# Logistic Regression & Multi-class classification

binary linear classification $z = w^T x + b$. $y = \begin{cases} 1, & z \geq r \\ 0, & z < r \end{cases}$

$w^T x + b \geq r \Rightarrow w^T x + b - r \geq 0$.

Recall that we can incorporate $b$ into the weight matrix $w$
$\Rightarrow z = w^T x$; $x_0 = 1$ & $x \in \mathbb{R}^{D+1}$

**NOT**

| $x_0$ | $x_1$ | $t$ |
|-------|-------|-----|
| 1 | 0 | 1 |
| 1 | 1 | 0 |

$x_0$ is always 1 because it's the dummy feature we use to incorporate $b$ into $w$.

when $x_1 = 0$: $w_0 x_0 + w_1 x_1 \geq 0 \Rightarrow x_0 \geq 0$    there are many

$x_1 = 1$: $w_0 x_0 + w_1 x_1 < 0 \Rightarrow w_0 + w_1 < 0$.   possible solutions

**AND**

| $x_0$ | $x_1$ | $x_2$ | $t$ |
|-------|-------|-------|-----|
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 0 |
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 1 |

$z = w_0 x_0 + w_1 x_1 + w_2 x_2$

$w_0 < 0$

$w_0 + w_2 < 0$    again there are many

$w_0 + w_1 < 0$     possible solutions

$w_0 + w_1 + w_2 > 0$

For NOT      $w_1$          an alternative picture:

$w_1 = -w_0$

$w_0 > 0$

$w_0$

$w_0 + w_1 < 0$

$x > 0$

$x$

$y + x < 0$

$y = x$;

$y + x = 0$

1

$$\mathcal{L}_{0,1}(y=t) = \begin{cases} 0 & y=t \\ 1 & y \neq t \end{cases} \Rightarrow \mathcal{L}_{0,1}(y,t) = \mathbb{I}(y \neq t)$$

$$\tilde{J} = \frac{1}{N} \sum_{i}^{N} \mathcal{L}_{0,1} \boxed{\mathbb{I}(y_i \neq t_i)} \qquad \frac{\partial \mathcal{L}_{0,1}}{\partial w_j} = \frac{\partial \mathcal{L}_{0,1}}{\partial z} \cdot \frac{\partial z}{\partial w_j}$$

$\dfrac{\partial \mathcal{L}_{0,1}}{\partial z}$ = (almost) zero anywhere it's defined.   $\Rightarrow$ unable to do
(refer to the figure on slide 14)      gradient descent

$\mathcal{L}_{SE} = \frac{1}{2}(z-t)^2$ ; let's set the final prediction threshold to be 0.5.

$\qquad \Rightarrow$ if $z \geq 0.5$, predict positive

$\qquad$ if $z < 0.5$, predict negative

Example:  ① a sample is "very positive" c the model predicts positive with high confidence. $\Rightarrow$ z is large ( z = 1000, for example)

② a sample is "not very positive" c the model predicts positive with low confidence. $\Rightarrow$ z is small but still above the threshold. ( z = 0.6, for example)

in ① : $\mathcal{L}_{SE} = \frac{1}{2}(\overset{z}{\cancel{y}} - t)^2 = \frac{1}{2} 999.5^2$  $\Rightarrow$ the loss function hates
in ② : $\mathcal{L}_{SE} = \frac{1}{2}(z-t)^2 = \frac{1}{2} 0.1^2$.    when you make c often time
                                     correct predictions with
                                     high confidence.

Logistic activation function. $\sigma(z) = \dfrac{1}{1+e^{-z}}$

$\Rightarrow z = \omega^T x$, $y = \sigma(z)$. $\mathcal{L}_{SE} = \dfrac{1}{2}(y-t)^2$

the logistic activation function converts an arbitrary big/small $z$ into the range $[0, 1]$. the bigger the $z$ is, $\sigma(z) = y$ approaches 1, and the smaller the $z$ is, vice versa.

$$\frac{\partial \mathcal{L}}{\partial w_j} = \frac{\partial \mathcal{L}}{\partial z} \cdot \frac{\partial z}{\partial w_j} \Rightarrow \text{differentiable now!}$$

Cross entropy loss $\quad \mathcal{L}_{CE} = -t\log y - (1-t)\log(1-y)$

$\mathcal{L}_{LCE} = \mathcal{L}_{CE}(\sigma(z), t) = -t\log\sigma(z) - (1-t)\log(1-\sigma(z))$

$\quad = -t\log\dfrac{1}{1+e^{-z}} - (1-t)\log\left(1 - \dfrac{1}{1+e^{-z}}\right)$

$\quad = -t[\log 1 - \log(1+e^{-z})] - (1-t)\left[\log\dfrac{1+e^{-z}}{1+e^{-z}}\right]$

