

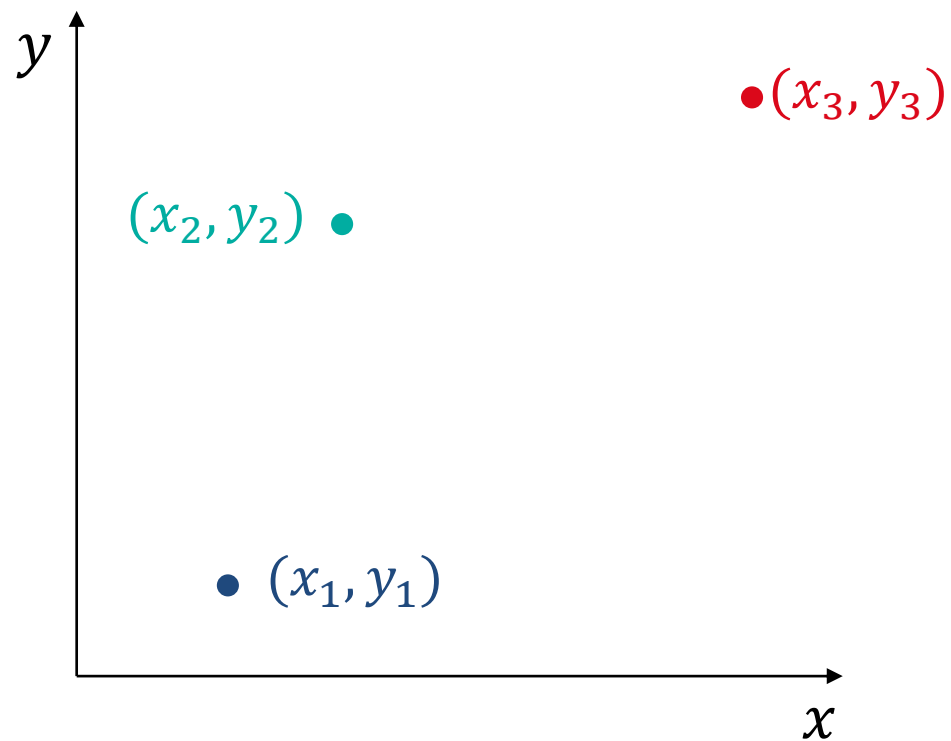
Linear regression

Eung-Hee Kim

ehkim@sunmoon.ac.kr

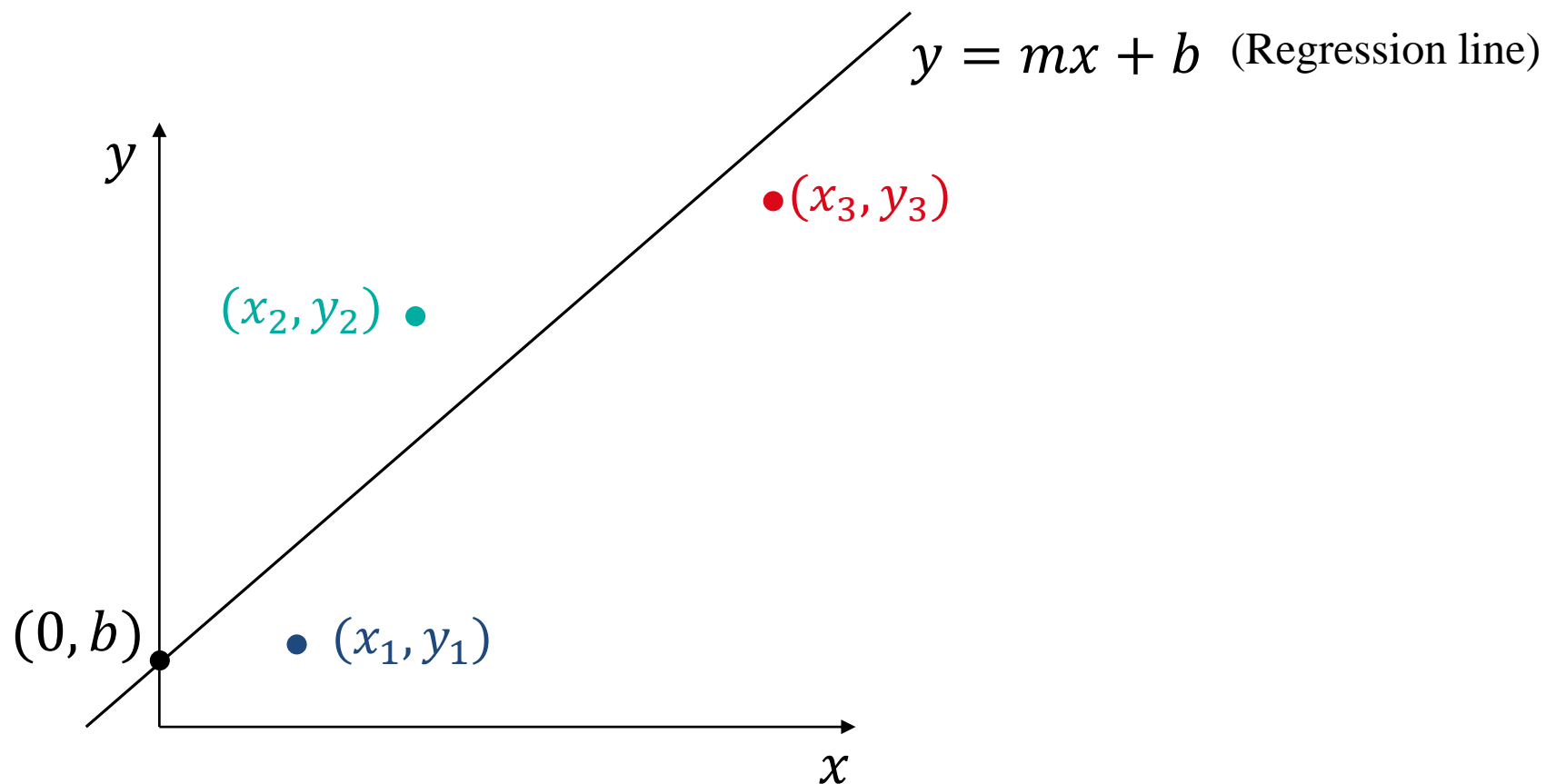
In almost every textbook dealing with the topics related to ‘linear regression’,

- Let x and y be
 - x : independent variable (독립 변수: feature in machine learning)
 - y : dependent variable (종속 변수: indicator or class in machine learning)



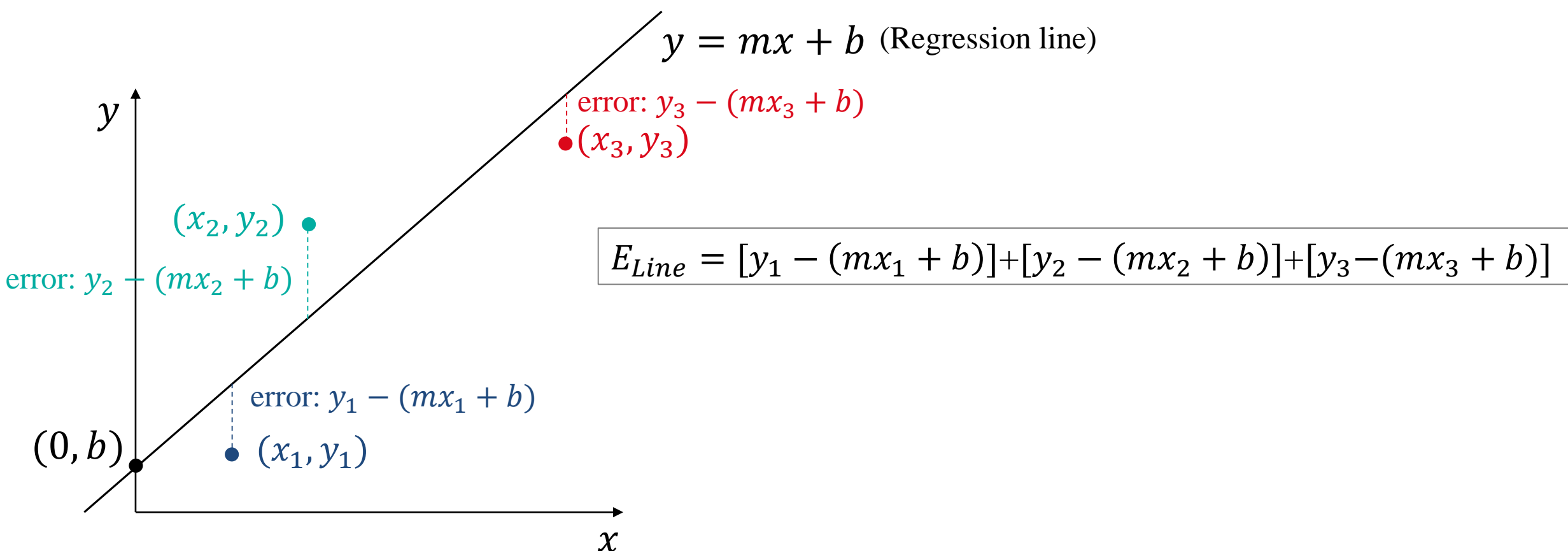
In almost every textbook dealing with the topics related to ‘linear regression’,

- Let x and y be
 - x : independent variable (독립 변수: feature in machine learning)
 - y : dependent variable (종속 변수: indicator or class in machine learning)



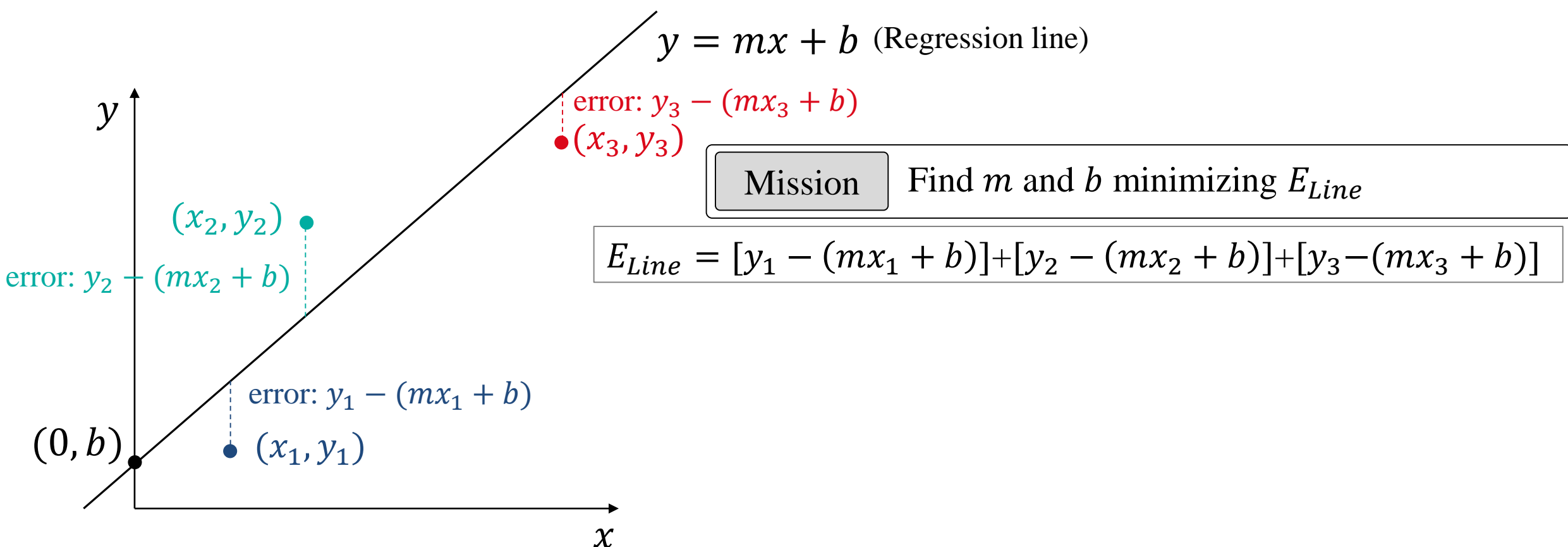
In almost every textbook dealing with the topics related to 'linear regression',

- Let x and y be
 - x : independent variable (독립 변수: feature in machine learning)
 - y : dependent variable (종속 변수: indicator or class in machine learning)



