# Linear Regression

- Intro to linear regression

- Intro to least square method

**DAY**

## Content

- Problem

- Linear regression

- Least Square Method

- Variance and Covariance

# Go-Car

- Go-car adalah taxi online yang sedang popular di Indonesia.

- Banyak orang naik Go-car, karena tarifnya sudah jelas di depan. Namun cara menghitung tarif **belum jelas** karena berbagai factor yang digunakan Go-car dalam menentukan tarifnya.

- Faktor utama adalah **Jarak**, namun ada faktor lain yang tidak kita ketahui yang juga mempengaruhi (mungkin **jumlah permintaan**, **cuaca**, **kemacetan**, dll).*





*Adopted from Machine Learning course IF-TelU

▶ "JalanJalan" adalah **perusahaan travel** yang sangat sering menggunakan jasa go-car untuk mengantar tamu-tamu mereka. Karena mereka harus menghitung anggaran bulanan, maka mereka harus **menghitung perkiraan biaya Go-car** jauh hari sebelum digunakan, sehingga pengecekan langsung pada aplikasi go-car dianggap bukan cara yang tepat.

▶ "JalanJalan" memiliki DATA history biaya go-car selama beberapa bulan terakhir.

▶ Pertanyaan: Bagaimana cara kita untuk menentukan MODEL dari data history biaya go-car untuk memprediksi tarif bagi perusahaan "JalanJalan" tersebut?
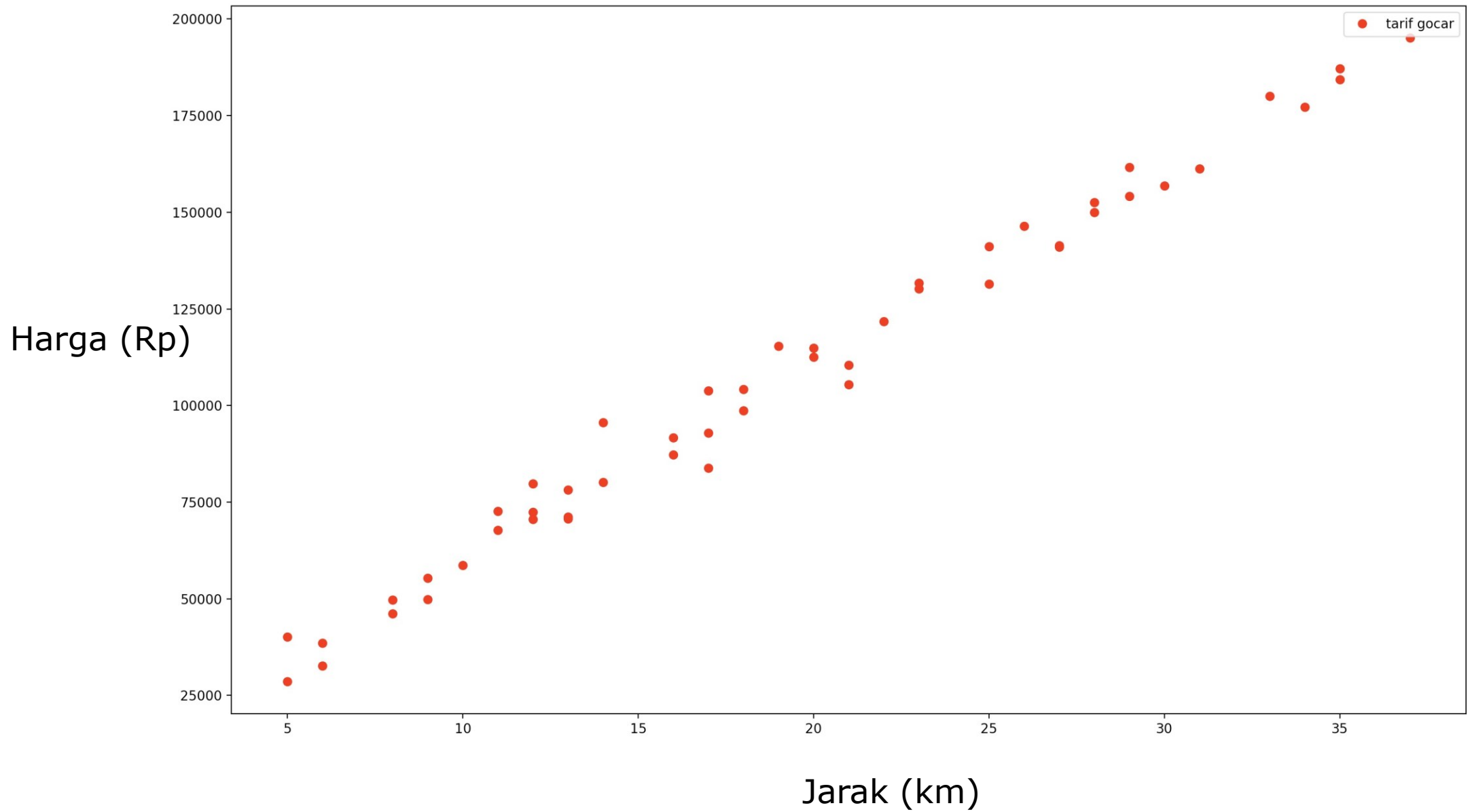
*Adopted from Machine Learning course IF-TelU

