

learning from data: Predictive Modelling

Linear Regression

Input features, targets.

Quantitative variables, Qualitative (Categorical) Variables

Predictive Modelling: Input features \rightarrow target

1. We want to predict if a patient is likely to require C-section delivery.

What information will be useful to predict this?

Input features: Patient Age, Patient history for example if this is the first time pregnancy? If the patient has undergone

Target: requires caesarian delivery vs. does not require caesarian delivery

The target here is a categorical variable. Yes, No?

The input feature "Patient Age" is a quantitative or N.

2. We want to predict the value of a stock in the stock market.

Input features: Past history of stock, Current financial situation of company, future growth prospective.

Stock. This is an example of regression problem

Prediction :

$$f: X \rightarrow Y \in \mathbb{R}$$

\mathbb{R} : Real number.

Categorial

Solve

Delivery : This is an example of classification problem.

$$f: X \rightarrow Y \quad Y \in \{0, 1, 2, \dots, 2\}$$

where the target/label is a category.

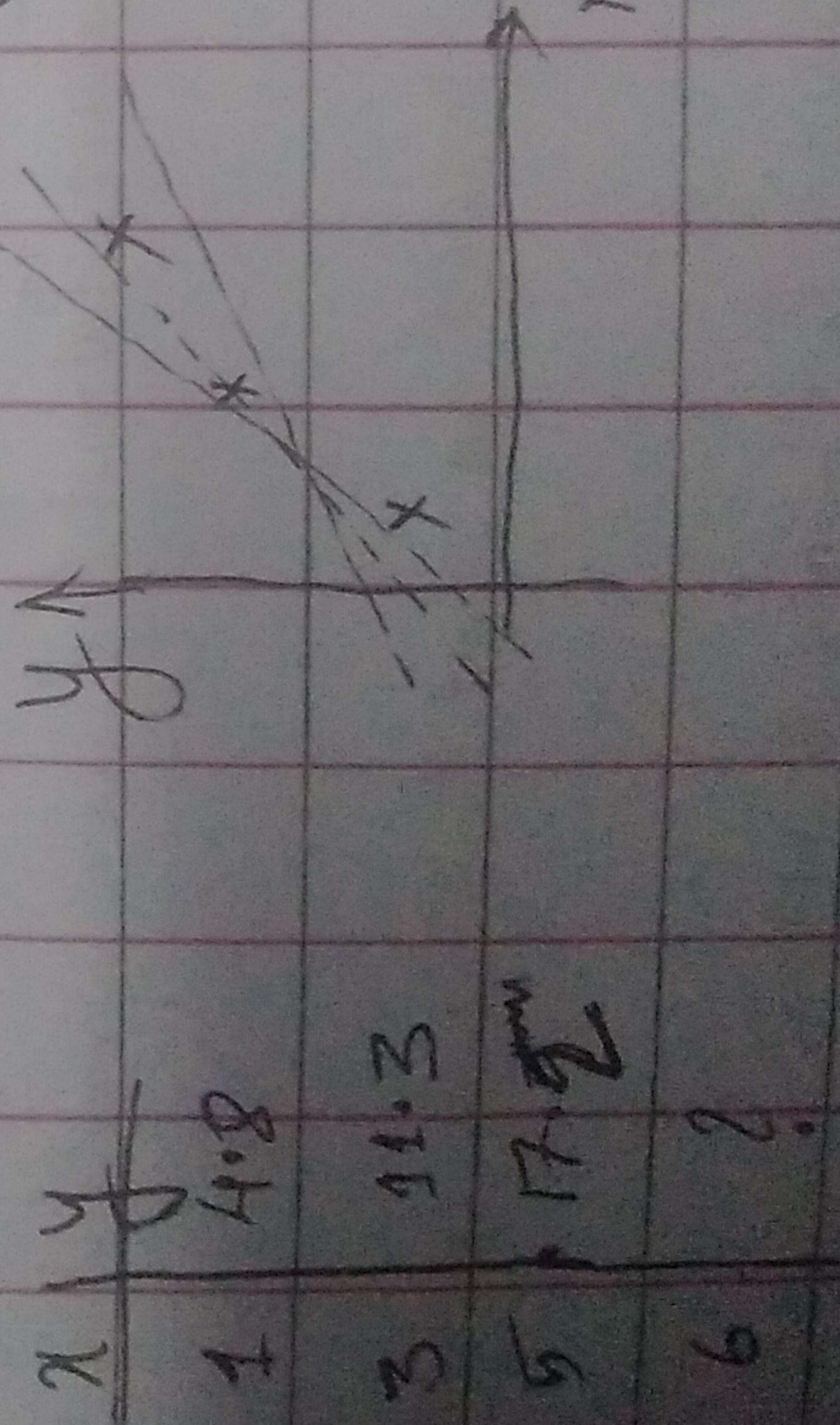
Linear Regression

$$f: X \rightarrow Y \quad Y \in \mathbb{R}$$

~~f~~) f is linear.