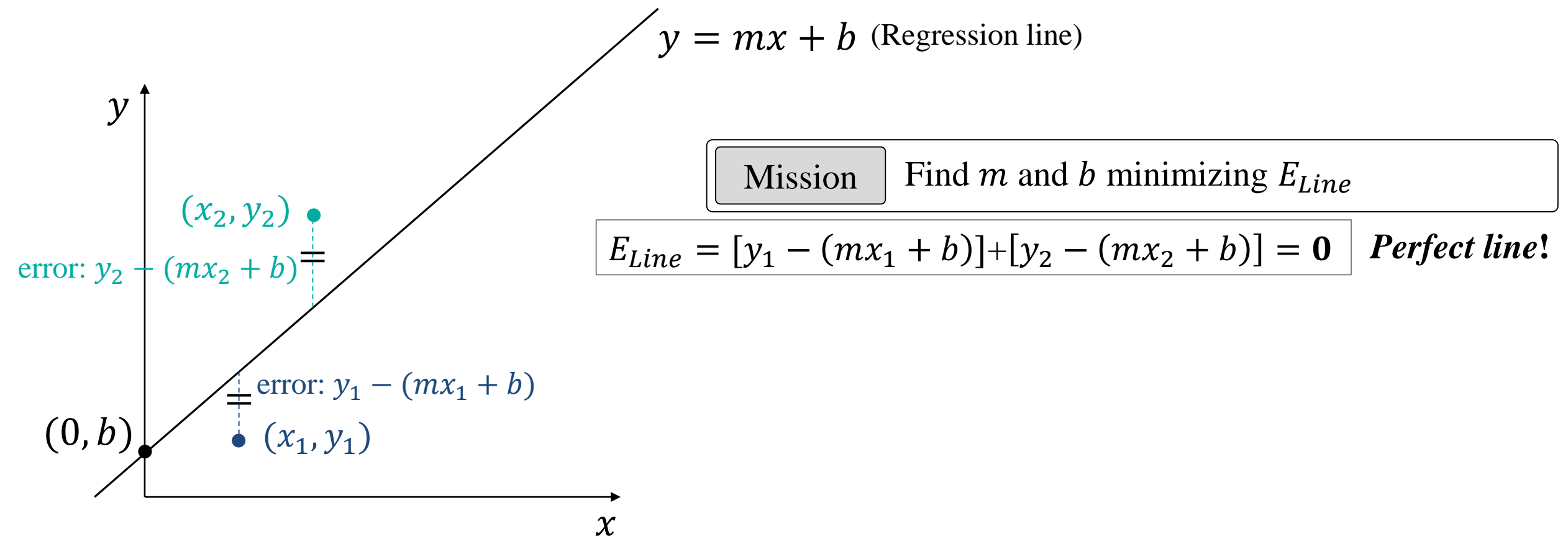
In almost every textbook dealing with the topics related to 'linear regression',

- Let x and y be
 - x : independent variable (독립 변수: feature in machine learning)
 - y : dependent variable (종속 변수: indicator or class in machine learning)



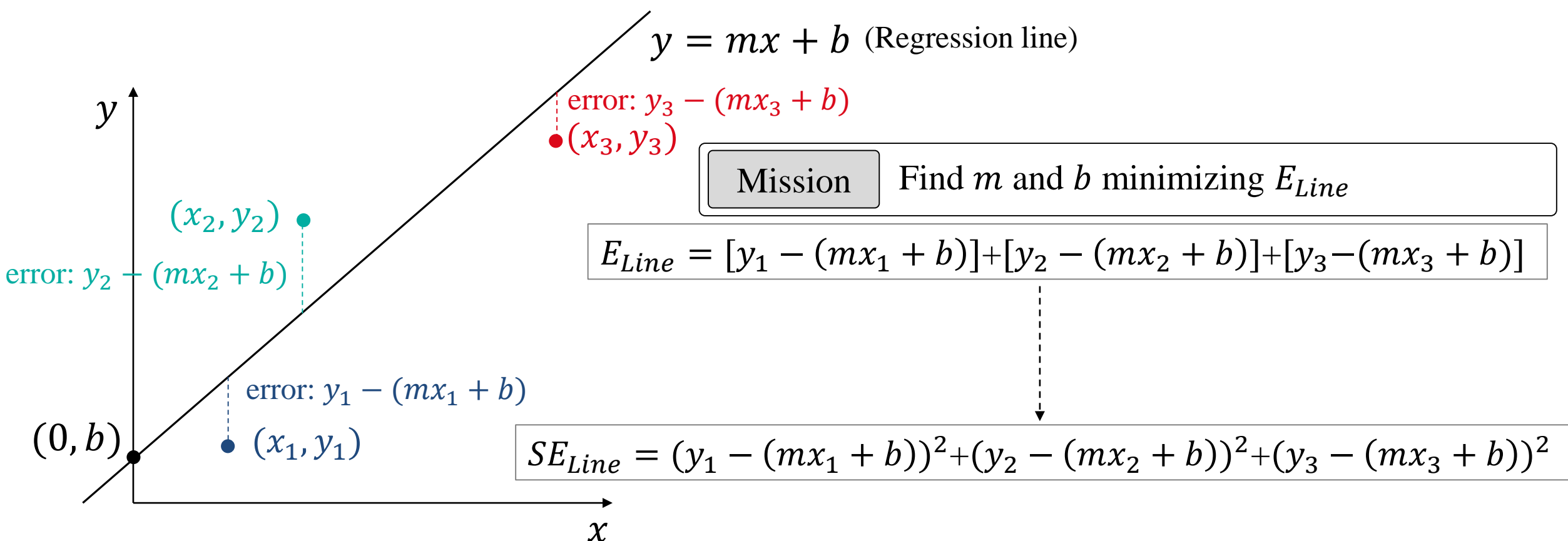
In almost every textbook dealing with the topics related to ‘linear regression’,

- Let x and y be
 - x : independent variable (독립 변수: feature in machine learning)
 - y : dependent variable (종속 변수: indicator or class in machine learning)



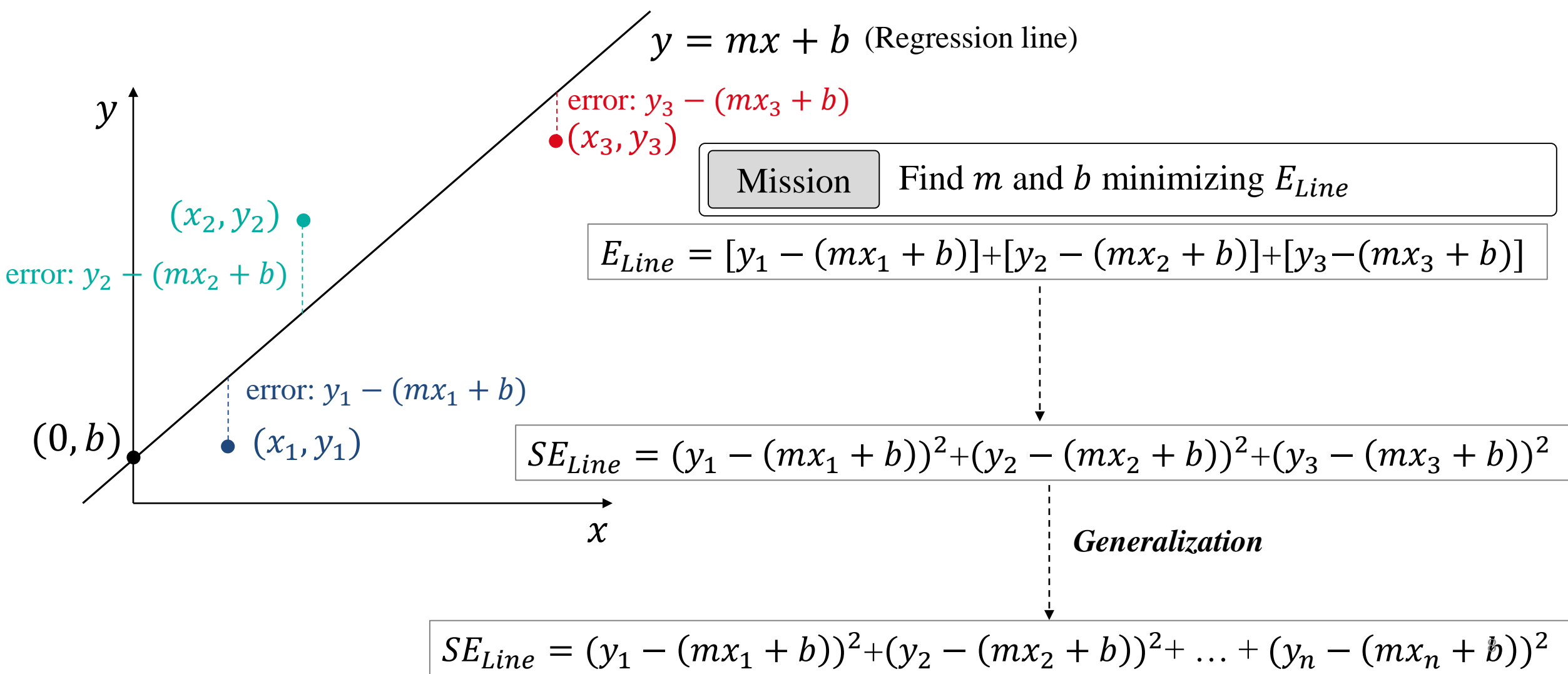
In almost every textbook dealing with the topics related to 'linear regression',

- Let x and y be
 - x : independent variable (독립 변수: feature in machine learning)
 - y : dependent variable (종속 변수: indicator or class in machine learning)



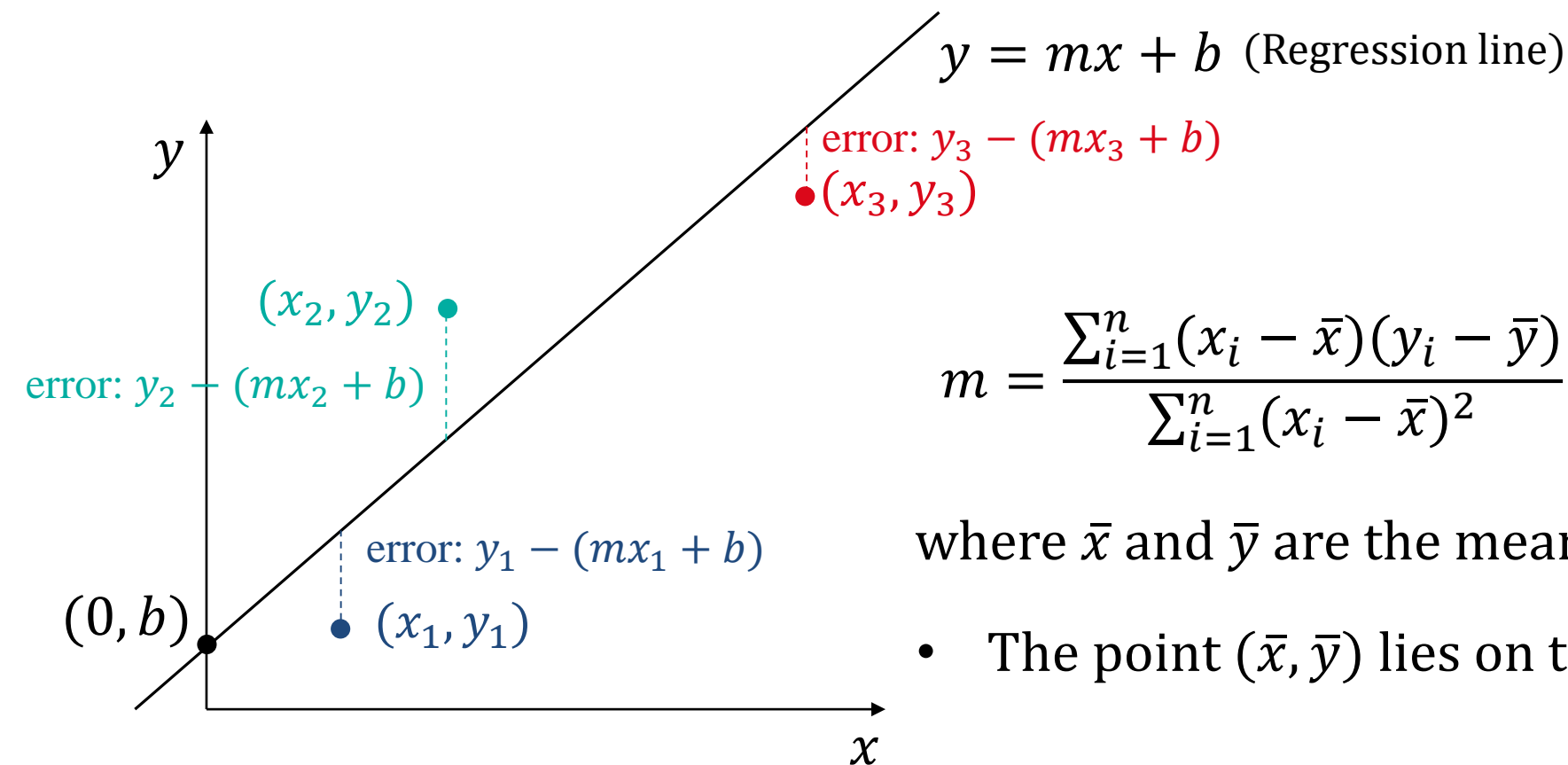
In almost every textbook dealing with the topics related to 'linear regression',

