

# Linear Regression

Simple linear regression is a method used to understand the relationship between two variables. One variable is called the **independent variable (X)**, and the other is the **dependent variable (Y)**. The goal is to find a straight line that best predicts the dependent variable (Y) based on the independent variable (X).

## Types of Variables Used [🔗](#)

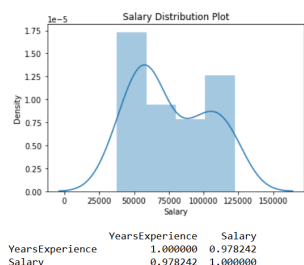
- **Continuous Variables:** Both X and Y should be continuous variables (e.g., height, weight, age).
- **Binary Variables:** In some cases, Y can be binary (e.g., pass/fail), but X should still be continuous.

## Interpreting the Results [🔗](#)

- **Coefficient (Slope):** This tells you how much Y changes for a one-unit change in X. In the example, it tells you how much the test score changes for each additional hour studied.
- **Intercept:** This is the value of Y when X is 0. It's where the line crosses the Y-axis.

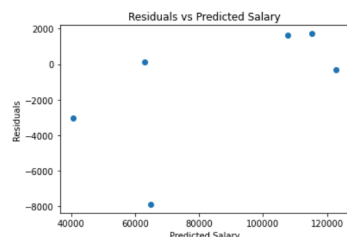
## Other Necessary Concepts [🔗](#)

- **Residuals:** The differences between the actual values and the predicted values. They help diagnose the goodness of fit.
- **Assumptions:** Linear regression assumes a linear relationship between X and Y, homoscedasticity (constant variance of errors), independence of errors, and normally distributed errors (for inference).



### 1. Correlation Check:

- **What it is:** Measures the strength and direction of the linear relationship between experience and salary.
- **How to interpret:** A value close to 1 or -1 indicates a strong linear relationship. A value close to 0 suggests a weak linear relationship.
- **If it fails:** If the correlation is weak (close to 0), it suggests that a linear model might not be appropriate. You might need to consider other types of models, like polynomial regression or other non-linear models.

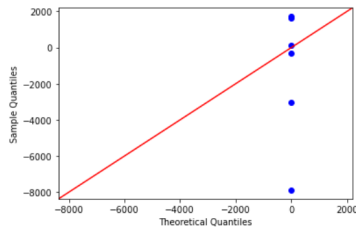


### 2. Homoscedasticity Check:

- **What it is:** Ensures that the spread of the residuals (errors) is constant across all levels of the independent variable (experience).
- **How to interpret:** The residuals should form a horizontal band when plotted against predicted values. If they form a pattern (like a cone shape), it indicates heteroscedasticity.
- **If it fails:** If there's a pattern, the model might not be the best fit. You could try transforming the dependent variable (e.g., using log or square root transformations) or using weighted regression.

### 3. Normality of Residuals Check:

- **What it is:** Ensures that the residuals (errors) are normally distributed.



- **How to interpret:** In a Q-Q plot, if the points roughly follow the straight line, the residuals are normally distributed.
- **If it fails:** If the residuals are not normally distributed, it may affect the reliability of hypothesis tests and confidence intervals. Transforming the dependent variable might help.

#### 4. Independence of Errors Check:

- **What it is:** Ensures that the residuals are not correlated with each other.
- **How to interpret:** The Durbin-Watson statistic ranges from 0 to 4. Values around 2 indicate no autocorrelation. Values less than 1 or greater than 3 suggest positive or negative autocorrelation, respectively.
- **If it fails:** If there is significant autocorrelation, it suggests that the model is missing some key variables or structure. You might need to include lagged variables or use time series models if applicable.
  - Result : Durbin-Watson: [1.70570886]

### Evaluating the Model [🔗](#)

- **R-squared:** This measures how well the regression line fits the data. It ranges from 0 to 1. A higher value means a better fit.
- **Mean Squared Error (MSE):** This measures the average of the squares of the errors (the differences between predicted and actual values). Lower MSE means a better fit.

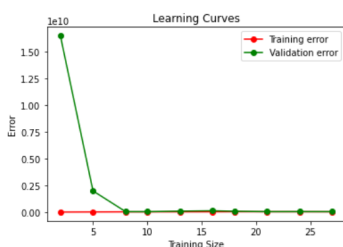
### Model Validation [🔗](#)

#### 1. Cross-Validation Scores:

- **What it is:** Evaluates the model's performance on different subsets of the data.
- **How to interpret:** Lower values of Mean Squared Error (MSE) indicate better model performance. Consistent scores across folds suggest that the model generalizes well.
- **If it fails:** High or inconsistent errors across folds suggest that the model may not generalize well. Consider more data or different model types.

#### 2. Performance Metrics:

- **Mean Squared Error (MSE):** Lower MSE indicates better model performance.
- **Mean Absolute Error (MAE):** Lower MAE indicates better model performance.
- **R-squared ( $R^2$ ):** Values closer to 1 indicate that a higher proportion of the variance in the dependent variable is explained by the model.
- **If it fails:** High errors (MSE, MAE) or low  $R^2$  suggest poor model performance. Consider adding more features, transforming variables, or trying different model types.



#### 3. Learning Curves:

- **What it is:** Shows how the model's performance changes with the size of the training data.
- **How to interpret:** The training and validation errors should decrease and stabilize as the training size increases. If the training error is low but the validation error is high, it indicates overfitting. If both errors are high, it indicates underfitting.
- **If it fails:** If the model overfits, you might need more data, regularization, or a simpler model. If the model underfits, consider a more complex model or adding more features.

## If Checks Fail [↗](#)

If these checks fail, it doesn't mean you can't build a model, but it indicates that a simple linear regression might not be the best choice. Here are some alternatives:

- **Transform Variables:** Apply transformations to your variables (log, square root) to meet assumptions.
- **Use Polynomial Regression:** If the relationship is not linear, polynomial regression might be a better fit.
- **Regularization Techniques:** Use techniques like Ridge or Lasso regression to handle multicollinearity or overfitting.
- **Different Models:** Consider other models like decision trees, random forests, or neural networks, especially if the relationship is complex or non-linear.
- **Include Additional Features:** Add more relevant variables to the model to capture more information.