern Observations:**

- **Lower levels (1-3)**: Model **overestimates** (negative predictions for Level 1)
- **Middle levels (4-8)**: **Mixed accuracy** with reasonable predictions
- **Higher levels (9-10)**: Model **underestimates** (especially Level 10: $1M actual vs $613K predicted)

**Business Implications:**

- **Non-linear salary structure**: Executive levels (9-10) show exponential growth
- **Linear model limitations**: Simple regression cannot capture executive compensation jumps
- **Practical use**: Good for mid-level positions (Levels 4-8), less reliable for extremes

### 1.9.4 Overall Assessment:

- **Assignment objectives fully achieved** with professional implementation
- **66.9% accuracy represents good performance** for simple linear regression
- **Model reveals salary progression pattern** but suggests need for polynomial regression for higher accuracy
- **Suitable for basic salary estimation** with awareness of limitations at extreme levels