- Let x and y be
 - x : independent variable (독립 변수: feature in machine learning)
 - y : dependent variable (종속 변수: indicator or class in machine learning)



In almost every textbook dealing with the topics related to 'linear regression',

- Let x and y be
 - x : independent variable (독립 변수: feature in machine learning)
 - y : dependent variable (종속 변수: indicator or class in machine learning)



$$m = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \text{and} \quad b = \bar{y} - m\bar{x}$$

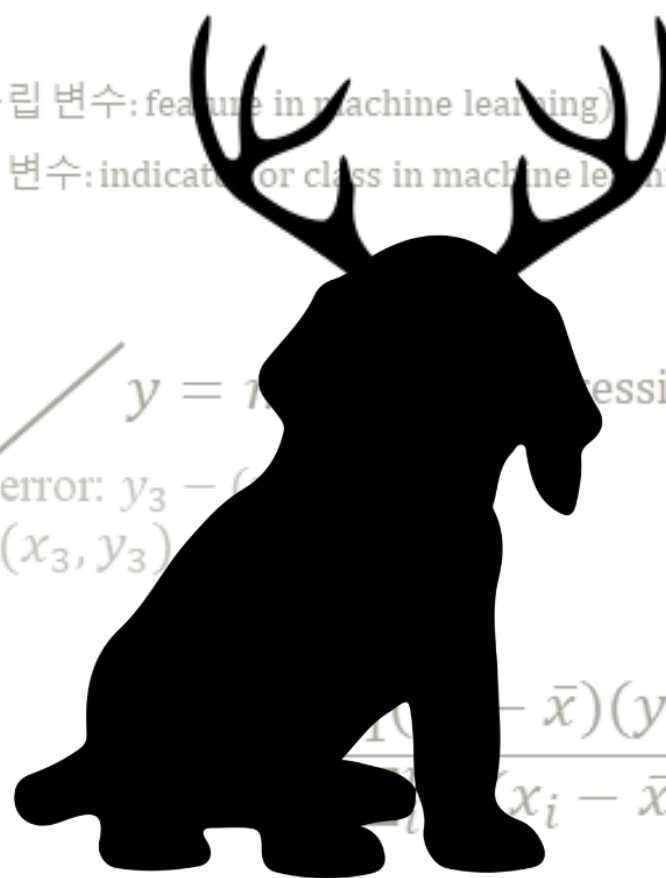
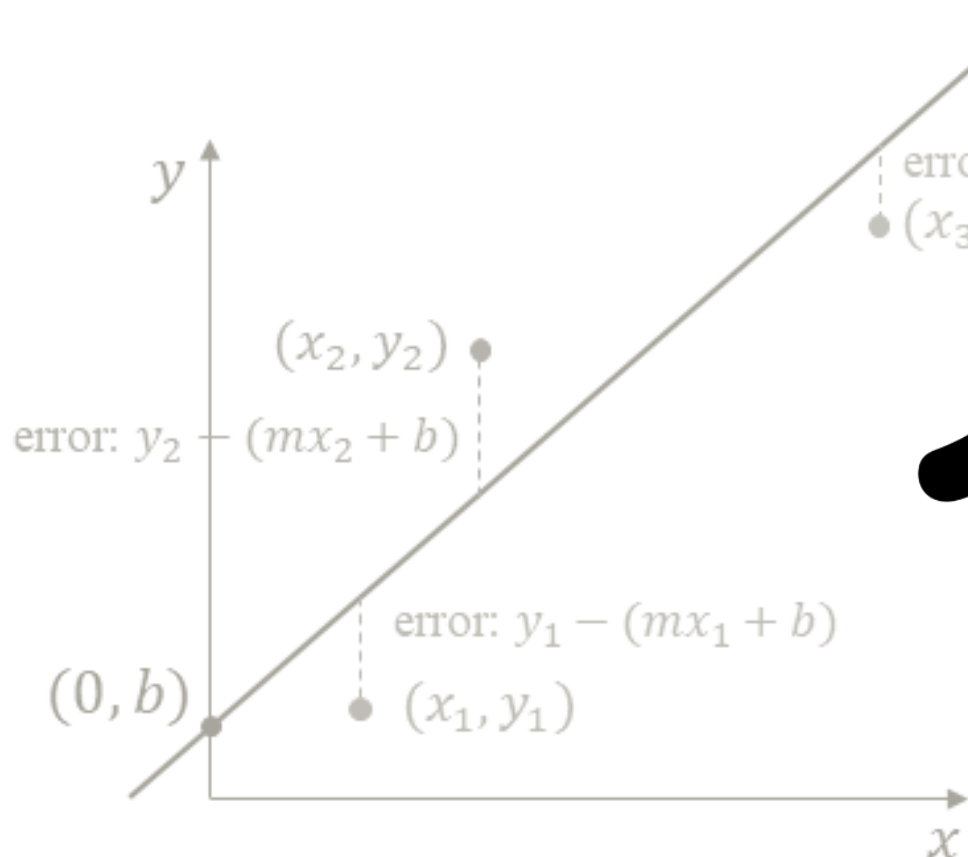
where \bar{x} and \bar{y} are the mean of x and y respectively.

- The point (\bar{x}, \bar{y}) lies on the regression line.



In almost every textbook dealing with the topics related to 'linear regression',

- Let x and y be
 - x : independent variable (독립 변수: feature in machine learning)
 - y : dependent variable (종속 변수: indicator or class in machine learning)



$y = mx + b$ (regression line)

error: $y_3 - (mx_3 + b)$
 (x_3, y_3)

(x_2, y_2)
error: $y_2 - (mx_2 + b)$

error: $y_1 - (mx_1 + b)$

$(0, b)$
 (x_1, y_1)

$$m = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad \text{and} \quad b = \bar{y} - m\bar{x}$$

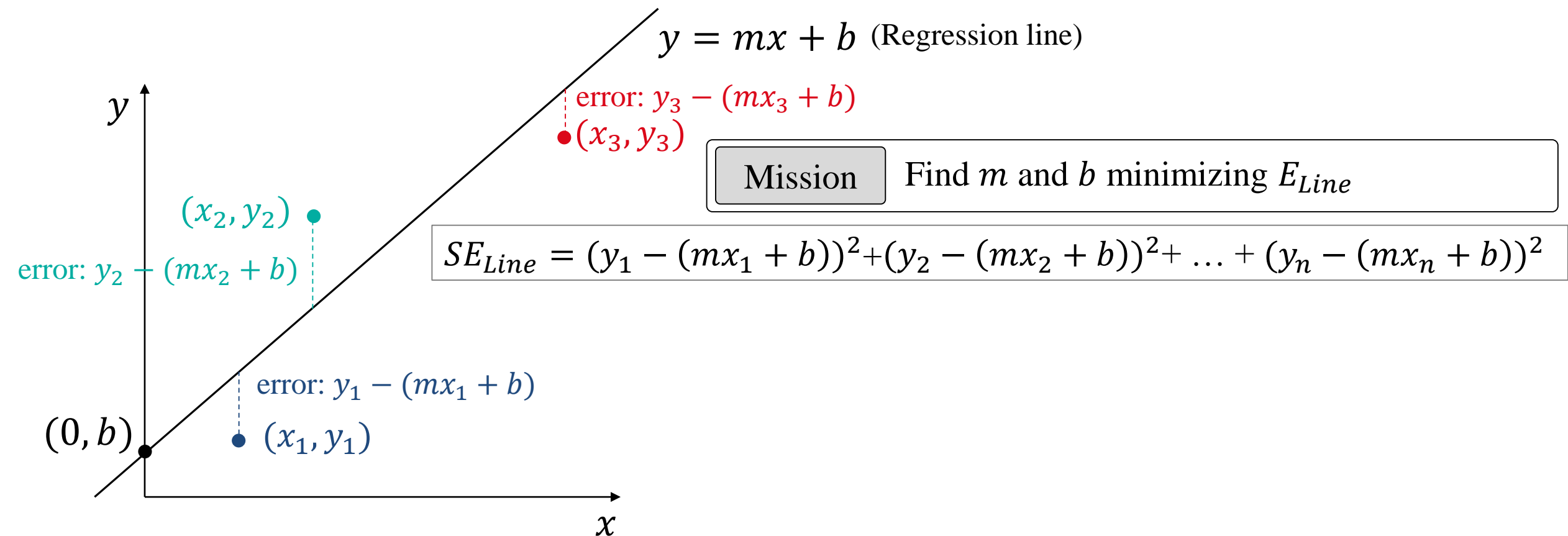
where \bar{x} and \bar{y} are the mean of x and y respectively.

The point (\bar{x}, \bar{y}) lies on the regression line.

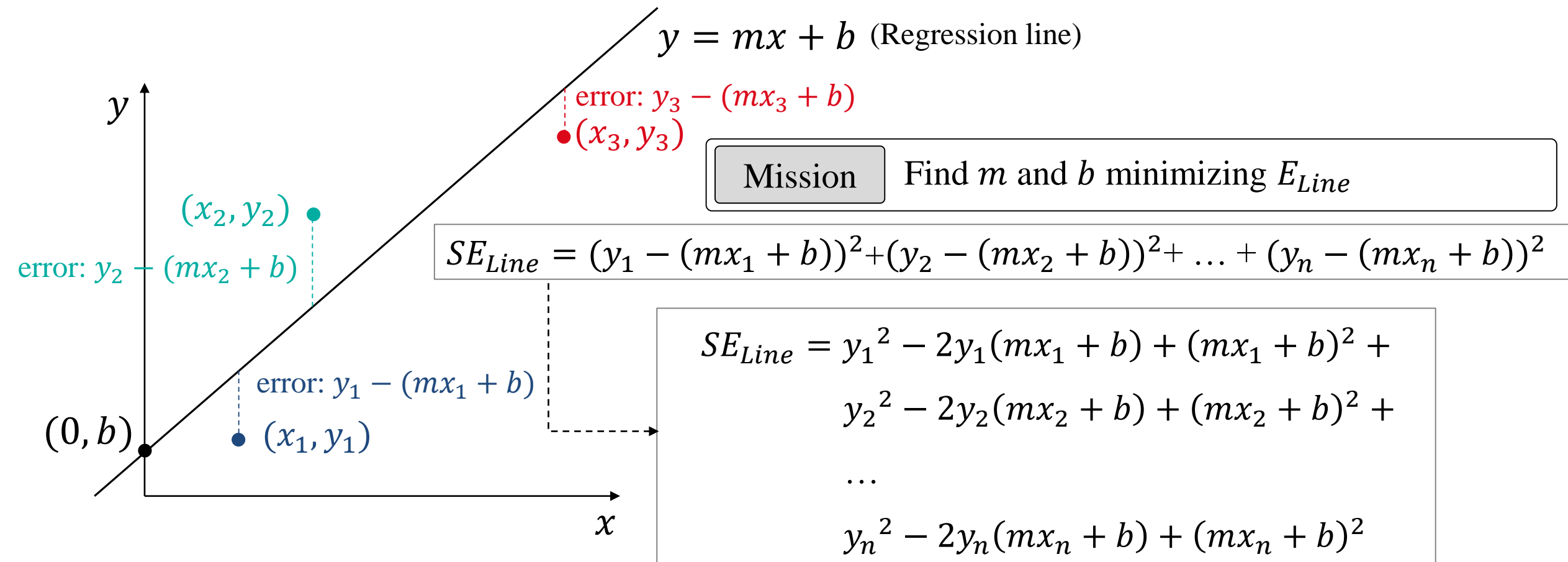


Proof

- Let x and y be
 - x : independent variable (독립 변수: feature in machine learning)
 - y : dependent variable (종속 변수: indicator or class in machine learning)

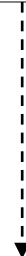


Proof (cont.)