$\quad = -t(0 - \log(1+e^{-z})) - (1-t)\log\dfrac{e^{-z}}{1+e^{-z}}$

$\quad = t\log(1+e^{-z}) - (1-t)[\log e^{-z} - \log(1+e^{-z})]$

$\quad = t\log(1+e^{-z}) - (1-t)[-z + \log(1+e^{-z})]$

$\quad = t\log(1+e^{-z}) + z + \log(1+e^{-z}) - tz - t\log(1+e^{-z})$

$\quad = z - tz + \log(1+e^{-z})$

$$\mathcal{L}_{LCE}(6(z),t) = -t\log\left(\frac{1}{1+e^{-z}}\right) - (1-t)\log\left(1 - \frac{1}{1+e^{-z}}\right)$$

$$= -t\left[\log 1 - \log(1+e^{-z})\right] - (1-t)\log\frac{\cancel{1}e^{-z}\cancel{1}}{1+e^{-z}}$$

$$= -t\left[0 - \log(1+e^{-z})\right] - (1-t)\log\frac{(e^{-z})\cdot e^z}{(1+e^{-z})\cdot e^z}$$

$$= t\log(1+e^{-z}) - (1-t)\log\frac{1}{e^z+1}$$

$$= t\log(1+e^{-z}) - (1-t)\left[\log 1 - \log(e^z+1)\right]$$

$$= t\log(1+e^{-z}) + (1-t)\log(e^z+1)$$

Gradient Descent of Logistic Regression.

$$\mathcal{L}(CE) = -t\log y - (1-t)\log(1-y)$$
$$y = \frac{1}{1+e^{-z}}, \quad z = w^T x$$

$$\frac{\partial \mathcal{L}_{CE}}{\partial w_j} = \underset{①}{\frac{\partial \mathcal{L}_{CE}}{\partial y}} \cdot \underset{②}{\frac{\partial y}{\partial z}} \cdot \underset{③}{\frac{\partial z}{\partial w_j}}$$

① $$\frac{\partial \mathcal{L}_{CE}}{\partial y} = \frac{d}{dy}\left[-t\log y - (1-t)\log(1-y)\right]$$

$$= \frac{d}{dy}(-t\log y) - \frac{d}{dy}(1-t)\log(1-y)$$

$$= \frac{-t}{y} + \frac{(1-t)}{1-y}$$

② $\dfrac{\partial y}{\partial z} = \dfrac{\partial}{\partial z}\left(\dfrac{1}{1+e^{-z}}\right)$ $\qquad \left(\dfrac{u}{v}\right)' = \dfrac{u'v - v'u}{v^2}$, whereas $u=1$,
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad V = 1+e^{-z}$

$\qquad\qquad = \dfrac{0 - (1+e^{-z})'}{(1+e^{-z})^2} = \dfrac{e^{-z}}{(1+e^{-z})^2}$ $\qquad$ Same

$y - y^2 = \dfrac{1}{1+e^{-z}} - \dfrac{1}{(1+e^{-z})^2} = \dfrac{1+e^{-z}-1}{(1+e^{-z})^2} = \dfrac{e^{-z}}{(1+e^{-z})^2}$

As a result, $\dfrac{\partial y}{\partial z} = y - y^2 = y(1-y)$

③ $\dfrac{\partial z}{\partial w_j} = \dfrac{\partial}{\partial w_j}\, \vec{w}^T x = \dfrac{\partial}{\partial w_j}(w_0 x_0 + w_1 x_1 + \cdots + w_j x_j + \cdots + w_{D+1}\, x_{D+1})$

$\qquad\qquad\qquad\qquad = x_j.$

$\Rightarrow \dfrac{\partial \mathcal{L}_{CE}}{\partial w_j} = \dfrac{\partial \mathcal{L}_{CE}}{\partial y}\cdot\dfrac{\partial y}{\partial z}\cdot\dfrac{\partial z}{\partial w_j}$

$\qquad\qquad = \left(-\dfrac{t}{y} + \dfrac{(1-t)}{(1-y)}\right)\cdot y(1-y)\cdot x_j$

$w_j \leftarrow w_j - \alpha\cdot\dfrac{\partial \mathcal{J}}{\partial w_j} = w_j - \alpha\cdot\dfrac{\partial}{\partial w_j}\dfrac{1}{N}\sum_i^N \mathcal{L}_{CE}^i$

$\qquad = w_j - \alpha\cdot\dfrac{1}{N}\sum_i^N \left(-\dfrac{t_i}{y_i} + \dfrac{1-t_i}{1-y_i}\right)\cdot y_i(1-y_i)\, x_j^i$

$\qquad = w_j \dfrac{\alpha}{N}\sum_i^N \left(-\dfrac{t_i}{y_i} + \dfrac{1-t_i}{1-y_i}\right)\cdot y_i(1-y_i)\cdot x_j^i$

# Multi-class linear classification

for the $k^{th}$ class, do a linear classifier $z_k = \sum_j^D w_{kj} \cdot x_j + b_k$

whereas your $k \in [1, \ldots, K]$ whereas there are $K$ classes.

$$y_i = \begin{cases} 1, & \text{if } i = \underset{k}{\text{argmax}} \; z_k \\ 0, & \text{if otherwise} \end{cases}$$

We want $\sum_k^K y_k = 1 \Rightarrow y_k = \text{softmax}(z_1 \cdots z_k)_k = \dfrac{e^{z_k}}{\sum_{k'} e^{z_{k'}}}$

$\mathcal{L}_{CE} = -\sum_k^K t_k \log y_k = -t^T \log y$

6.