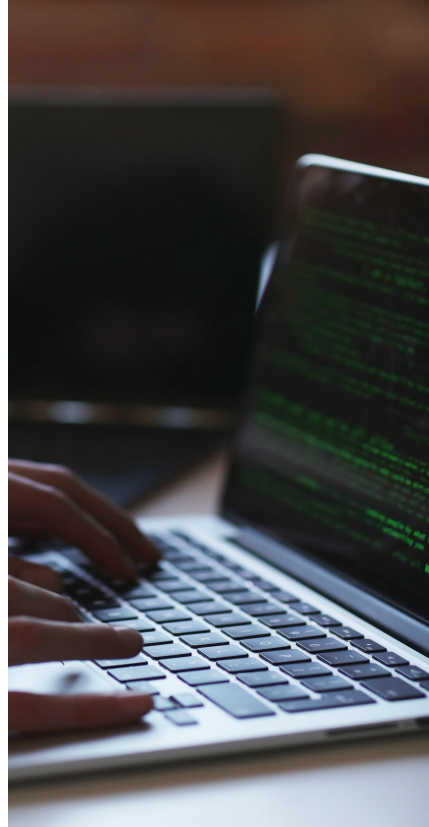


Introduction to Regression





APA YANG AKAN KITA PELAJARI



Regression and applications

- What is Regression
- Regression Intuition
- Regression Algorithm
- Applications of Regression
- Regression Cases

Apa itu Regression Mengapa Kita Pelajari di Data Science ?



Categorical



Type of Data



Numeric

Contoh Data Numeric



House Price
Prediction

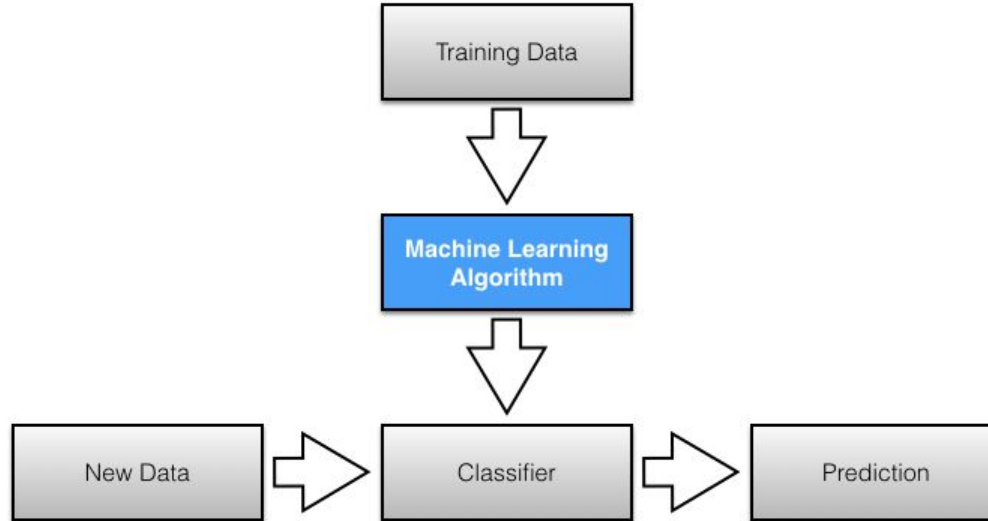


Stock Price
Prediction



Insurance Price
Prediction

Supervised Learning Process



Regression Algorithm



1. Linear Regression
2. Multiple Linear Regression
3. Polynomial Regression
4. **Decision Tree Regression**
5. **Random Forest Regression**
6. **XGBoost Regression**

Regression Intuition



A Glance on Regression Data Example



▲ Suburb	▲ Address	# Rooms	▲ Type	# Price
Abbotsford	49 Lithgow St	3	h	1490000
Abbotsford	59A Turner St	3	h	1220000
Abbotsford	119B Yarra St	3	h	1420000
Aberfeldie	68 Vida St	3	h	1515000
Airport West	92 Clydesdale Rd	2	h	670000

X(Independent) :

- Suburb
- Address
- Rooms
- Type
- ...

y(Dependent):