Proof (cont.)

$$\begin{aligned} SE_{Line} = & y_1^2 - 2y_1(mx_1 + b) + (mx_1 + b)^2 + \\ & y_2^2 - 2y_2(mx_2 + b) + (mx_2 + b)^2 + \\ & \dots \\ & y_n^2 - 2y_n(mx_n + b) + (mx_n + b)^2 \end{aligned}$$



$$\begin{aligned} SE_{Line} = & y_1^2 - 2y_1mx_1 - 2y_1b + m^2x_1^2 + 2mx_1b + b^2 + \\ & y_2^2 - 2y_2mx_2 - 2y_2b + m^2x_2^2 + 2mx_2b + b^2 + \\ & \dots \\ & y_n^2 - 2y_nmx_n - 2y_nb + m^2x_n^2 + 2mx_nb + b^2 \end{aligned}$$

Proof (cont.)

$$\begin{aligned} SE_{Line} = & y_1^2 - 2y_1mx_1 - 2y_1b + m^2x_1^2 + 2mx_1b + b^2 + \\ & y_2^2 - 2y_2mx_2 - 2y_2b + m^2x_2^2 + 2mx_2b + b^2 + \\ & \dots \\ & y_n^2 - 2y_nmx_n - 2y_nb + m^2x_n^2 + 2mx_nb + b^2 \end{aligned}$$



$$\begin{aligned} SE_{Line} = & (y_1^2 + y_2^2 + \dots + y_n^2) - 2m(y_1x_1 + y_2x_2 + \dots + y_nx_n) - 2b(y_1 + y_2 + \dots + y_n) + \\ & m^2(x_1^2 + x_2^2 + \dots + x_n^2) + 2mb(x_1 + x_2 + \dots + x_n) + nb^2 \end{aligned}$$

Proof (cont.)

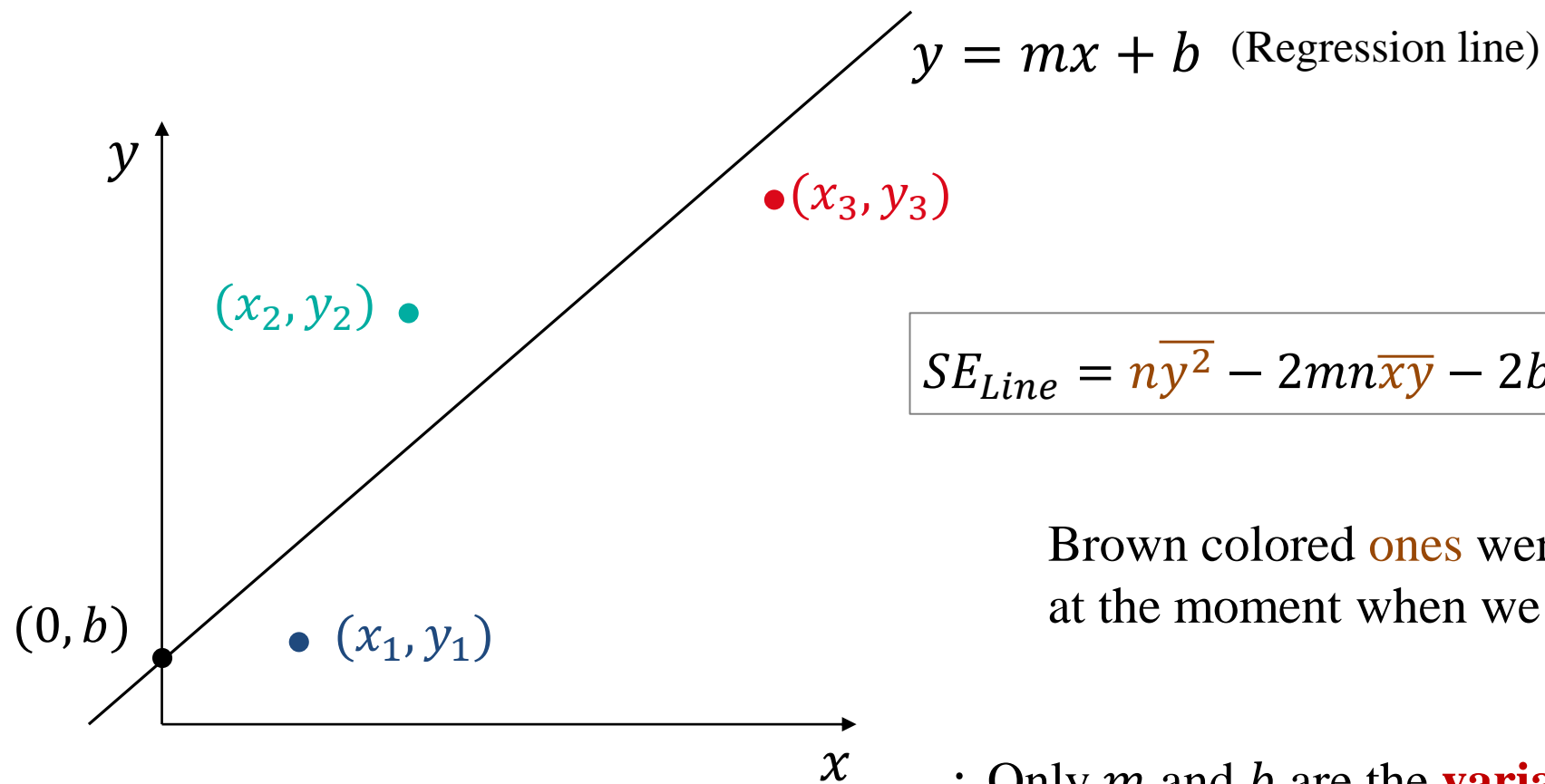
- Notation for average: $\bar{\alpha} = \frac{\alpha_1 + \alpha_2 + \dots + \alpha_n}{n}$
- $\therefore \alpha_1 + \alpha_2 + \dots + \alpha_n = n\bar{\alpha}$

$$SE_{Line} = (y_1^2 + y_2^2 + \dots + y_n^2) - 2m(y_1x_1 + y_2x_2 + \dots + y_nx_n) - 2b(y_1 + y_2 + \dots + y_n) + m^2(x_1^2 + x_2^2 + \dots + x_n^2) + 2mb(x_1 + x_2 + \dots + x_n) + nb^2$$



$$SE_{Line} = n\bar{y}^2 - 2mn\bar{x}\bar{y} - 2bn\bar{y} + m^2n\bar{x}^2 + 2mbn\bar{x} + nb^2$$

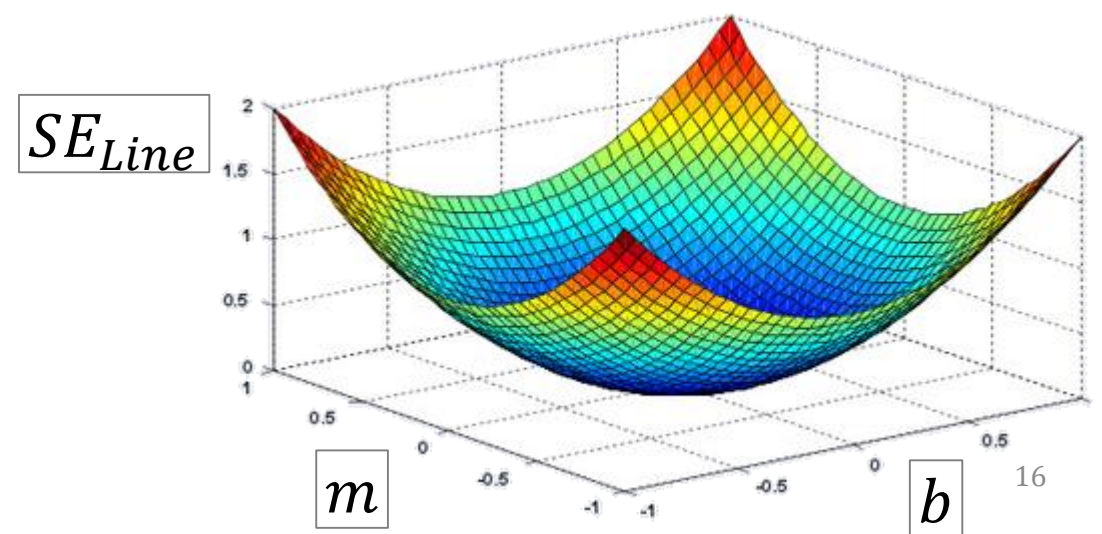
Proof (cont.)



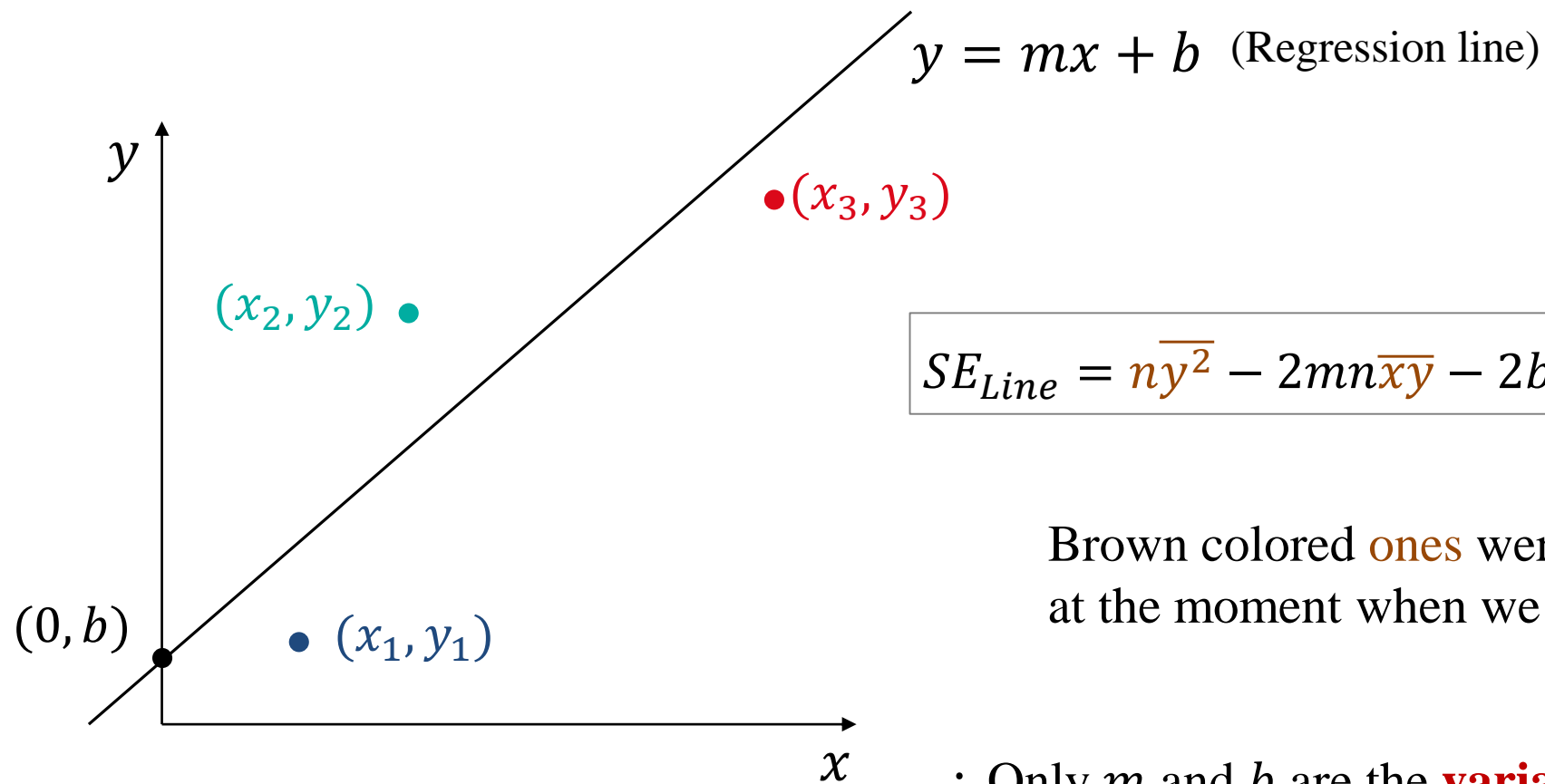
$$SE_{Line} = \overbrace{ny^2} - 2mn\overbrace{xy} - 2bn\overbrace{y} + m^2n\overbrace{x^2} + 2mbn\overbrace{x} + \overbrace{nb^2}$$

Brown colored **ones** were **decided** (constants)
at the moment when we got **our data**.

\therefore Only m and b are the **variables** of the function SE_{Line} .