# DATA historis Go-car

‣ Sebagai contoh data go-car dari 15 kali naik Go-car sebagai berikut:

‣ What can you do with this history data ?

| NO | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| KM | 5 | 5 | 6 | 6 | 8 | 8 | 9 | 9 | 10 | 11 | 11 | 12 | 12 | 12 | 13 |
| TARIF | 28600 | 40100 | 32600 | 38600 | 49700 | 46100 | 55300 | 49800 | 58700 | 67700 | 72600 | 79800 | 72400 | 70600 | 70700 |

Data plot

# Persamaan Garis?

- **y = a + mx**

- **Cari m**

- m = slope = $(y_2-y_1)/(x_2-x_1)$

- = (160000-34000)/(30-5)

- = 126000/25 = 5040

- **Cari Persamaan Garis dari 1 titik**

- m = $(y-y_1)/(x-x_1)$

- y = $y_1+m(x-x_1)$

- y = 34000+5040(x-5)

- y = 34000 + 5040x − 25200

- **y = 8800 + 5040x**

❯ Linear regression is a linear approach for modeling the relationship between a scalar dependent variable $y$ and one or more explanatory variables (or independent variables) denoted $X$

❯ The goal is to make quantitative (real valued) predictions on the basis of a (vector of) features or attributes

❯ We consider the case $x \in \mathbb{R}^d$ throughout this chapter

❯ Modelling the relation that best fit the data using a single line (linear function)

- $\beta_0$ : Population Y-Intercept (intercept, bias, ...)
- $\beta_1$ : Population slope (weight vector,...)
- $\varepsilon$ : Random error

$$y = \beta_0 + \beta_1 x + \varepsilon \qquad (1)$$

❯ Also often written

- $f(\boldsymbol{x}) = \boldsymbol{w} \cdot \boldsymbol{x} = \Sigma_{j=1}^{d} w_j x_j$
- $f(\boldsymbol{x}) = \boldsymbol{w} \cdot \boldsymbol{x} + a$

(1a)

(1b)

# Linear Regression

› Linear regression assumes that…

- The relationship between X and Y is **linear**

- Y is distributed normally at each value of X

- The variance of Y at every value of X is the same (homogeneity of variances)

- The **observations are independent**

› The learning problem is to determine the parameters $w$ and $a$ based on data

‣ $X$ is one-dimensional case (d = 1)

