$\sim y = a_0 + a_1 x$

$(x_N, y_N)$

$(x_1, y_1)$

In linear regression, we try to get the best fit line, which is a straight line.

Same approach in → all statistical learning

Standard approach to solve linear regression
1) Deriving closed form solution
2) Applying gradient descent

In linear regression, the model consists of linear function
$$y = \sum_j w_j x_j + b \quad —①$$

Loss function is defined as
$$\ell(y, t) = \frac{1}{2}(y - t)^2 \quad —②$$

Considering ② in ①, we get    cost function

$$C(w_1, w_2, \ldots, w_D, b) = \frac{1}{N} \sum_{i=1}^{N} \ell(y^{(i)}, t^{(i)})$$

cost function $= \dfrac{1}{2N} \sum_{i=1}^{N} (y^{(i)} - t^{(i)})^2$

$$= 1/2N \sum_{i=1}^{N} \left( \sum_j w_j x_j^{(i)} + b - t^{(i)} \right)^2 \quad —(3)$$

└——┘ — Here the choice of $w_i$ & $b$ is optimized to reduce $C$.

for a optimization problem, a good place to start is to compute the partial derivative of cost function

(3) ⇒ Applying chain rule

$$\delta C / \delta w_j = 1/N \sum_{i=1}^{N} x_j^{(i)} \left( \sum_{j'} w_{j'} x_{j'}^{(i)} + b - t^{(i)} \right)$$

$$\delta C / \delta b = 1/N \sum_{i=1}^{N} \left( \sum_{j'} w_{j'} x_{j'}^{(i)} + b - t^{(i)} \right)$$

The following can be rewritten as

$$\delta C / \delta w_j = (1/N) \sum_{i=1}^{N} x_j^{(i)} (y^{(i)} - t^{(i)}) \quad —⑤$$

$$\delta C / \delta b = 1/N \sum_{i=1}^{N} y^{(i)} - t^{(i)}$$

If a function is differentiable, then, $\delta f / \delta x_i$ is 0 at minimum.

$\delta f / x_i \quad \rightarrow -ve, \quad$ increase $x_i$ slightly
$\delta f / x_i \quad \rightarrow +ve, \quad$ decrease $x_i$ slightly

Critical point $\rightarrow$ partial derivative is zero

So, equating ⑤, to 0, to get the critical point value,

$$\delta C / \delta w_j = 1/N \left( \sum_{i=1}^{N} x_j^{(i)} (y^{(i)} - t^{(i)}) \right) = 0 \quad —Eq 6$$

Solving for above will provide the relevant weight
The above is direct method.

## Gradient Descent Method

In gradient descent, the direction of steepest ascent of a $f^n$ is determined. Gradient descent considers partial derivative of the variables.

$$\frac{\delta E}{\delta w} = \begin{pmatrix} \frac{\delta E}{\delta w_1} \\ \vdots \\ \delta E / \delta w_D \end{pmatrix}$$

Using update parameter,

$$w \leftarrow w - \alpha \frac{\delta E}{\delta w}$$