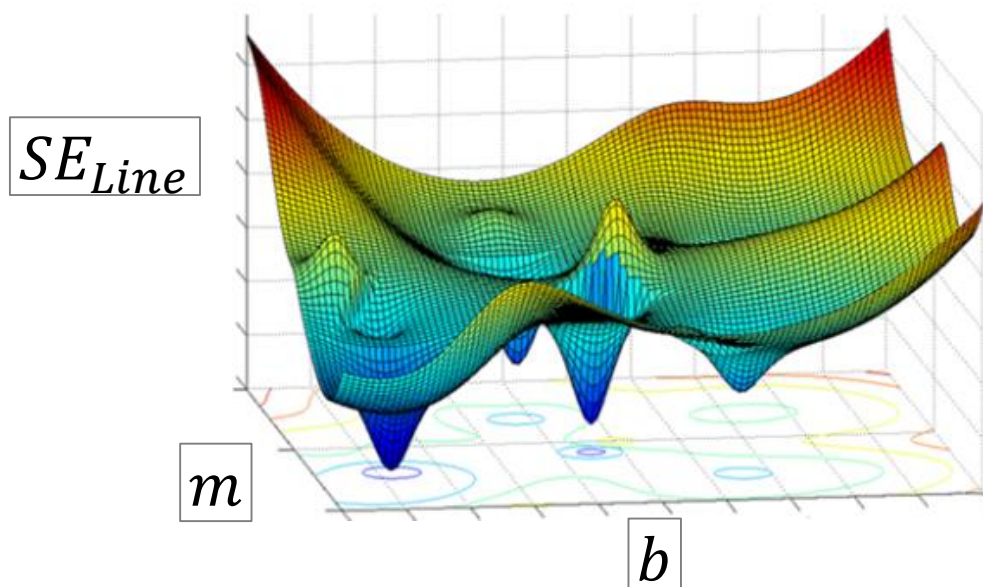
Proof (cont.)



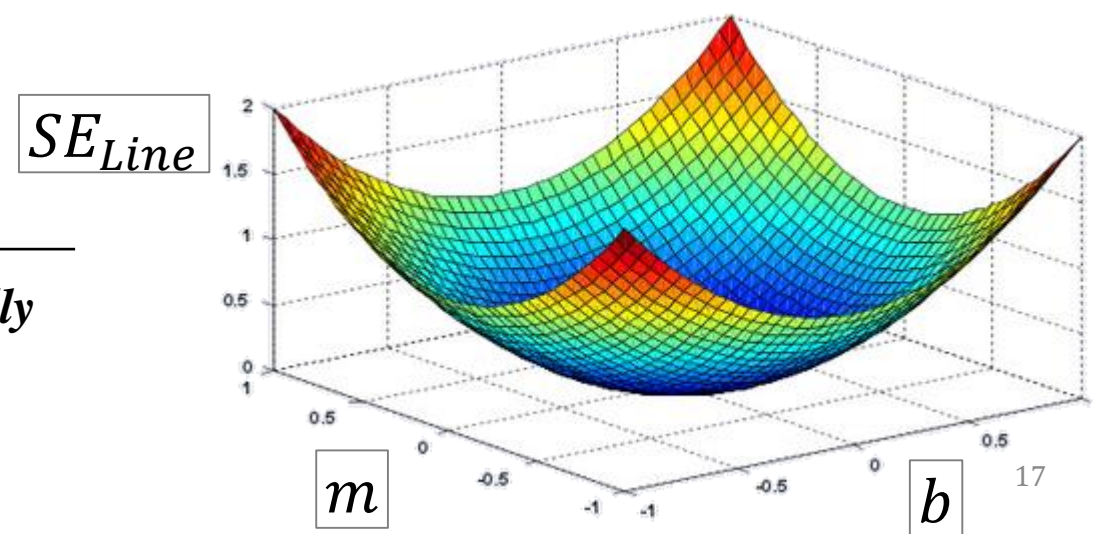
$$SE_{Line} = \overbrace{ny^2} - 2mn\overbrace{xy} - 2bn\overbrace{y} + m^2n\overbrace{x^2} + 2mbn\overbrace{x} + \overbrace{nb^2}$$

Brown colored **ones** were **decided** (actual numbers) at the moment when we got **our data**.

\therefore Only m and b are the **variables** of the function SE_{Line} .



\leftarrow
theoretically possible



Proof (cont.)

$$SE_{Line} = n\bar{y}^2 - 2mn\bar{x}\bar{y} - 2bn\bar{y} + m^2n\bar{x}^2 + 2mbn\bar{x} + nb^2$$

$$\frac{\partial SE_{Line}}{\partial m} = 0$$

$$\frac{\partial SE_{Line}}{\partial b} = 0$$

Proof (cont.)

$$SE_{Line} = n\bar{y}^2 - 2mn\bar{x}\bar{y} - 2bn\bar{y} + m^2n\bar{x}^2 + 2mbn\bar{x} + nb^2$$

$$\frac{\partial SE_{Line}}{\partial m} = 0$$

$$\frac{\partial SE_{Line}}{\partial b} = 0$$

$$\frac{\partial SE_{Line}}{\partial m} = 0 \equiv -2n\bar{x}\bar{y} + 2mn\bar{x}^2 + 2bn\bar{x} = 0$$

$$\frac{\partial SE_{Line}}{\partial b} = 0 \equiv -2n\bar{y} + 2mn\bar{x} + 2nb = 0$$

Proof (cont.)

$$SE_{Line} = n\bar{y}^2 - 2mn\bar{x}\bar{y} - 2bn\bar{y} + m^2n\bar{x}^2 + 2mbn\bar{x} + nb^2$$

$$\frac{\partial SE_{Line}}{\partial m} = 0$$

$$\frac{\partial SE_{Line}}{\partial b} = 0$$

$$\frac{\partial SE_{Line}}{\partial m} = 0 \equiv -2n\bar{x}\bar{y} + 2mn\bar{x}^2 + 2bn\bar{x} = 0$$

$$\frac{\partial SE_{Line}}{\partial b} = 0 \equiv -2n\bar{y} + 2mn\bar{x} + 2nb = 0$$

$$\frac{\partial SE_{Line}}{\partial m} = 0 \equiv -\bar{x}\bar{y} + m\bar{x}^2 + b\bar{x} = 0$$

$$\frac{\partial SE_{Line}}{\partial b} = 0 \equiv -\bar{y} + m\bar{x} + b = 0$$

Proof (cont.)

$$SE_{Line} = n\overline{y^2} - 2mn\overline{xy} - 2bn\overline{y} + m^2n\overline{x^2} + 2mbn\overline{x} + nb^2$$

$$\frac{\partial SE_{Line}}{\partial m} = 0 \equiv -\overline{xy} + m\overline{x^2} + b\overline{x} = 0$$

$$\frac{\partial SE_{Line}}{\partial b} = 0 \equiv -\overline{y} + m\overline{x} + b = 0$$

$$\overline{xy} - \overline{xy} + m\overline{x^2} + b\overline{x} = \overline{xy} + 0$$

$$m\overline{x^2} + b\overline{x} = \overline{xy}$$

Proof (cont.)

$$SE_{Line} = n\overline{y}^2 - 2mn\overline{x}\overline{y} - 2bn\overline{y} + m^2n\overline{x}^2 + 2mbn\overline{x} + nb^2$$

$$\frac{\partial SE_{Line}}{\partial m} = 0 \equiv -\overline{x}\overline{y} + m\overline{x}^2 + b\overline{x} = 0$$

$$\overline{x}\overline{y} - \overline{x}\overline{y} + m\overline{x}^2 + b\overline{x} = \overline{x}\overline{y} + 0$$

$$m\overline{x}^2 + b\overline{x} = \overline{x}\overline{y}$$

$$\frac{\partial SE_{Line}}{\partial b} = 0 \equiv -\overline{y} + m\overline{x} + b = 0$$

$$\overline{y} - \overline{y} + m\overline{x} + b = \overline{y} + 0$$

$$m\overline{x} + b = \overline{y}$$

Proof (cont.)

$$SE_{Line} = n\overline{y^2} - 2mn\overline{xy} - 2bn\bar{y} + m^2n\overline{x^2} + 2mbn\bar{x} + nb^2$$

$$m\overline{x^2} + b\bar{x} = \overline{xy}$$

$$m\bar{x} + b = \bar{y}$$

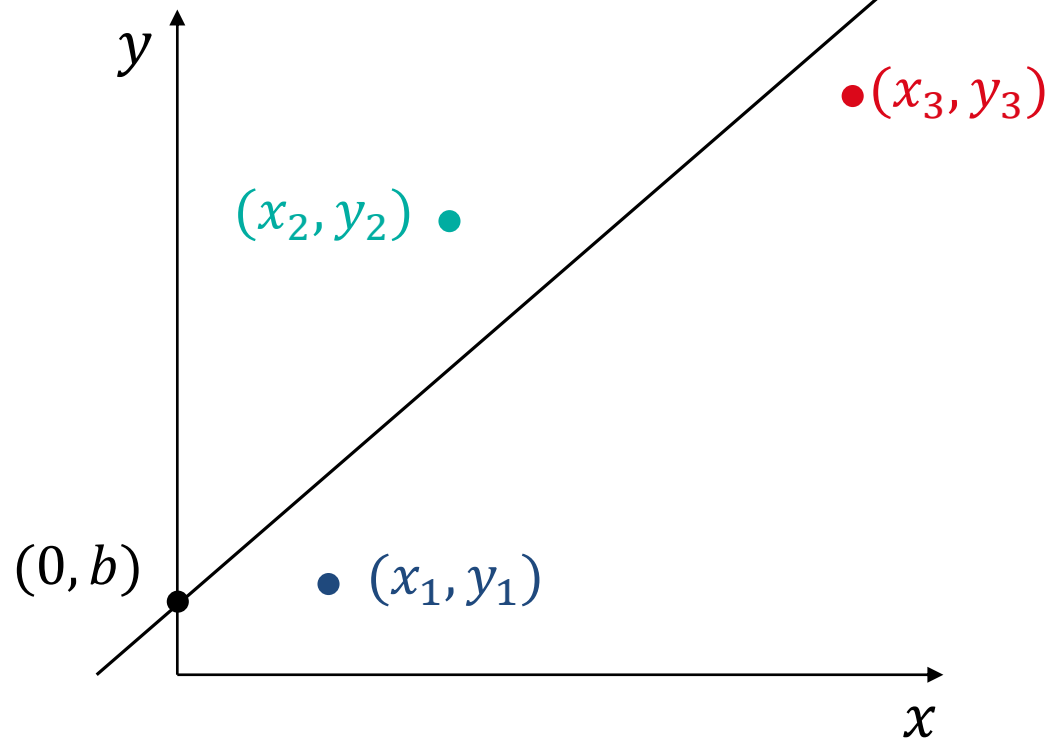
Proof (cont.)

$$SE_{Line} = n\bar{y}^2 - 2mn\bar{x}\bar{y} - 2bn\bar{y} + m^2n\bar{x}^2 + 2mbn\bar{x} + nb^2$$

$$m\bar{x}^2 + b\bar{x} = \bar{x}\bar{y}$$

$$m\bar{x} + b = \bar{y}$$

$$y = mx + b \text{ (Regression line)}$$



① (\bar{x}, \bar{y}) lies on the regression line.

Proof (cont.)