Let's take a simple example of learning to predict from data where input and target are scalar variables.

We will fit a line of the form $y = \theta_0 + \theta_1 x$.



All these dotted lines

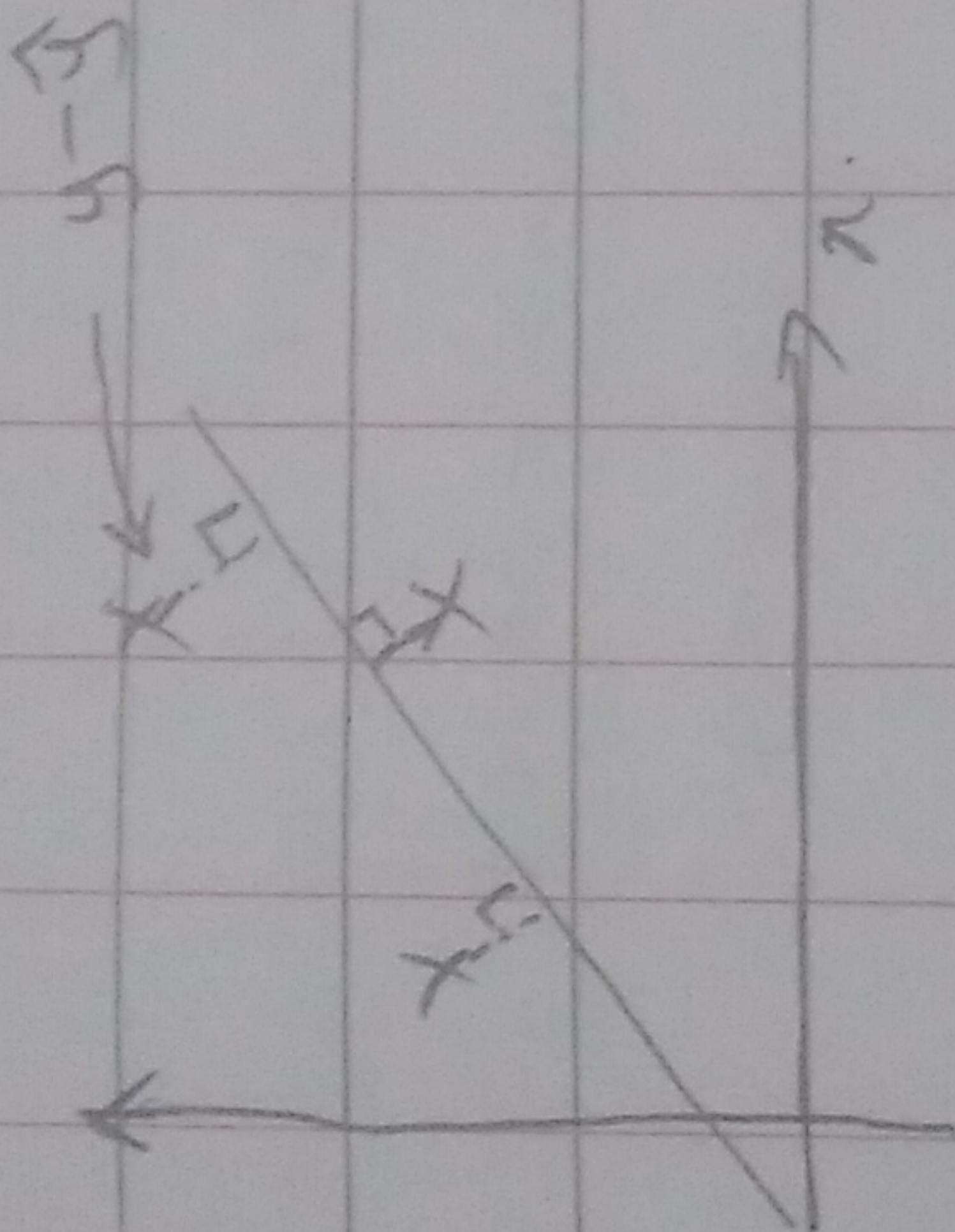
look like good enough

n fit for the data.

We want to pick a unique solution.

What would be a good criteria of choosing among all these possible lines?

Solution: One reasonable solution is pick a line s.t. Mean square error (MSE) is smallest.



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

(also called L2 norm)
of the error.

