

David Freilich

1. Derivation of Soft-max

$$y_j(x) = \text{softmax}(f(x)) = \frac{e^{f(x_j)}}{\sum_k e^{f(x_k)}}$$

$$f(x) = Wx + b$$

$$(1) \frac{\partial y_i}{\partial f_j} = \frac{\partial e^{f(x_i)}}{\sum_k e^{f(x_k)}}$$

$$(2) \text{Through quotient rule, get: } \frac{e^{f(x_i)} \sum_k e^{f(x_k)} - e^{f(x_i)} e^{f(x_j)}}{(\sum_k e^{f(x_k)})^2}$$

$$(3) \text{if } i=j, \text{ get: } \frac{e^{f(x_i)} \cdot \sum_k e^{f(x_k)} - e^{f(x_i)} e^{f(x_i)}}{\sum_k e^{f(x_k)}} = y_i \cdot (1 - y_i)$$

$$(4) \text{if } i \neq j, \frac{\partial e^{f(x_i)}}{\partial f_j} = 0, \text{ so get: } \frac{0 \cdot \sum_k e^{f(x_k)} - e^{f(x_i)} e^{f(x_j)}}{(\sum_k e^{f(x_k)})^2} = \frac{-e^{f(x_i)} e^{f(x_j)}}{\sum_k e^{f(x_k)}} = -y_i \cdot y_j$$

2. Derivation of Loss

$$(1) L(X, Y; W, b) = -\frac{1}{N} \sum_j \sum_i \overset{\text{ground truth}}{Y'_{ji}} \log Y_{ji} \quad \left| \begin{array}{l} Y'_{ji} = 1 \text{ iff image } j \in \text{class } i \\ Y_{ji} = \text{prob that image } j \in \text{class } i \end{array} \right.$$

$$(2) \frac{\partial L}{\partial f_j} = - \sum_i \frac{\partial Y'_{ji} \cdot \log(Y_{ji})}{\partial x_i} = - \sum_i Y'_{ji} \cdot \frac{\partial \log(Y_{ji})}{\partial f_j}$$

$$(3) \text{through log derivation, we get: } - \sum_j Y'_{ji} \cdot \frac{1}{Y_{ji}} \cdot \frac{\partial Y_{ji}}{\partial f_j} = - \left(\sum_j \frac{Y'_{ji}}{Y_{ji}} \cdot \frac{\partial Y_{ji}}{\partial f_j} \right)$$

$$(4) \text{Split summation, taking out case } j=i \text{ so we get: } - \frac{Y'_{ii}}{Y_{ii}} \frac{\partial Y_{ii}}{\partial f_i} - \sum_{j \neq i} \frac{Y'_{ji}}{Y_{ji}} \cdot \frac{\partial Y_{ji}}{\partial f_i}$$

$$(5) \text{Fill in } \frac{\partial Y_{ji}}{\partial x_i} \text{ softmax, so we get: } \frac{-Y'_{ii} \cdot y_i (1 - y_i)}{y_i} - \sum_{j \neq i} \frac{Y'_{ji}}{y_j} \cdot (-y_i y_j)$$

$$(6) \frac{\partial L}{\partial f_j} = \frac{-y'_i}{y_i} (y_i(1-y_i)) - \sum_{j \neq i} \frac{y'_i}{y_i} (-y_i y_j) = -y'_i + y'_i y_i + \sum_{j \neq i} y'_i y_j$$

(7) Combine $j \neq i$ w/ summation, so we get:

$$\frac{\partial L}{\partial f_j} = -y'_i + \sum_{j=1} y'_i y_j = -y'_i + y_i \sum_{j=1} y'_j$$

(8) Since y' is one-hot-encoded, $\sum y'_j = 1$, so we get:

$$\frac{\partial L}{\partial f_j} = y_i - y'_i$$

3. Derive weight gradients

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial f_j} \cdot \frac{\partial f_j}{\partial W_j} = (y_i - y'_i) \cdot \frac{\partial (Wx+b)}{\partial W} = (y_i - y'_i) (x)$$

$f(x) = Wx+b$

$$\frac{\partial L}{\partial b_j} = \frac{\partial L}{\partial f_j} \cdot \frac{\partial f_j}{\partial b_j} = (y_i - y'_i) \cdot \frac{\partial (Wx+b)}{\partial b} = (y_i - y'_i)$$