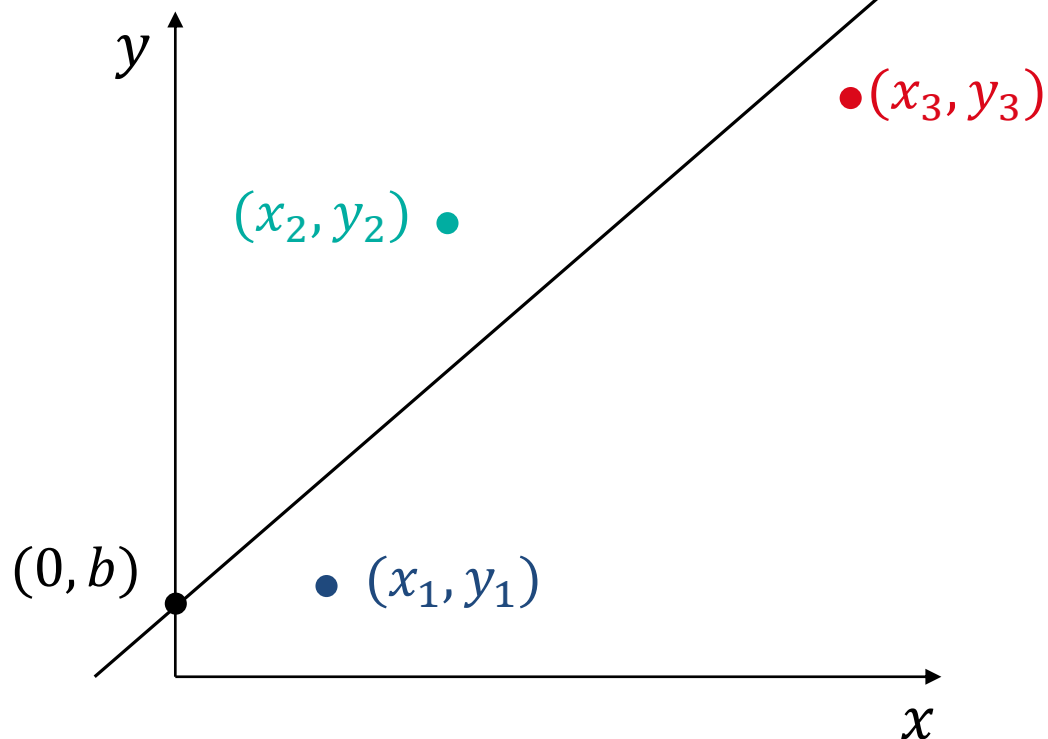
$$SE_{Line} = n\overline{y^2} - 2mn\overline{xy} - 2bn\overline{y} + m^2n\overline{x^2} + 2mbn\overline{x} + nb^2$$

$$m\overline{x^2} + b\overline{x} = \overline{xy}$$

$$m\overline{x} + b = \overline{y}$$

$$m\frac{\overline{x^2}}{\overline{x}} + b = \frac{\overline{xy}}{\overline{x}}$$

$y = mx + b$ (Regression line)

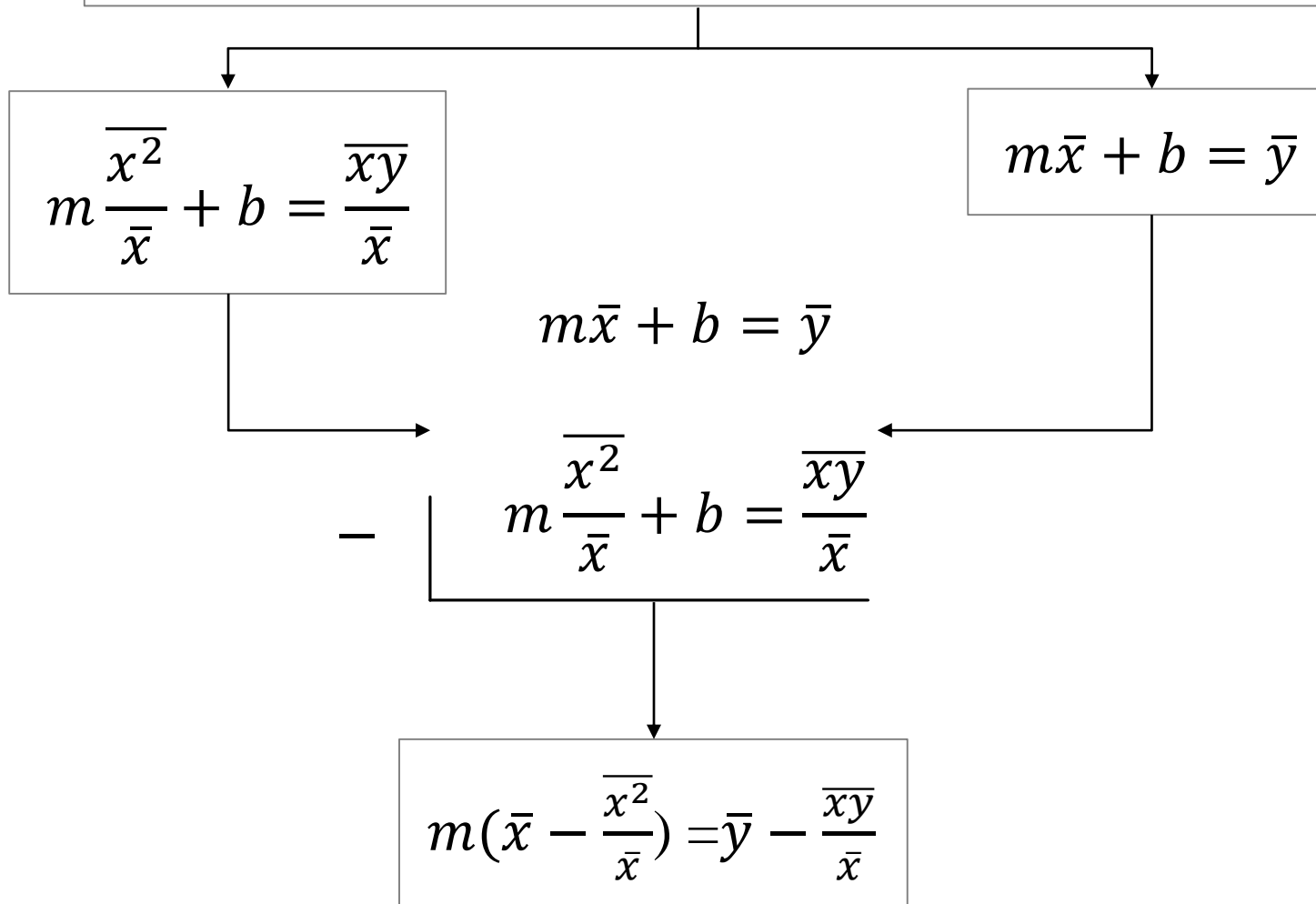


① $(\overline{x}, \overline{y})$ lies on the regression line.

② $(\frac{\overline{x^2}}{\overline{x}}, \frac{\overline{xy}}{\overline{x}})$ lies on the regression line.

Proof (cont.)

$$SE_{Line} = n\overline{y^2} - 2mn\overline{xy} - 2bn\overline{y} + m^2n\overline{x^2} + 2mbn\overline{x} + nb^2$$



Proof (cont.)

$$SE_{Line} = n\overline{y^2} - 2mn\overline{xy} - 2bn\overline{y} + m^2n\overline{x^2} + 2mbn\overline{x} + nb^2$$

$$m \frac{\overline{x^2}}{\overline{x}} + b = \frac{\overline{xy}}{\overline{x}}$$

$$m\overline{x} + b = \overline{y}$$

$$m\overline{x} + b = \overline{y}$$

$$- \quad m \frac{\overline{x^2}}{\overline{x}} + b = \frac{\overline{xy}}{\overline{x}}$$

$$m\left(\overline{x} - \frac{\overline{x^2}}{\overline{x}}\right) = \overline{y} - \frac{\overline{xy}}{\overline{x}}$$

$$m = \frac{\overline{y} - \frac{\overline{xy}}{\overline{x}}}{\overline{x} - \frac{\overline{x^2}}{\overline{x}}}$$

$$m = \frac{\overline{y} - \frac{\overline{xy}}{\overline{x}}}{\overline{x} - \frac{\overline{x^2}}{\overline{x}}} \cdot \frac{\overline{x}}{\overline{x}} = \frac{\overline{x}\overline{y} - \overline{xy}}{(\overline{x})^2 - \overline{x^2}}$$

Proof (cont.)

$$SE_{Line} = n\overline{y^2} - 2mn\overline{xy} - 2bn\overline{y} + m^2n\overline{x^2} + 2mbn\overline{x} + nb^2$$

$$m \frac{\overline{x^2}}{\overline{x}} + b = \frac{\overline{xy}}{\overline{x}}$$

$$m = \frac{\overline{x}\overline{y} - \overline{xy}}{(\overline{x})^2 - \overline{x^2}}$$

$$m\overline{x} + b = \overline{y}$$

$$b = \overline{y} - m\overline{x}$$

Proof (cont.)

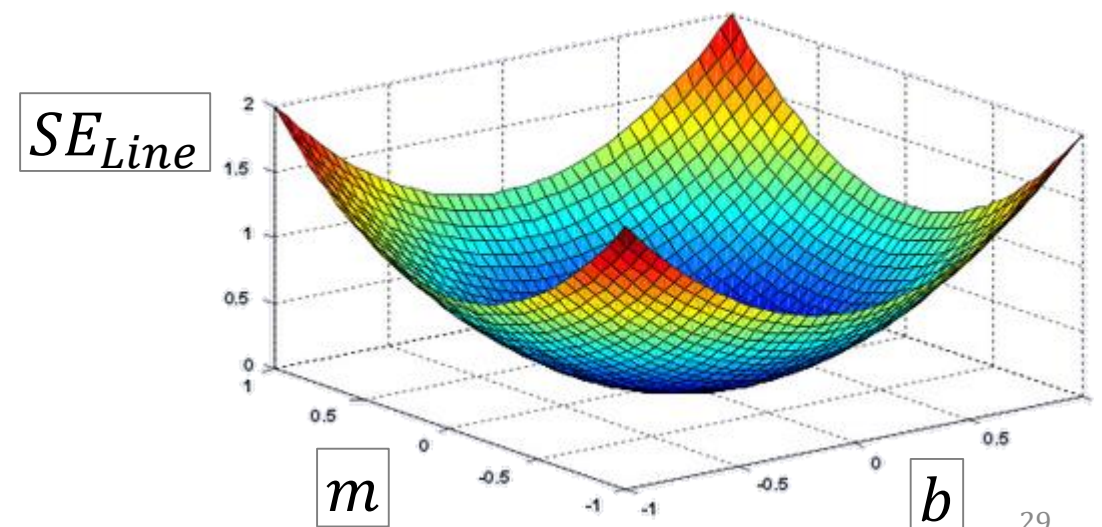
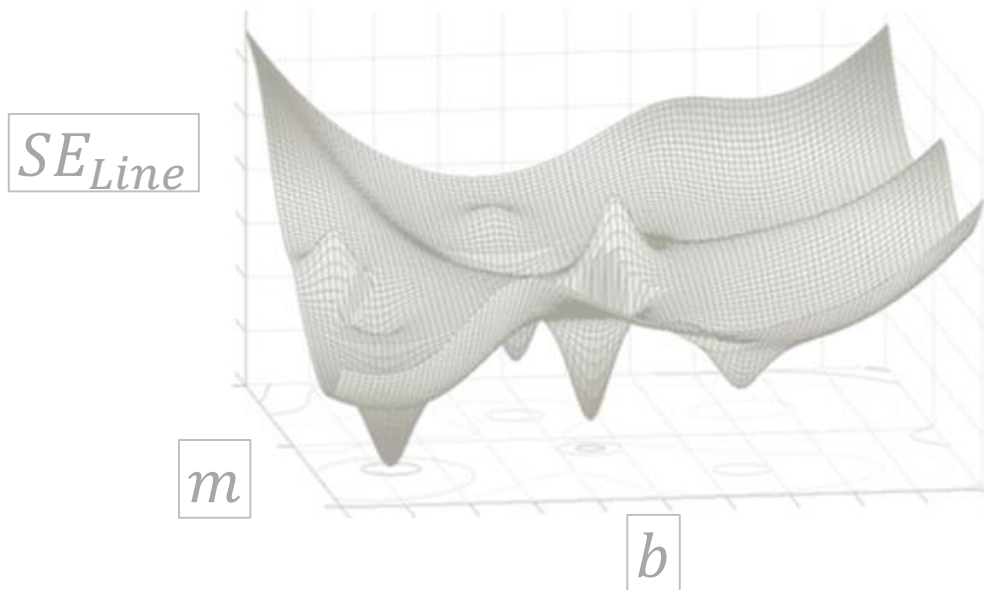
$$SE_{Line} = n\bar{y}^2 - 2mn\bar{x}\bar{y} - 2bn\bar{y} + m^2n\bar{x}^2 + 2mbn\bar{x} + nb^2$$

$$m \frac{\bar{x}^2}{\bar{x}} + b = \frac{\bar{x}\bar{y}}{\bar{x}}$$

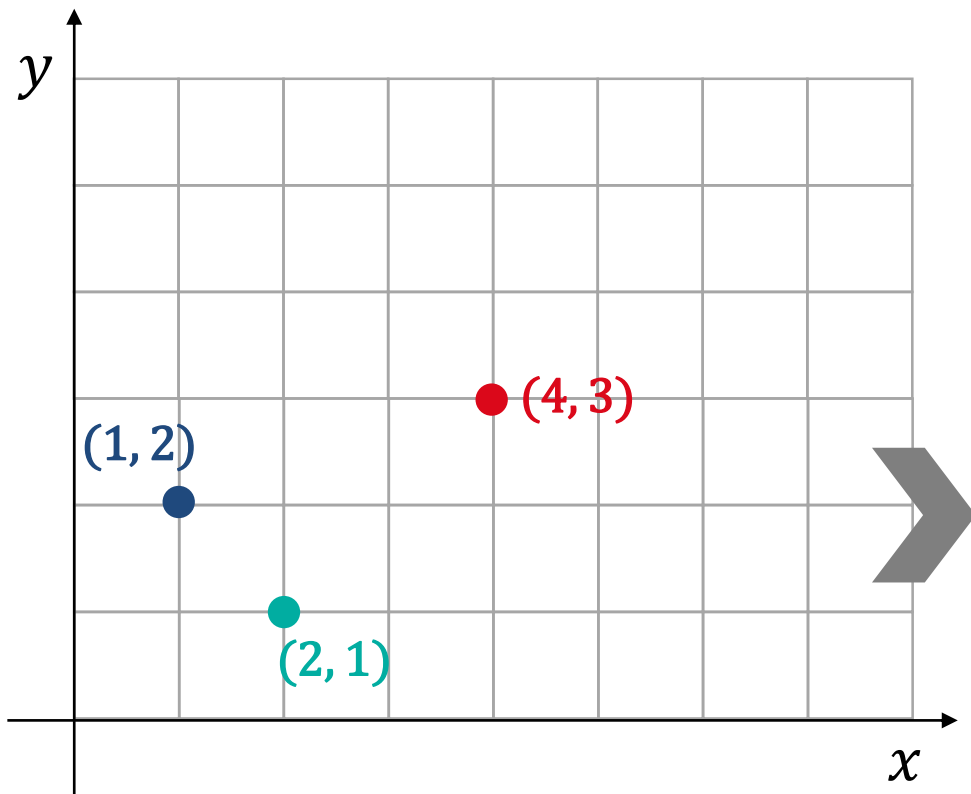
$$m = \frac{\bar{x}\bar{y} - \bar{x}\bar{y}}{(\bar{x})^2 - \bar{x}^2}$$

$$m\bar{x} + b = \bar{y}$$

$$b = \bar{y} - m\bar{x}$$



Regression line example



$$m = \frac{\bar{x}\bar{y} - \overline{xy}}{(\bar{x})^2 - \overline{x^2}}$$

$$b = \bar{y} - m\bar{x}$$

$$\bar{x} = \frac{1 + 2 + 4}{3} = \frac{7}{3} \quad \bar{y} = \frac{2 + 1 + 3}{3} = 2$$