- Price

Regression Formula



Linear Regression: Single Variable

$$\boxed{\hat{y}} = \beta_0 + \beta_1 \boxed{x} + \boxed{\epsilon}$$

Predicted output Coefficients Input Error

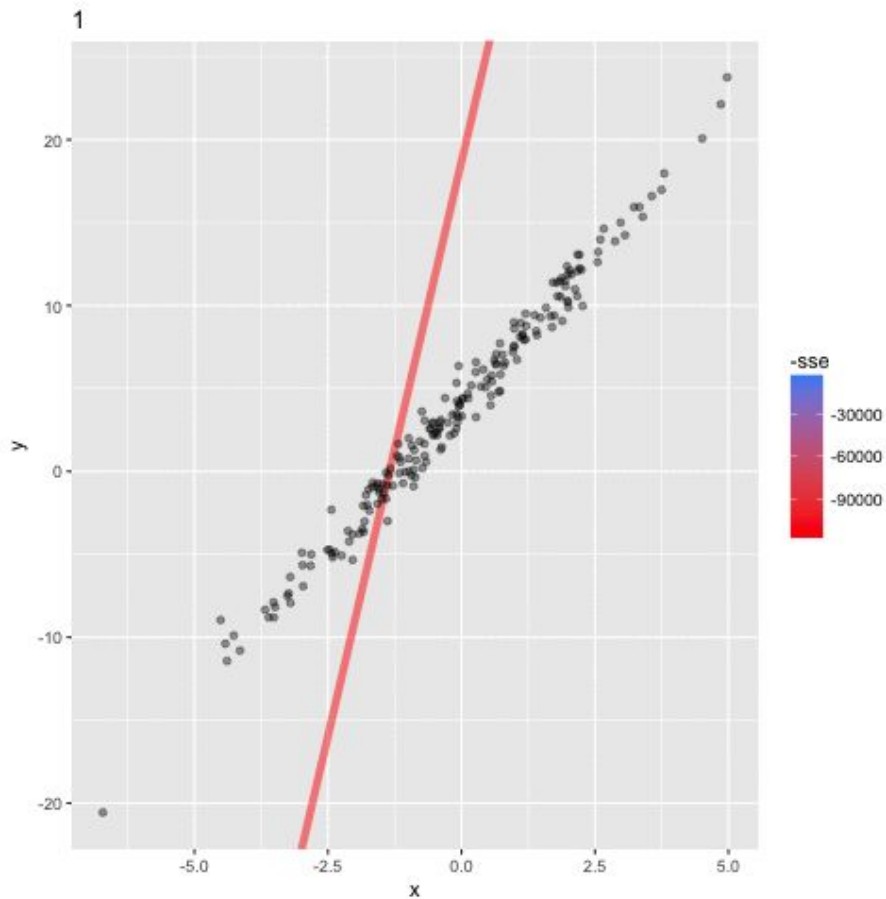
Linear Regression: Multiple Variables

$$\boxed{\hat{y}} = \beta_0 + \beta_1 \boxed{x_1} + \dots + \beta_p \boxed{x_p} + \boxed{\epsilon}$$

Regression Formula



Regression Model	Equation
Simple linear	$Y = a + bX$
Quadratic	$Y = a + bX + bX^2$
Logarithmic	$Y = a + b \log X$
Exponential	$Y = ae^{bx}$ $e = 2.7183$



Cost Function



Cost function is a measure of how wrong the model is in terms of its ability to estimate the relationship between X and y

The objective of a ML model, therefore, is to find parameters, weights or a structure that minimises the cost function.



Cost Function:
$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Goal:
$$\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_0, \theta_1)$$

n = number of data

y_i = actual data (test data)