Other reasonable solutions

Mean absolute error (MAE)

To obtain the parameters (θ_0, θ_1)
of the line that satisfies
the least square error, use the following

formula

$$\theta_1 = \frac{\bar{xy} - \bar{x} \cdot \bar{y}}{\bar{x}^2 - (\bar{x})^2} \quad \text{--- (1)}$$
$$\theta_0 = \bar{y} - \theta_1 \cdot \bar{x}$$

To do: derivation of the equation (1).

Applying eqn(1), we have

x	y	xy	x^2
4	4.8	4.8	1
3	11.3	33.9	9
5	17.2	86	25

Average:

$$\bar{x} = 3, \bar{y} = 11.1, \bar{xy} = 41.57, \bar{x^2} = 11.67.$$

$$\theta_1 = \frac{\bar{xy} - \bar{x}\bar{y}}{\bar{x^2} - (\bar{x})^2} = \frac{41.57 - 3 \times 11.1}{11.67 - 9}$$

$$= 3.1.$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x} = 11.1 - 3.1 \times 3.$$

$$= 1.8.$$

Thus, the best fit line for the NSS criteria is

$$y = f(x) = 1.8 + 3.1x.$$

Derivation of Equation 1.

Where did this magical equation ① come from?

lost function

Given a dataset $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$, find a functional relationship of the form

$$f(x; \theta) = f_\theta(x) = \theta_0 + \theta_1 x$$

such that $f_\theta(x)$ minimizes mean squared error
i.e cost function $L(f_\theta(x), D)$

Setup:

Use the idea that $\frac{\partial L(f_\theta(x), D)}{\partial \theta} = 0$ is the point

where $\theta = \theta^*$

$$\begin{aligned}\frac{\partial L(f_\theta(x), D)}{\partial \theta} &= \frac{1}{N} \sum_{i=1}^N (y_i - (\theta_0 + \theta_1 x_i))^2 \\ &= \frac{1}{N} \sum_{i=1}^N y_i - (\theta_0 + \theta_1 x_i)^2\end{aligned}$$

Differentiating the cost function,

$$\frac{\partial L}{\partial \theta_0} = \frac{1}{N} \sum_{i=1}^N (-2y_i + 2\theta_0 + 2\theta_1 x_i) \quad \text{--- ②}$$

$$\frac{\partial L}{\partial \theta_1} = \frac{1}{N} \sum_{i=1}^N (-2x_i y_i + 2\theta_0 + 2\theta_1 x_i^2) \quad \text{--- ③}$$

$$\text{Equation } \frac{\partial L}{\partial v_0} = 0$$

$$\text{via term, } \frac{1}{v_1} \sum_{i=1}^n (y_i + v_0 + v_1 x_i) = 0$$

$$v_1 - \bar{y}_1 + v_0 + v_1 \bar{x} = 0$$

$$v_1 - v_0 + v_1 \bar{x} - \bar{y}_1 = 0 \quad \textcircled{4}$$

$$\text{Equation } \frac{\partial L}{\partial v_1} = 0$$

$$\text{via term, } \frac{1}{v_1} \sum_{i=1}^n (y_i + v_0 + v_1 x_i)^2 = 0.$$

$$v_1 - \bar{y}_1 + v_0 \bar{x} + v_1 \bar{x}^2 = 0$$

$$v_1 - v_0 \bar{x} + v_1 \bar{x}^2 - \bar{y}_1 = 0 \quad \textcircled{5}$$

Equation eqn \textcircled{4} and \textcircled{5},

$$v_0 + v_1 \bar{x} - \bar{y}_1 = 0$$

$$v_0 \bar{x} + v_1 \bar{x}^2 - \bar{y}_1 = 0,$$

$$v_1 \left[\frac{1}{\bar{x}} \right] [0, 1] \rightarrow \left[\frac{\bar{y}_1}{\bar{x}} \right]$$

$$f(x) = \frac{1}{2}x^2 - (\bar{x})^2$$

$$= \frac{1}{2}x^2 - \frac{1}{2}x^2 + \bar{x}\bar{x}$$

$$= \bar{x}\bar{x}$$

$$= \frac{1}{x^2 - 5x} = \frac{1}{x(x-5)} = \frac{1}{x} + \frac{1}{x-5}$$

Home 91 *
 ~~$x^2 - x \cdot y$~~
 $x^2 - (x)^2$

— $\bar{x}_2 \cdot \bar{y} + x_1 \cdot y$ —
thus can be defined
— $\bar{x}_1 \cdot \bar{y} + x_2 \cdot y$ —
— $\bar{x}_1 \cdot \bar{y} + x_1 \cdot y$ —

$$x^2 - y^2 = (x+y)(x-y)$$

$$\frac{1}{x} - \frac{1}{x^2}$$

卷之三