

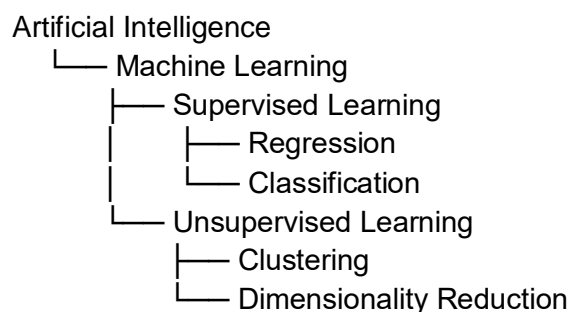
Type	Description	Examples
1 <input type="checkbox"/> <b>Supervised Learning</b>	Model learns from labeled data (data with known outputs)	Regression, Classification
2 <input type="checkbox"/> <b>Unsupervised Learning</b>	Model learns from unlabeled data (no known outputs)	Clustering, Dimensionality Reduction

---

☐ **Inside Supervised Learning**, we have:

Subtype	Output Type	Examples
<b>Regression</b>	Continuous output	Linear Regression, Polynomial Regression
<b>Classification</b>	Categorical output	Logistic Regression, Decision Tree, Random Forest, SVM, etc.

---




## Linear Regression

Linear Regression is a method used to find a linear relationship between independent variable(s) (X) and a dependent variable (Y).

Formula

$$Y = mx + c$$

 If there are multiple variables (multiple linear regression):

$$y = b_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n$$

Symbols	Meaning	Analogy
<b>Y</b>	Dependent variable (the value you're predicting)	e.g., Life Expectancy
<b>X</b>	Independent variable (the input)	e.g., GDP
<b>b<sub>0</sub></b>	Intercept → value of Y when X = 0	Like "c" in $y = mx + c$
<b>b<sub>1</sub></b>	Slope → how much Y changes when X increases by 1 unit	Like "m" in $y = mx + c$
<b>ε</b>	Error term → difference between actual and predicted Y	The leftover or noise

B0 = Intercept(c)

B1 = the slope(m)

Type	Description
Simple Linear Regression	One independent variable
Multiple Linear Regression	More than one independent variable
Polynomial Linear Regression	Uses higher powers of X (like $X^2$ , $X^3$ ) to model curved relationships

m = coefficient(s) → how much Y changes when X changes

c = intercept → value of Y when X = 0

---

**Best fit line** : The **best fit line** is also called the **prediction line**, because it represents the model's predicted values.

The closer the actual (dependent variable) points are to this line, the **better the model fits** the data — meaning the **prediction error is low** and the model is **accurate**.

---

## Regression Evaluation Metrics

Metric	What It Means (Simple Words)	Good Score Means
<b>MSE (Mean Squared Error)</b>	It shows how far your predictions are from the real values — it squares the errors, so big mistakes count more.	<input type="checkbox"/> <b>Lower is better</b> (closer to 0 means your model is accurate)
<b>RMSE (Root Mean Squared Error)</b>	Same as MSE, but in the same units as your target (e.g., years). It tells how much your predictions differ from the real values on average.	<input type="checkbox"/> <b>Lower is better</b>
<b>MAE (Mean Absolute Error)</b>	It's the average of all absolute differences between predicted and actual values — simple average error.	<input type="checkbox"/> <b>Lower is better</b>
<b>R<sup>2</sup> (R-Square)</b>	It tells how much of the variation in the actual data your model can explain. (e.g., R <sup>2</sup> = 0.85 means your model explains 85% of the data.)	<input type="checkbox"/> <b>Higher is better</b> (closer to 1 means more accurate)
<b>Adjusted R<sup>2</sup></b>	Similar to R <sup>2</sup> , but it gives a small penalty when you add useless features. Helps in comparing models with different numbers of predictors.	<input type="checkbox"/> <b>Higher is better</b> (should be close to R <sup>2</sup> , not much lower)

$$\text{Adjusted R square} = 1 - (1 - R^2) * (N - 1) / N - p - 1$$

## Where:

Symbol	Meaning
$R^2$	Normal R-squared value
$N$	Total number of observations (rows)
$p$	Number of independent variables (features/predictors)

## □ Why we need Adjusted $R^2$

You said it perfectly —

**$R^2$  always increases when you add more independent variables**, even if those variables are **not actually useful** in prediction.

□ So,  $R^2$  can **mislead you** — it'll look like your model is improving, but in reality, the model might just be becoming more complex and overfitted.

## □ How Adjusted $R^2$ fixes that

- Adjusted  $R^2$  **penalizes unnecessary variables** (via the term  $p$ ).
- If you add a new variable that **actually helps**, Adjusted  $R^2$  will **increase**.
- If you add a variable that **doesn't help**, Adjusted  $R^2$  will **decrease**.

This makes it a **better, fairer measure** of model performance when multiple predictors are used.

Variable	$R^2$	Adjusted $R^2$
1 Predictor	0.75	0.74
3 Predictors	0.80	0.77

6 Predictors   0.82   0.76

□ Notice: even though  $R^2$  keeps rising,  
Adjusted  $R^2$  starts dropping because not all predictors are useful.

## □ In short:

Metric	Meaning	Problem / Solution
$R^2$	% of variation in Y explained by X	Always increases with new variables
<b>Adjusted <math>R^2</math></b>	Penalizes unnecessary predictors	Only increases if the new variable adds real value

## □ Case 1 – Few useful variables

Model	Independent Variables (p)	What they represent	$R^2$	Adjusted $R^2$
Model 1	1 → <b>engine_size</b>	Bigger engine → higher price	0.70	0.69
Model 2	2 → <b>engine_size, mileage</b>	Mileage also affects price	0.82	0.81
Model 3	3 → <b>engine_size, mileage, brand_rating</b>	Brand also important	0.88	0.87

✓ Here every new p (variable) **adds real information**,  
so  $R^2 \uparrow$  and **Adjusted  $R^2 \uparrow$** .  
The model truly got better.

## □ Case 2 – Adding useless variables

Now you start adding random columns like the color of the dashboard, number of cup holders, or serial number.

Model	Independent Variables (p)	Are they useful?	R <sup>2</sup>	Adjusted R <sup>2</sup>
Model 4	+ dashboard_color	✗ No relation to price	0.885	0.86
Model 5	+ cup_holders, serial_number	✗ Still no relation	0.89	0.84

□ R<sup>2</sup> **keeps increasing a little** (because adding any variable can always fit the data a tiny bit more),  
but **Adjusted R<sup>2</sup> drops** — it's punishing you for adding useless features.

When we add more independent variables (**p increases**),  
the denominator  $N - p - 1$  **decreases**.  
As a result, the fraction value **increases**.

Since this entire fraction is **subtracted from 1**,  
the overall Adjusted R<sup>2</sup> **decreases** — unless the new variable actually improves R<sup>2</sup> a lot.

This means Adjusted R<sup>2</sup> **penalizes** the model for adding too many variables that don't truly help in prediction.

Case 2 = If p decreases denominator  $N - p - 1$  increases and the fraction gets smaller,  
hence R square goes up.

## □ In short:

Change in p	Effect on Denominator	Effect on Fraction	Effect on Adjusted R <sup>2</sup>
p increases	Denominator ↓	Fraction ↑	Adjusted R <sup>2</sup> ↓
p decreases	Denominator ↑	Fraction ↓	Adjusted R <sup>2</sup> ↑