  – $X \in \mathbb{R}^1$ , $\boldsymbol{y} \in \mathbb{R}^1$,

‣ If we have $n$ data available,

  – $X = [n \times 1]$, $\boldsymbol{y} = [n \times 1]$
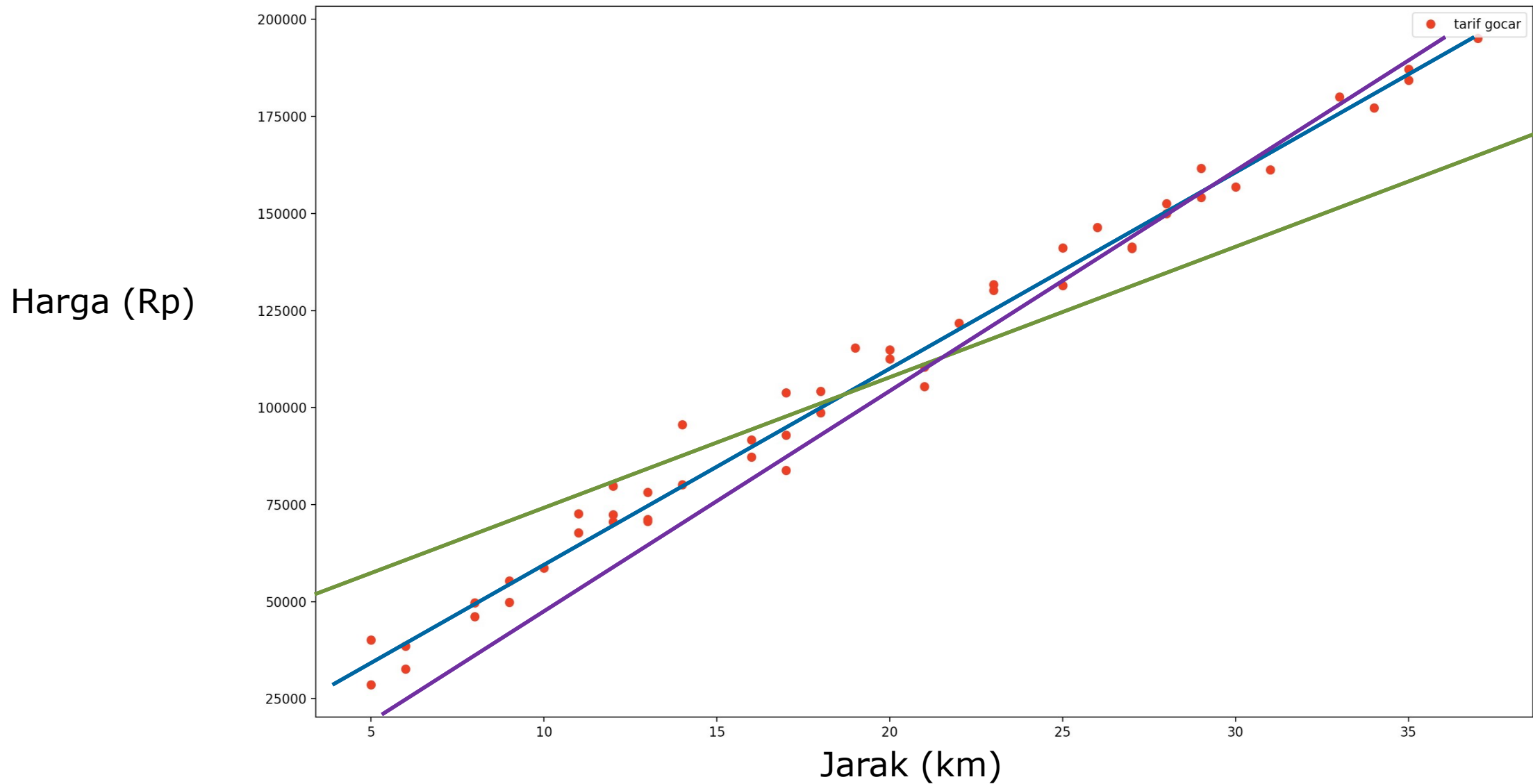