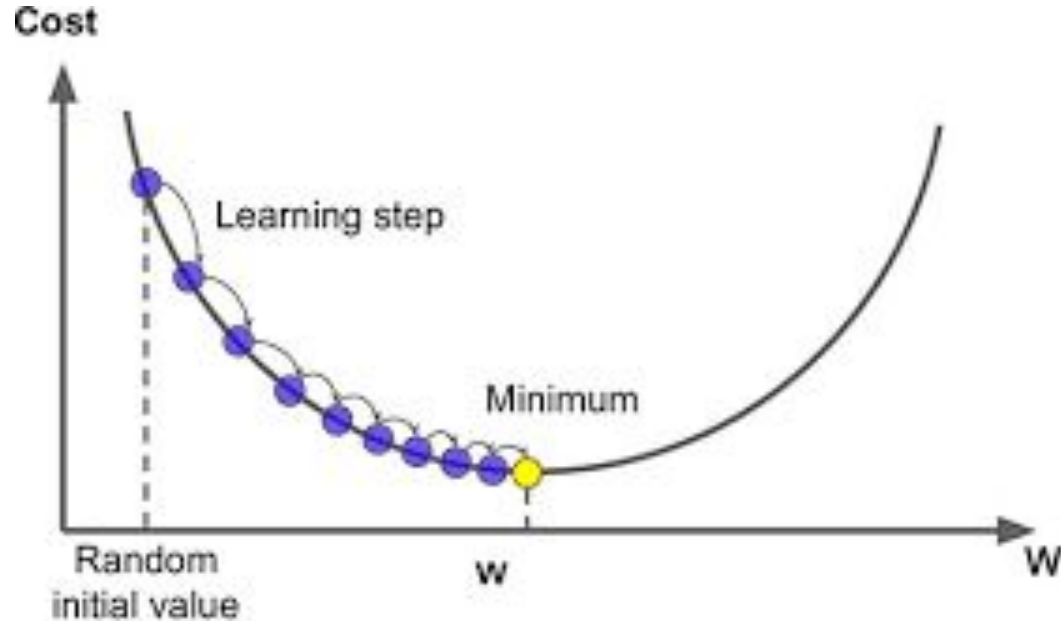
\hat{y}_i = prediction data

Error = Prediction - Actual

How to Minimize Cost Function?



Gradient Descent



Gradient Descent



repeat until convergence {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

}

Regression Evaluation



Data yang akan di-training.

fitur_1	fitur_2	fitur_3	fitur_4	fitur_5	target
...	0.95
...	0.32
...	0.05
...	2.34
...	1.65
...	1.34

Setelah fase training, diperoleh model yang dapat memprediksi target.

fitur_1	fitur_2	fitur_3	fitur_4	fitur_5	target	prediksi_target
...	0.95	0.90
...	0.32	0.32
...	0.05	0.07
...	2.34	1.23
...	1.65	1.64
...	1.34	1.24

Target dan Hasil prediksi dievaluasi.
MAE, MSE, RMSE (Error Based)
 R^2

fitur_1	fitur_2	fitur_3	fitur_4	fitur_5	target	prediksi_target	Evaluasi (selisih)
...	0.95	0.90	0.05
...	0.32	0.32	0.00
...	0.05	0.07	0.02
...	2.34	1.23	1.11
...	1.65	1.64	0.01
...	1.34	1.24	0.01

Regression Evaluation



Metrics for Regression :

1. R-Squared
2. Mean Squared Error (MSE)
3. Root Mean Squared Error (RMSE)
4. Mean Absolute Error (MAE)

R-Squared

$$R^2 = 1 - \frac{\text{sum squared regression (SSR)}}{\text{total sum of squares (SST)}},$$
$$= 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}.$$



Interval Nilai : -inf s.d. 1

buruk baik

Root Mean Squared Error (RMSE)



Interval Nilai : 0 s.d. +inf

baik

buruk

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

n = number of data

y_i = actual data (test data)

\hat{y}_i = prediction data

Error = Prediction - Actual

Mean Absolute Error (MAE)



$$MAE = \frac{1}{n} \sum \left| y - \hat{y} \right|$$

Diagram annotations:

- Divide by the total number of data points (points to $\frac{1}{n}$)
- Actual output value (points to y)
- Predicted output value (points to \hat{y})
- Sum of (points to the summation symbol \sum)
- The absolute value of the residual (points to the absolute value bars $| \cdot |$)

n = number of data
 y_i = actual data (test data)
 y_i_hat = prediction data

Error = Prediction - Actual

Interval Nilai : 0 s.d. +inf

baik **buruk**

Mean Squared Error (MSE)



Interval Nilai : 0 s.d. +inf
 baik **buruk**

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \tilde{y}_i)^2$$

n = number of data

y_i = actual data (test data)

\tilde{y}_i = prediction data

Error = Prediction - Actual