# GRADIENT OF A LINEAR CLASSIFIER WITH CROSS-ENTROPY LOSS

$$\begin{bmatrix} \underline{\quad} & w_1^T & \rightarrow \\ \underline{\quad} & w_2^T & \rightarrow \\ & \vdots & \\ \underline{\quad} & w_m^T & \rightarrow \\ & \vdots & \\ \underline{\quad} & w_{y_i}^T & \rightarrow \end{bmatrix} \begin{bmatrix} \\ \\ x_i \\ \\ \\ \end{bmatrix} = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_m \\ \vdots \\ f_{y_i} \\ \vdots \end{bmatrix} \leftarrow \text{ground-truth class}$$

$$W$$

$$f_m = w_m^T x_i \qquad\qquad w_m: \text{ m-th row of } W$$

$$\begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_m \\ \vdots \\ f_{y_i} \\ \vdots \end{bmatrix} \xrightarrow[\text{softmax}(\cdot)]{} \begin{bmatrix} p_1 \\ p_2 \\ \vdots \\ p_m \\ \vdots \\ p_{y_i} \\ \vdots \end{bmatrix} \leftarrow \text{probability of ground-truth class} \qquad p_m = \frac{e^{f_m}}{\sum e^{f_j}}$$

cross-entropy loss for training sample $(x_i, y_i)$:

$$= L_i = -\log\left(\frac{e^{f_{y_i}}}{\sum e^{f_j}}\right) = -\log\left(p_{y_i}\right)$$

Want $\quad \nabla_W L_i = \begin{bmatrix} & \text{---} & \nabla_{W_1} L_i & \longrightarrow & \\ & \text{---} & \nabla_{W_2} L_i & \text{---} & \\ & & \vdots & & \\ & \text{---} & \nabla_{W_m} L_i & \text{---} & \\ & \text{---} & \nabla_{W_{y_i}} L_i & \text{---} & \end{bmatrix}$

Note: $\nabla_W L_i$ is of same dimension as $W$

$$\frac{\partial L_i}{\partial W_m} = \frac{\partial}{\partial W_m}\left[ -\log(P_{y_i}) \right]$$

$$= \underbrace{\frac{\partial\left[ -\log(P_{y_i}) \right]}{\partial P_{y_i}}}_{①} \quad \underbrace{\frac{\partial P_{y_i}}{\partial f_m}}_{②} \quad \underbrace{\frac{\partial f_m}{\partial W_m}}_{③} \qquad \text{chain rule}$$

$①: \quad \frac{\partial}{\partial P_{y_i}}\left[ -\log(P_{y_i}) \right] = \frac{-1}{P_{y_i}}$

$③: \quad \frac{\partial f_m}{\partial W_m} = \frac{\partial}{\partial W_m}\left[ W_m^T x_i \right] = x_i$

②: $\dfrac{\partial p_{y_i}}{\partial f_m} = \dfrac{\partial}{\partial f_m}\left[\dfrac{e^{f_{y_i}}}{\sum e^{f_j}}\right]$

(i) $m \neq y_i$    i.e. $f_m$ and $f_{y_i}$ are different variables

Recall: $\dfrac{\partial}{\partial x}\left[\dfrac{g(x)}{h(x)}\right] = \dfrac{g'(x)h(x) - g(x)h'(x)}{[h(x)]^2}$    quotient rule

$\dfrac{\partial p_{y_i}}{\partial f_m} = \dfrac{\partial}{\partial f_m}\left[\dfrac{e^{f_{y_i}}}{\sum e^{f_j}}\right]$

$\phantom{\dfrac{\partial p_{y_i}}{\partial f_m}} = \dfrac{\left(e^{f_{y_i}}\right)'\sum e^{f_j} - e^{f_{y_i}}\left[\sum e^{f_j}\right]'}{\left[\sum e^{f_j}\right]^2}$

$\phantom{\dfrac{\partial p_{y_i}}{\partial f_m}} = \dfrac{0 \cdot \sum e^{f_j} - e^{f_{y_i}} e^{f_m}}{\left[\sum e^{f_j}\right]^2}$

$\phantom{\dfrac{\partial p_{y_i}}{\partial f_m}} = \dfrac{-e^{f_{y_i}} e^{f_m}}{\left[\sum e^{f_j}\right]^2}$

Note: $\dfrac{\partial\left[e^{f_{y_i}}\right]}{\partial f_m} = 0$
∴ $f_m, f_{y_i}$ are different variable.

$\dfrac{\partial p_{y_i}}{\partial f_m} = - p_{y_i}\, p_m$    $\left(\text{when } m \neq y_i\right)$

③

(ii) $m = y_i$ i.e. $f_m$, $f_{y_i}$ are the same variable

$$\frac{\partial P_{y_i}}{\partial f_m} = \frac{\partial}{\partial f_m} \left[ \frac{e^{f_{y_i}}}{\sum e^{f_j}} \right]$$

$$= \frac{(e^{f_{y_i}})' \sum e^{f_j} - e^{f_{y_i}} \left[ \sum e^{f_j} \right]'}{\left[ \sum e^{f_j} \right]^2}$$

$$= \frac{e^{f_m} \sum e^{f_j} - e^{f_{y_i}} e^{f_m}}{\left[ \sum e^{f_j} \right]^2}$$

$$= \frac{e^{f_m}}{\sum e^{f_j}} - \frac{e^{f_{y_i}} e^{f_m}}{\sum e^{f_j} \sum e^{f_j}}$$

$$= P_m - P_{y_i} P_m$$

$$\frac{\partial P_{y_i}}{\partial f_m} = P_{y_i} (1 - P_{y_i}) \qquad \left( \begin{array}{c} \text{when} \\ m = y_i \end{array} \right)$$

$$\frac{\partial L_i}{\partial w_m} = \begin{cases} \frac{-1}{p_{y_i}} \left( -p_{y_i} \, p_m \right) x_i & m \neq y_i \\ \\ \frac{-1}{p_{y_i}} \, p_{y_i}(1 - p_{y_i}) \, x_i & m = y_i \end{cases}$$

$$\left( \begin{array}{c} \text{see} \\ \text{p.2} \end{array} \right)$$

$$= \begin{cases} p_m x_i & m \neq y_i \\ \\ (p_{y_i} - 1) \, x_i & m = y_i \end{cases}$$

$$\nabla_w L_i = \begin{bmatrix} \underline{\quad} \; p_1 x_i \; \underline{\quad} \\ \underline{\quad} \; p_2 x_i \; \underline{\quad} \\ \vdots \\ \underline{\quad} \; p_m x_i \; \underline{\quad} \\ \underline{\quad} \; (p_{y_i} - 1) x_i \; \underline{\quad} \\ \vdots \end{bmatrix}$$

* Remember
this is the
gradient matrix
for only one
training sample
(hence the
index $i$ )

Gradient descent:

$$W' = W - \gamma \, \nabla_w L_i$$

$\because p_m x_i \geq 0$

$\therefore w_m$ will decrease, so as
to decrease $f_m$

$\because (p_{y_i} - 1) x_i \leq 0$

$\therefore w_{y_i}$ will increase, so as t...