$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad \boldsymbol{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$
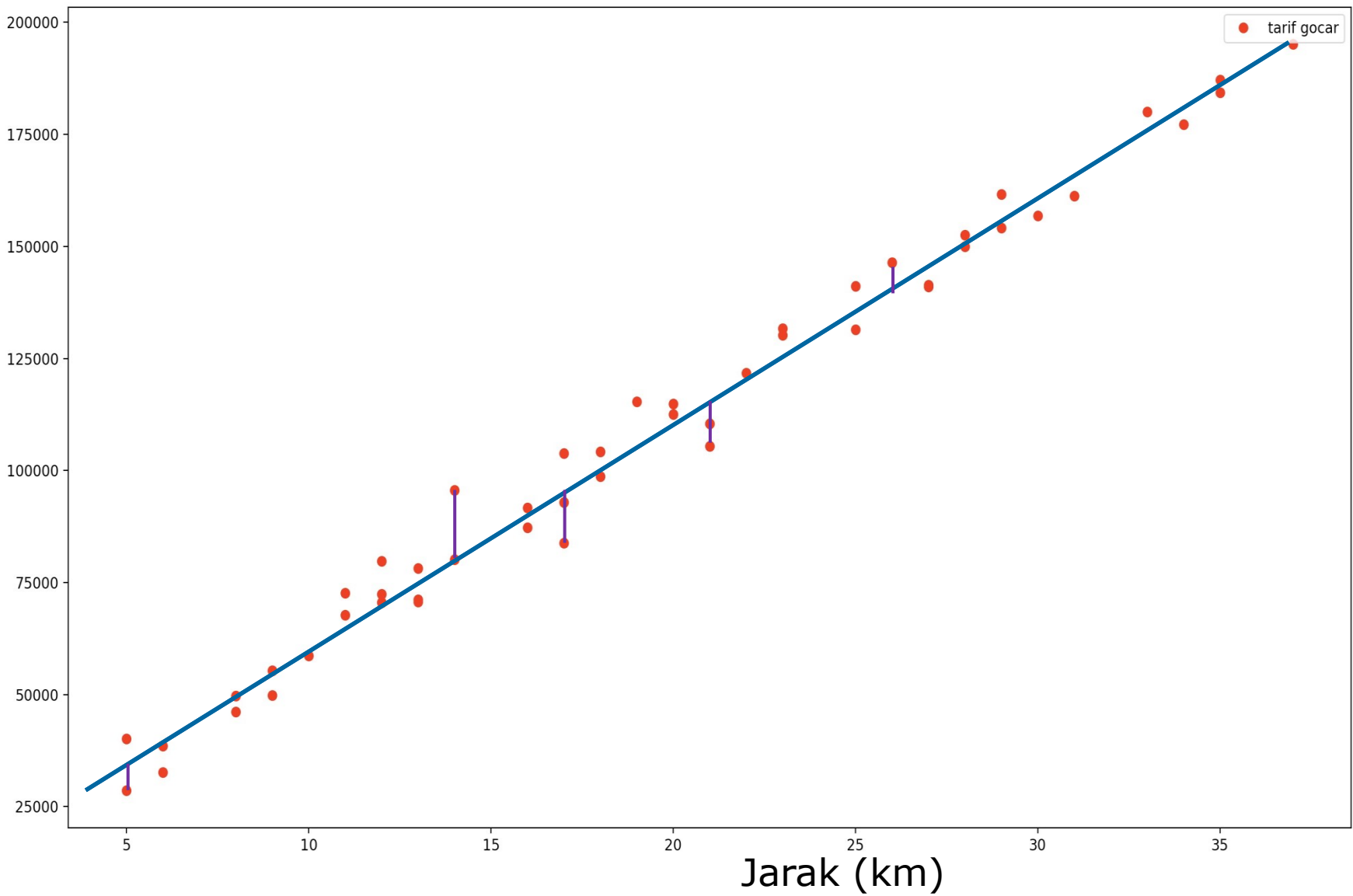
‣ Then the regression for each $x_i$ data point:

$$f(x) = w\,x + a \qquad\qquad w, a \in \mathbb{R}\ (scalar)$$
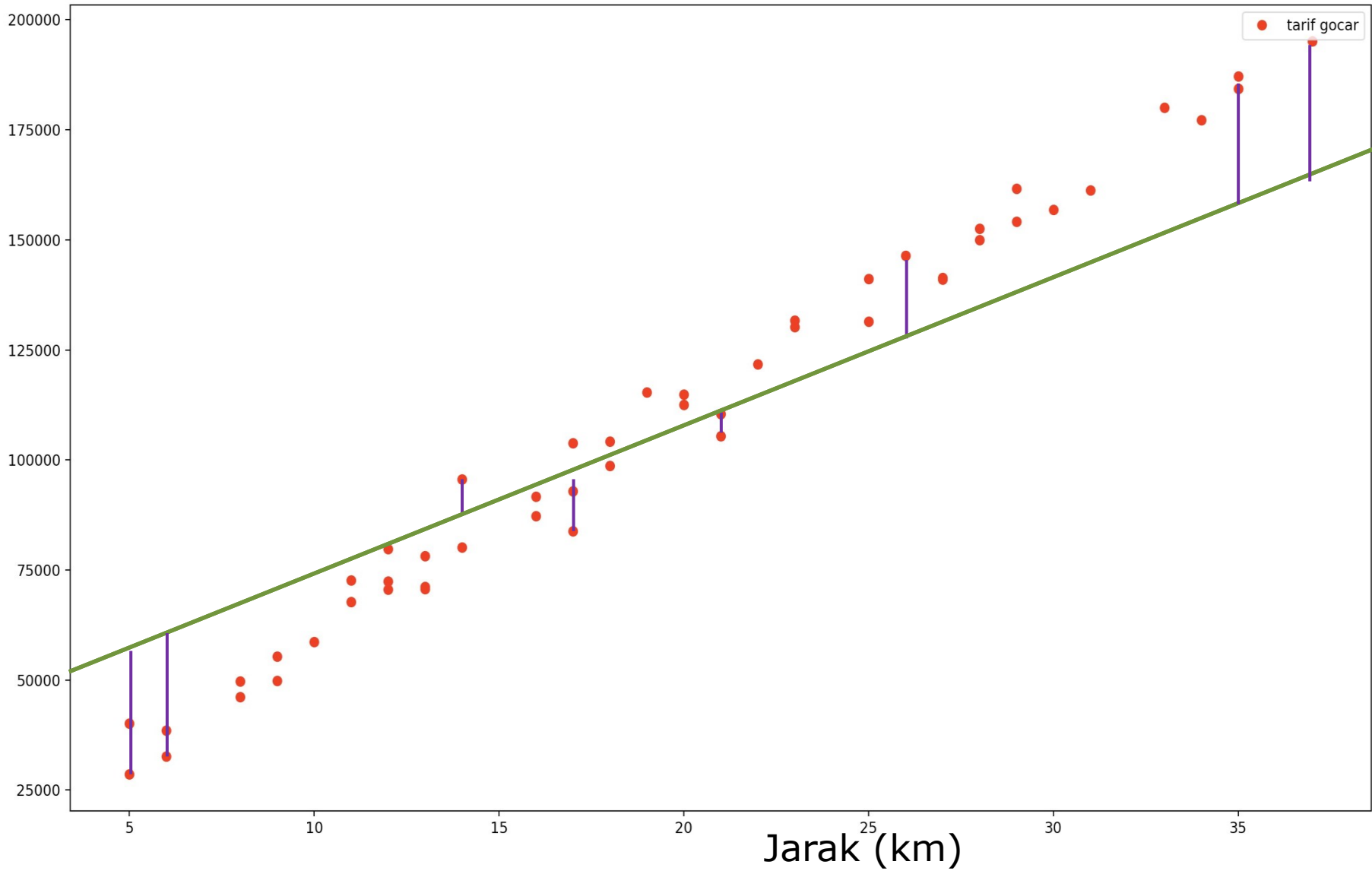
# Performance Measure : How to choose The Best line ?

