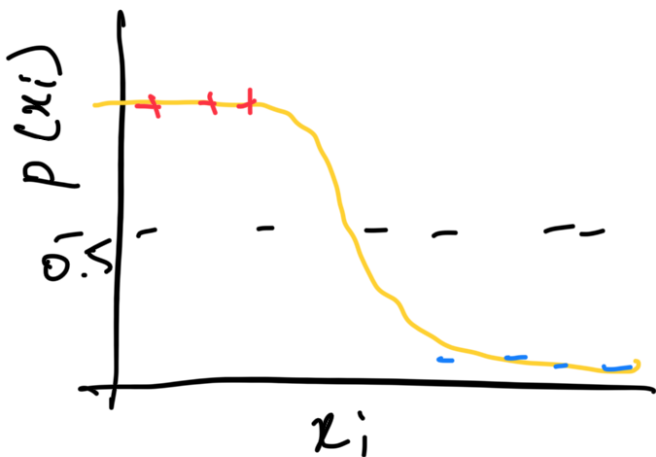# Logistic Regression.

→ classification ; binary cross-entropy loss.

$$\hat{y}_i = \sigma(wx_i + b)$$

$\quad\quad\quad \hookrightarrow$ sigmoid : $\sigma(z) = \dfrac{1}{1 + e^{-z}}$

→ Now loss, here it is not straightforward,

so, lets see how loss came



→ Lets, say we are having $y_i \{1, -1\}$.

→ And $P(x_i) \to \hat{y}_{(i)}$ is the prob of the class.

$\quad\quad\quad\quad\quad\quad\quad\quad$ req model to

$\rightarrow$ We want our $y_i \rightarrow$ log, ~0

maximize the observed prob.. either be

tve or -ve.

$\Rightarrow$ The total prob. of data , since each

observation is i.i.l.d

$$\Pi \ (\text{data points})$$

$$\Rightarrow \quad \Pi \begin{cases} g(i) \quad \text{when } y_i = +1 \\ 1 - g(i) \quad \text{when } y_i = -1 \end{cases}$$

$g(i) \rightarrow$ prob. of $x_i$ being $+1$, that's why

we use $1 - g(i)$ when $x_i \neq +1$, to maximize

overall prob.

$$\prod_{i=1}^{N} \ (g(i))^{\overset{\rightarrow y_i}{1\{y_i = +1\}}} \cdot (1-g(i))^{\overset{\rightarrow 1-y_i}{1\{y_i = -1\}}}$$

$\rightarrow$ multiplying all maximized prob.

$$\Rightarrow \prod^{N} \hat{y}_i^{\ y_i} \cdot (1-\hat{y}_i)^{1-y_i}$$

$$i = 1$$

→ the loss is -ve of above term & since multiply long prob ~0 and math. efficiency converting multiplication to addition we take log.

$$loss = -\log [\ ..\ ]$$

$$= y_i \cdot \log(\hat{y}) + (1-y_i) \log(1-\hat{y})$$

final loss $\hat{\ }$

Now, need to find.

$$\frac{\partial L}{\partial w} \quad \& \quad \frac{\partial L}{\partial b}$$

$$\frac{\partial L}{\partial w} = y_i \frac{\partial \log(\hat{y})}{\partial w} + (1-y_i) \frac{\partial}{\partial w} \log(1-\hat{y})$$

lets start with inner,

$$\frac{\partial (\log(\hat{y_i}))}{\partial w} = \frac{1}{\hat{y}} \cdot \frac{\partial \hat{y_i}}{\partial w}$$

Since, $\hat{y}_i = \sigma(z_i) = \dfrac{1}{1 + e^{-z_i}}$ ; $z_i = w x_i + b$

$$\frac{\partial \hat{y}_i}{\partial w} = \frac{\partial \sigma(z_i)}{\partial w} = \sigma'(z_i) \cdot \frac{\partial z_i}{\partial w}$$

$$\sigma'(z_i) = \sigma(z_i)(1 - \sigma(z_i)) = \hat{y}_i \cdot (1 - \hat{y}_i)$$

$$z_i = w x_i + b \; ; \quad \frac{\partial z_i}{\partial w} = x_i$$

$$\frac{\partial z_i}{\partial b} = 1$$

$$\therefore \quad \frac{\partial \hat{y}_i}{\partial w} = \hat{y}_i \cdot (1 - \hat{y}_i) \cdot x_i$$

$$\therefore \quad \frac{\partial (\log(\hat{y}_i))}{\partial w} = \frac{1}{\hat{y}} \cdot (\hat{y}_i) \cdot (1 - \hat{y}_i) \cdot x_i$$

$$= (1 - \hat{y}_i) x_i$$

Similarly,

$$\frac{\partial \left( \log(1 - \hat{y}_i) \right)}{\partial w} = \frac{1}{1 - \hat{y}_i} \cdot \frac{\partial (1 - \hat{y}_i)}{\partial w}$$

$$= \frac{1}{1 - \hat{y}} \cdot (-1) \cdot \hat{y}_i \cdot (1 - \hat{y}_i) x_i$$

$$= - \hat{y}_i x_i$$

Combining everything.

$$\frac{\partial L}{\partial w} = y_i (1 - \hat{y}_i) x_i + (-\hat{y}_i x_i)(1 - \hat{y}_i)$$

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

Vectorized,

$$\frac{1}{n} \cdot X^T (y - \hat{y})$$

$$-\frac{\hat{y}_i}{\partial w}$$

$$\frac{\partial \sigma}{\partial b}, \text{ not calculating }, \text{...},$$

$$\frac{1}{n} \sum (\hat{y} - y)$$