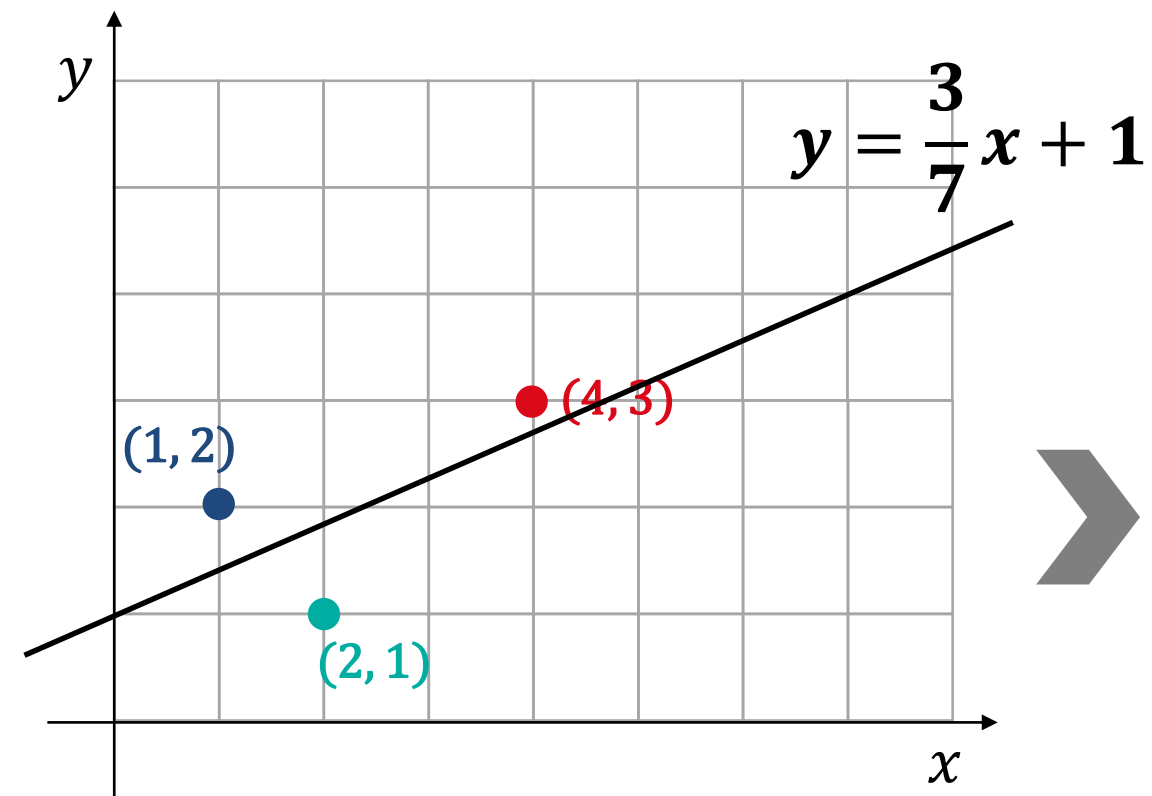
$$\overline{xy} = \frac{1 \cdot 2 + 2 \cdot 1 + 4 \cdot 3}{3} = \frac{16}{3}$$

$$\overline{x^2} = \frac{1^2 + 2^2 + 4^2}{3} = 7$$

$$m = \frac{\frac{7}{3} \cdot 2 - \frac{16}{3}}{(\frac{7}{3})^2 - 7} = \frac{\frac{14 - 16}{3}}{\frac{49 - 63}{9}} = \frac{-\frac{2}{3}}{-\frac{14}{9}} = \frac{18}{42} = \frac{3}{7}$$

$$b = 2 - \frac{3}{7} \cdot \frac{7}{3} = 1$$

Regression line example (cont.)



$$m = \frac{\bar{x}\bar{y} - \overline{xy}}{(\bar{x})^2 - \overline{x^2}}$$

$$b = \bar{y} - m\bar{x}$$

$$\bar{x} = \frac{1 + 2 + 4}{3} = \frac{7}{3} \quad \bar{y} = \frac{2 + 1 + 3}{3} = 2$$

$$\overline{xy} = \frac{1 \cdot 2 + 2 \cdot 1 + 4 \cdot 3}{3} = \frac{16}{3}$$

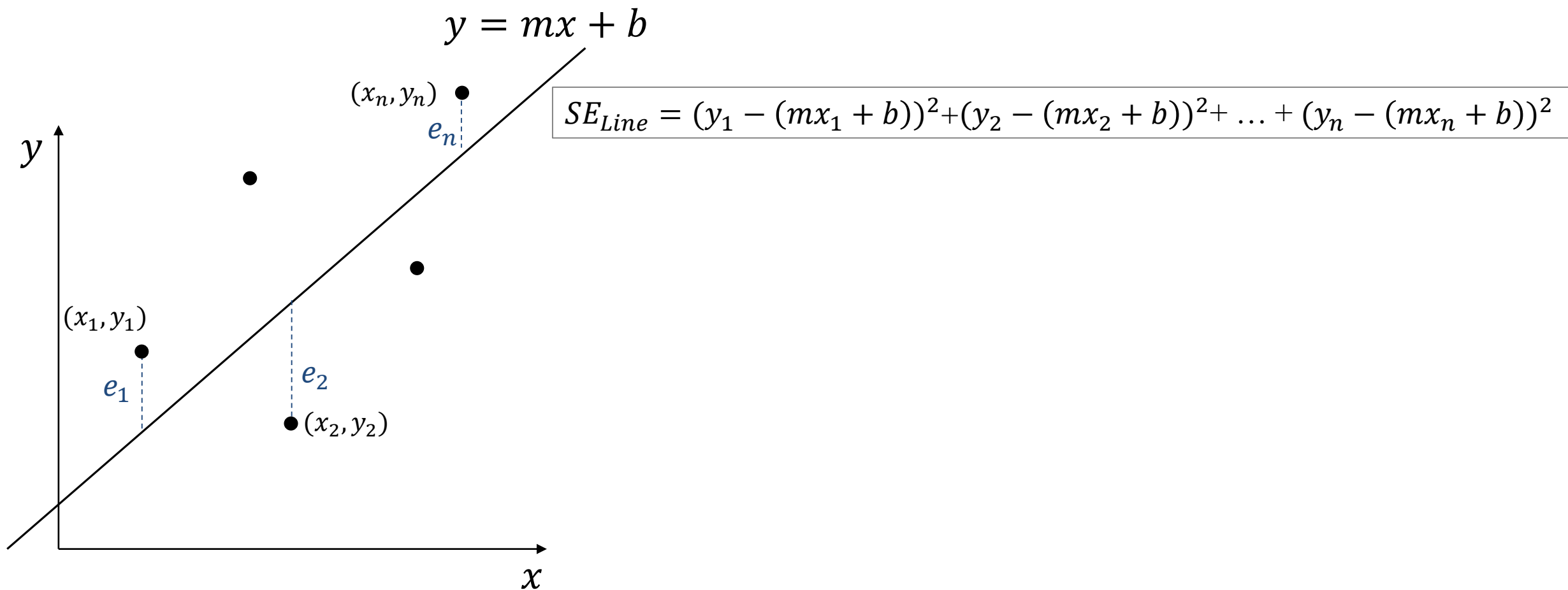
$$\overline{x^2} = \frac{1^2 + 2^2 + 4^2}{3} = 7$$

$$m = \frac{3}{7}$$

$$b = 1$$

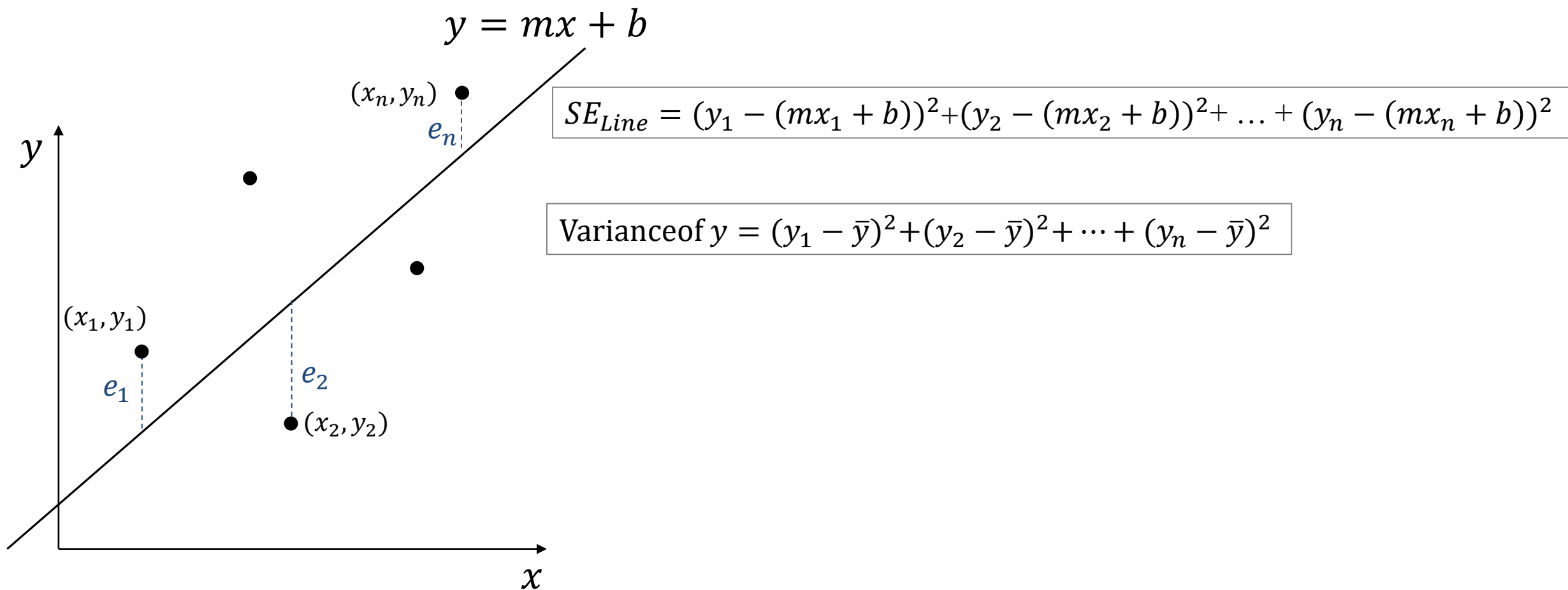
Coefficient of determination

- **Definition:** the proportion of the variance in the dependent variable that is predictable from independent variable(s).
 - What % of variation in y is described by the regression line (or by x).