Harga (Rp)

Jarak (km)

tarif gocar

**Error model-1**

Harga (Rp)

Jarak (km)

tarif gocar

**Error model-2**

Harga (Rp)

Which one is better??

Good Model ??

➢ Smallest error
➢ Which error?

- **y = 8800 + 5040x,** y adalah tarif **prediksi**

- Karena kita punya data Jarak dan Tarif (**sebenarnya**), yang juga biasa dinotasikan dengan x dan y, maka ada 2 y.

- Bedakan

- $\hat{y}$ : Nilai Prediksi (*predicted*)

- y : Nilai Sebenarnya (*observed*)

- Maka Error untuk data pertama adalah :

- $\varepsilon_1 = y_1 - \hat{y}_1$

# Error untuk semua ?

‣ Jarak/Gap antara *observed* dengan *predicted*

‣ **$\varepsilon_1 = y_1 − \hat{y}_1$, jika $y_1 > \hat{y}_1$ positive, jika $\hat{y}_1 > y_1$ negative**

‣ Error negative ? + Error Positive ? = 0 Error ? **Something is wrong**

‣ **$\varepsilon_1 = (y_1 − \hat{y}_1)^2$**

‣ **Total error/Sum Square Error (SSE)= $\varepsilon_1 + \varepsilon_2 + ... + \varepsilon_n$**

‣ **Performance Measure = SSE, semakin kecil semakin baik**

HOW??                    This week: Least Square Method!

›

› **(SSE)= $\varepsilon_1$ +$\varepsilon_2$ + ... +$\varepsilon_n$**

– Find $w$ and $a$ that minimize this function

$$E(w, a) = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 \qquad (2)$$

# Intermezzo:

ARE you familiar with:
➢ Minimization problem??
➢ Derivative??
➢ Find critical point of a function??

❯ Minimization ➜ derivative

$$E(w, a) = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

❯ Solve the minimization problem by setting
the partial derivatives to zero                    (2)

   – Denote the solution by $(\hat{w}, \hat{a})$

$$E(w, a) = \sum_{i=1}^{n} (y_i - wx_i - a)^2$$                    (2b)

› From the first derivative of (2b) w.r.t $a$, We have

$$\frac{\partial E(w,a)}{\partial a} = -2 \sum_{i=1}^{n}(y_i - wx_i - a)$$

› and setting (3) to zero gives

- And $\bar{y} = \frac{1}{n}\Sigma_i y_i$ and $\bar{x} = \frac{1}{n}\Sigma_i x_i$

(3)

$$\hat{a} = \bar{y} - w\bar{x}$$

(4)

Detail:

$0 = -2\sum_{i=1}^{n}(y_i - wx_i - a)$

$0 = \sum_{i=1}^{n}(y_i) - w\sum_{i=1}^{n}(x_i) - an$

$0 = \frac{1}{n}\Sigma_i y_i - w\frac{1}{n}\Sigma_i x_i - a$

❯ From the first derivative of (2b) w.r.t $w$, We have

❯ Plugging in $a = \hat{a}$ in (5) and setting the derivative to zero gives us

(5)

$$\frac{\partial E(w, a)}{\partial w} = -2 \sum_{i=1}^{n} x_i(y_i - wx_i - a)$$

$$0 = \sum_{i=1}^{n} x_i(y_i - wx_i - \bar{y} + w\bar{x})$$

(5b)

❯ For which we can solve (5b) to

$$\widehat{w} = \frac{\sum_{i=1}^{N} x_i (y_i - \bar{y})}{\sum_{i=1}^{N} x_i (x_i - \bar{x})}$$

(6)

Detail:

$$0 = \sum_{i=1}^{n} x_i (y_i - wx_i - \bar{y} + w\bar{x})$$

$$0 = \sum_{i=1}^{n} x_i (y_i - \bar{y} + w\bar{x} - wx_i)$$

$$0 = \sum_{i=1}^{n} x_i (y_i - \bar{y} + w(x_i - \bar{x}))$$

$$0 = \sum_{i=1}^{n} x_i (y_i - \bar{y}) - \sum_{i=1}^{n} wx_i (x_i - \bar{x})$$

›  Challenge : ubah (6) jadi (7)

$$\widehat{w} = \frac{\sum_{i=1}^{N}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{N}(x_i - \bar{x})^2}$$

(7)

❯ Since

$$cov(X, Y) = \sum_{i}^{N} \frac{(x_i - \bar{X})(y_i - \bar{Y})}{N - 1}$$

❯ Notice that from (8), we have

(0b)

$$\widehat{w} = \frac{cov(X, Y)}{cov(X, X)}$$

(9)