

Reference: <http://math.stackexchange.com/questions/945871/derivative-of-softmax-loss-function>

Inputs: x_i , where $i = 1, 2 \dots 785$

Hidden Layer: output is h_i and input is a_i , where $i = 1, 2 \dots 201$

Output Layer: output is g_i and input is b_i , where $i = 1, 2 \dots 26$

V is 200×785 weight matrix, V_{ij} is the weight from x_j to h_i

W is 26×201 weight matrix, W_{ij} is the weight from h_j to g_i

$h_i = t(a_i) = t(V_i x) = t(\sum_{j=1}^{785} V_{ij} x_j)$ where t is tanh function

$g_i = s(b_i) = s(W_i h) = s(\sum_{j=1}^{201} W_{ij} h_j)$ where s is softmax function

$$\frac{dt(a)}{da} = 1 - t^2(a)$$

$$\frac{ds(b)}{db} = s(b) \cdot (1 - s(b))$$

$$\frac{dL(y, s(b))}{db} = b - y$$

W :

$$\begin{aligned} \frac{dL}{dw_{ij}} &= \frac{dL}{db_i} \cdot \frac{db_i}{dw_{ij}} \\ &= (s(W_i h) - y_i) \cdot h_j \\ &= (s(W_i h) - y_i) \cdot t(V_j x) \end{aligned}$$

matrix form:

$$\begin{aligned} \frac{dL}{dW} &= (b - y)^T \cdot h \\ &= (Wh - y)^T \cdot h \\ &= (Wt(Vx) - y)^T \cdot t(Vx) \end{aligned}$$

V :

$$\begin{aligned} \frac{dL}{dv_{ij}} &= \frac{dL}{dh_i} \cdot \frac{dh_i}{dv_{ij}} = \frac{dh_i}{dv_{ij}} \cdot \sum_{j=1}^{26} \left(\frac{dL}{db_j} \cdot \frac{b_j}{h_i} \right) \\ &= x_j \cdot (1 - t^2(a_i)) \cdot \sum_{j=1}^{26} ((s(b_j) - y_j) \cdot w_{ji}) \end{aligned}$$

matrix form:

$$\begin{aligned} \frac{dL}{dV} &= W^T (g - y) (1 - h^2) x \\ &= W^T (s(Wh) - y) (1 - h^2) x \\ &= W^T (s(Wt(Vx)) - y) (1 - t(Vx)^2) x \end{aligned}$$

Stochastic Gradient Descent:

$$\begin{aligned} w_{ij} &= w_{ij} - \epsilon \cdot \frac{dL}{dw_{ij}} \\ v_{ij} &= v_{ij} - \epsilon \cdot \frac{dL}{dv_{ij}} \end{aligned}$$