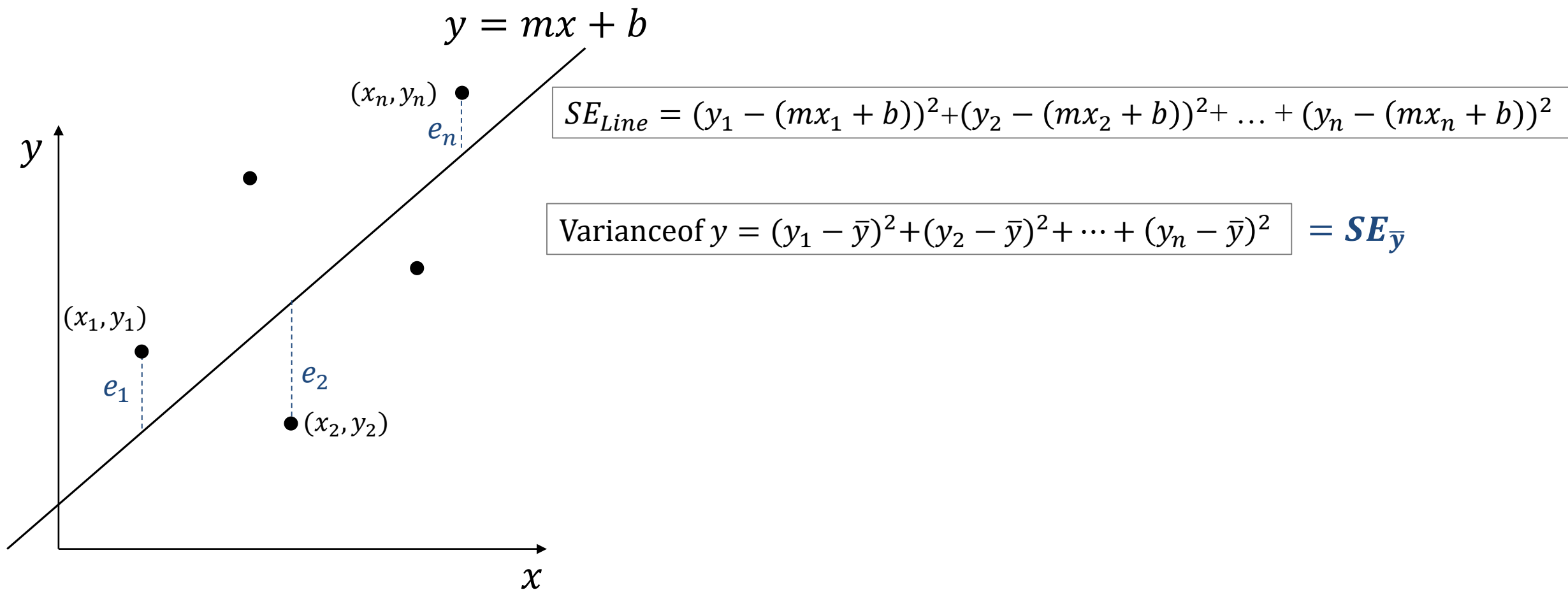
Coefficient of determination

- **Definition:** the proportion of the variance in the dependent variable that is predictable from independent variable(s).
 - What % of variation in y is described by the regression line (or by x).



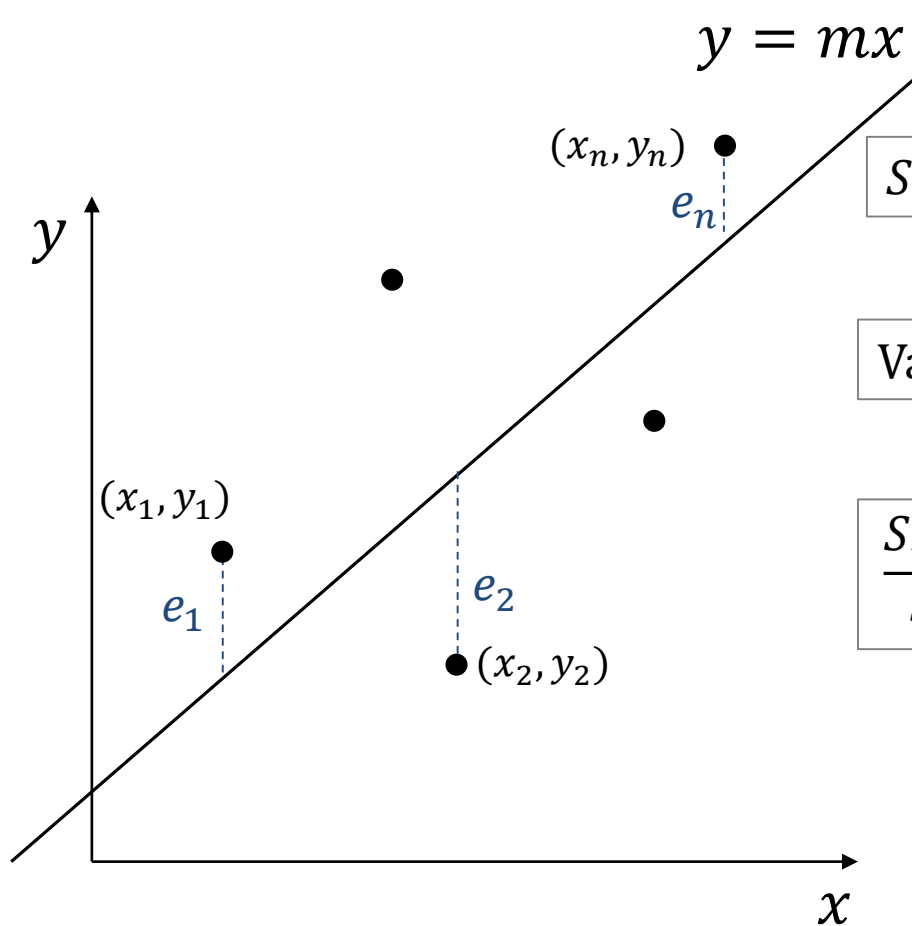
Coefficient of determination

- **Definition:** the proportion of the variance in the dependent variable that is predictable from independent variable(s).
 - What % of variation in y is described by the regression line (or by x).



Coefficient of determination

- **Definition:** the proportion of the variance in the dependent variable that is predictable from independent variable(s).
 - What % of variation in y is described by the regression line.



$$SE_{Line} = (y_1 - (mx_1 + b))^2 + (y_2 - (mx_2 + b))^2 + \dots + (y_n - (mx_n + b))^2$$

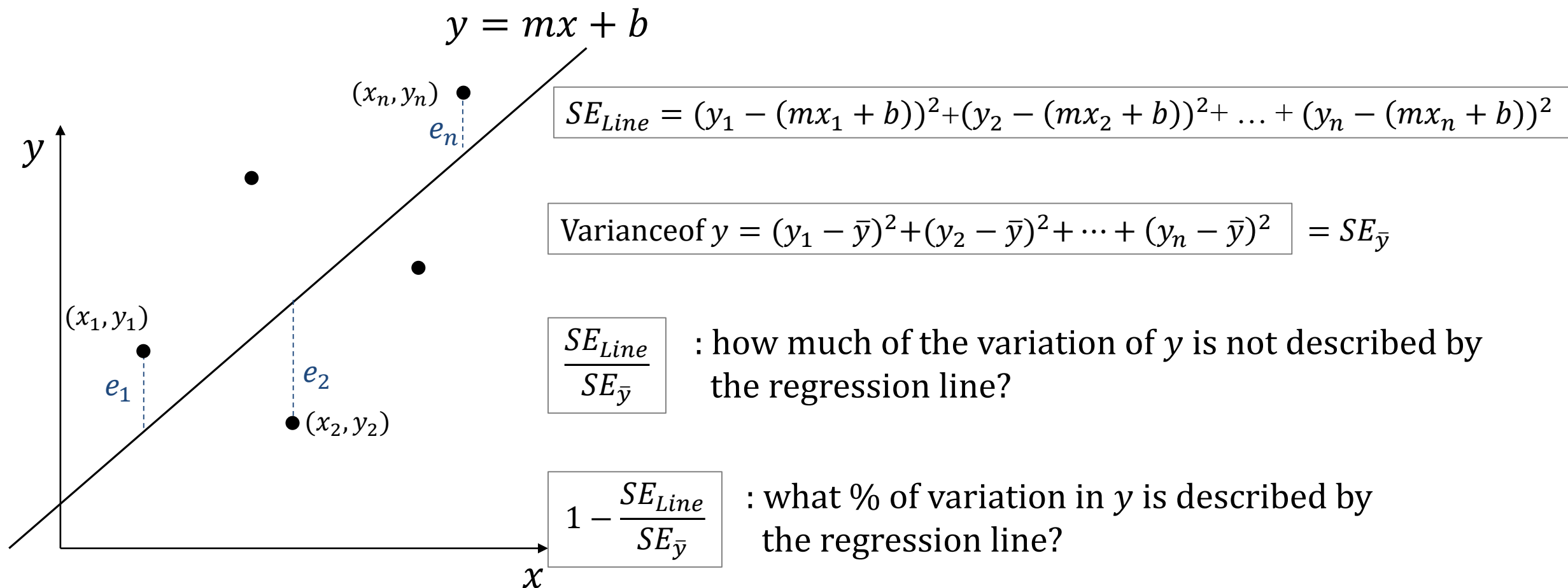
$$\text{Variance of } y = (y_1 - \bar{y})^2 + (y_2 - \bar{y})^2 + \dots + (y_n - \bar{y})^2 = SE_{\bar{y}}$$

$$\frac{SE_{Line}}{SE_{\bar{y}}}$$

: how much of the variation of y is not described by the regression line?

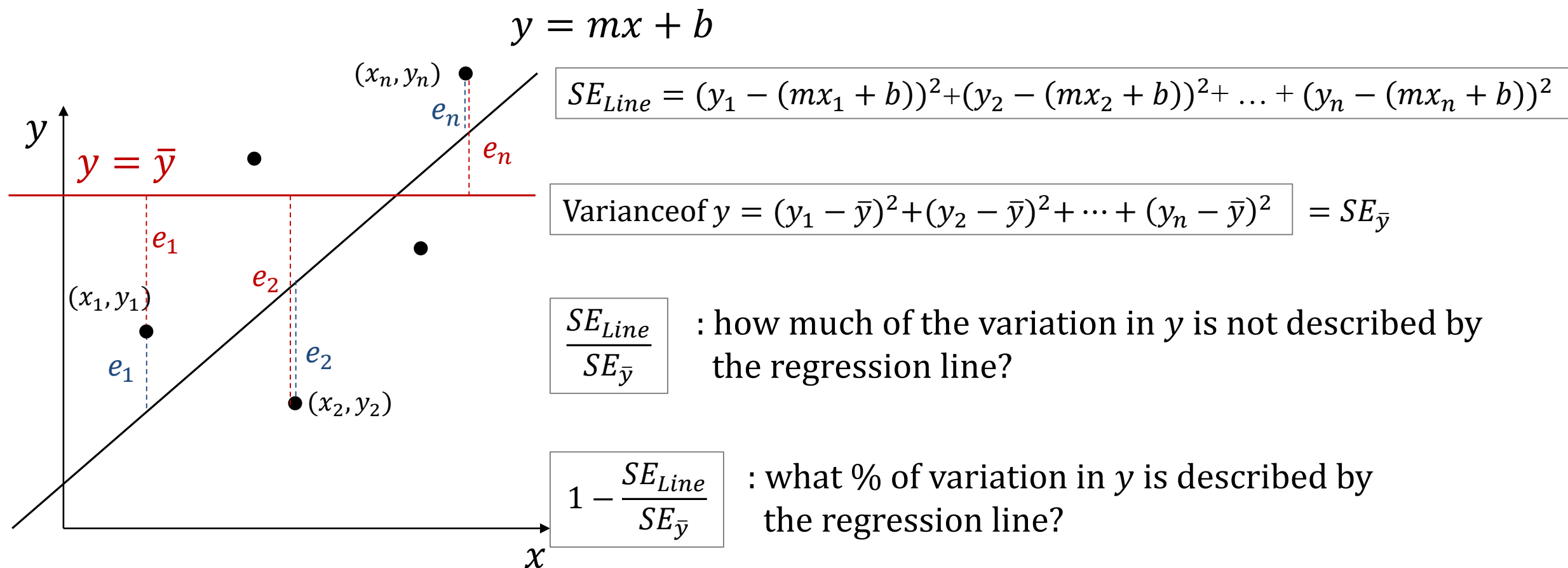
Coefficient of determination

- **Definition:** the proportion of the variance in the dependent variable that is predictable from independent variable(s).
 - What % of variation in y is described by the regression line?



Coefficient of determination

- **Definition:** the proportion of the variance in the dependent variable that is predictable from independent variable(s).
 - What % of variation in y is described by the regression line?



Q & A