

NLP Stanford - Assignment 01

1 Softmax

(a) Prove that $\text{softmax}(x) = \text{softmax}(x+c)$

$$\begin{aligned}\text{softmax}(x+c) &= \frac{e^{x+c}}{\sum_j e^{x_j+c}} \\ &= \frac{e^x e^c}{\sum_j e^{x_j} e^c} \\ &= \frac{e^x}{\sum_j e^{x_j}} \\ &= \text{softmax}(x)\end{aligned}$$

2 Neural Networks Basics

(a) Derive the Sigmoid function

$$\sigma(x) = \frac{1}{1+e^{-x}} \quad f(x) = \frac{y(x)}{h(x)} \quad f'(x) = \frac{y'(x)h(x) - y(x)h'(x)}{h(x)^2}$$

$$\sigma'(x) = \frac{y(x)}{h(x)} = \frac{1}{1+e^{-x}}$$

$$g(x) = \frac{\partial}{\partial x} 1 = 0 \quad h(x) = \frac{\partial}{\partial x} (1+e^{-x}) = \frac{\partial}{\partial x} 1 + \frac{\partial}{\partial x} e^{-x} = 0 + e^{-x} = -e^{-x}$$

$$\begin{aligned}\sigma'(x) &= \frac{0(1+e^{-x}) - 1(-e^{-x})}{(1+e^{-x})^2} = \frac{e^{-x}}{(1+e^{-x})^2} = \left(\frac{1}{1+e^{-x}}\right) \left(\frac{e^{-x}}{1+e^{-x}}\right) \\ &= \left(\frac{1}{1+e^{-x}}\right) \left(\frac{1+e^{-x}-1}{1+e^{-x}}\right) = \frac{1}{1+e^{-x}} \left(1 - \frac{1}{1+e^{-x}}\right) \\ &= \sigma(x)(1-\sigma(x))\end{aligned}$$

(b) Derive the Cross-Entropy function.

$$\mathcal{L}(y, \hat{y}) = - \sum_i y_i \log(\hat{y}_i)$$

$$\hat{y}_i = \text{softmax}(a_i) = \frac{e^{a_i}}{\sum_j e^{a_j}} \quad \log \frac{y}{z} = \log y - \log z$$

$$\mathcal{L}(y, \hat{y}) = - \sum_i y_i \log \left(\frac{e^{a_i}}{\sum_j e^{a_j}} \right) = - \sum_i y_i \left(\log e^{a_i} - \log \sum_j e^{a_j} \right)$$

$$(a = f(z)) = \sum_i y_i \log(\hat{y}_i)$$

$$\hat{y} = \text{softmax}(a)$$

$$\begin{aligned}\mathcal{L}(y, \hat{y}) &= - \sum_i y_i \log \left(\frac{e^{a_i}}{\sum_j e^{a_j}} \right) \\ &= - \sum_i y_i \left(\log e^{a_i} - \log \sum_j e^{a_j} \right) \\ &= - \sum_i y_i a_i + \sum_i y_i \log \sum_j e^{a_j} \\ &\quad (1) \quad (2)\end{aligned}$$

$$(1) \frac{\partial}{\partial a_i} \sum_j y_j a_j = y_i$$

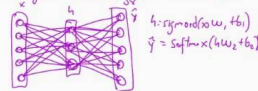
$$(2) \frac{\partial}{\partial a_i} \sum_j y_j \log \sum_k e^{a_k} = \frac{\partial}{\partial a_i} \log \sum_k e^{a_k} \quad f(a) = \log u \quad u = \sum_k e^{a_k}$$

$$\begin{aligned}\frac{\partial}{\partial a_i} \log \sum_k e^{a_k} &= \frac{\partial}{\partial u} \log u \cdot \frac{\partial u}{\partial a_i} \\ &= \frac{1}{u} \cdot \frac{\partial}{\partial a_i} \sum_k e^{a_k} \\ &= \frac{1}{\sum_k e^{a_k}} \cdot \frac{\partial}{\partial a_i} \sum_k e^{a_k} \\ &= \frac{1}{\sum_k e^{a_k}} e^{a_i} = \frac{e^{a_i}}{\sum_k e^{a_k}}\end{aligned}$$

$$\text{softmax}(a)_i = \hat{y}_i$$

$$\begin{aligned}\frac{\partial \mathcal{L}(y, \hat{y})}{\partial a_i} &= \frac{\partial}{\partial a_i} \left(- \sum_j y_j a_j + \sum_j y_j \log \sum_k e^{a_k} \right) \\ &= -y_i + \frac{e^{a_i}}{\sum_k e^{a_k}} = -y_i + \text{softmax}(a_i) = \hat{y}_i - y_i \\ &= \hat{y}_i - y_i\end{aligned}$$

$$\Rightarrow \frac{\partial}{\partial b} \mathcal{L}(y, \hat{y}) = \hat{y} - y$$

(c) Derivate the product with respect to the weights for an hidden layer neural network, $\frac{\partial \mathcal{L}}{\partial x}$ where $\mathcal{L} = \mathcal{L}(y, \hat{y})$ 

$$h_i = \text{sigmoid}(xw_1 + b_1) = \text{sigmoid}(z_i) \quad z_i = xw_1 + b_1$$

$$\hat{y} = \text{softmax}(hw_2 + b_2) = \text{softmax}(z_2) \quad z_2 = hw_2 + b_2$$

$$\frac{\partial \mathcal{L}(y, \hat{y})}{\partial x} = \frac{\partial \mathcal{L}(y, \hat{y})}{\partial z_2} \frac{\partial z_2}{\partial h} \frac{\partial h}{\partial z_1} \frac{\partial z_1}{\partial x} \quad \text{by chain rule}$$

$$\text{Just do for } i=1$$

$$\frac{\partial \mathcal{L}(y, \hat{y})}{\partial z_2} = \hat{y}_i - y_i$$

$$\text{So becomes } \frac{\partial \mathcal{L}(y, \hat{y})}{\partial z_2} \frac{\partial z_2}{\partial h} \frac{\partial h}{\partial z_1} \frac{\partial z_1}{\partial x} \rightarrow \mathbb{R}^m \rightarrow \mathbb{R}^m$$

$$\text{Jacobian} = \frac{\partial z_2}{\partial h} = \begin{bmatrix} \frac{\partial z_2^1}{\partial h_1} & \dots & \frac{\partial z_2^m}{\partial h_m} \\ \vdots & & \vdots \end{bmatrix} \quad \frac{\partial \mathcal{L}}{\partial z_2} \text{ (constant)} = z_2 = z^1$$

for convenience: $W_2 = W^{(2)}$; $W_3 = W^{(3)}$ weights of neurons hidden
 $b_1 = b^{(1)}$ for output
 $z_1 = \underbrace{[w_{11}x_1 + w_{12}x_2 + \dots + w_{1n}x_n + b_1]}_{z_1} \dots \underbrace{[w_{m1}x_1 + w_{m2}x_2 + \dots + w_{mn}x_n + b_m]}_{z_m}$

$$\frac{\partial z_1}{\partial x} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1n} \\ w_{21} & w_{22} & \dots & w_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m1} & w_{m2} & \dots & w_{mn} \end{bmatrix} = W_2^T$$

• sigmoid(z) = $\sigma(z) = \frac{1}{1+e^{-z}}$; $\frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1-\sigma(z))$

for convenience: $z_1 = z^{(1)}$
 $\frac{\partial z_1}{\partial z_1} = \sigma'(z_1)(1-\sigma'(z_1))$

• $z_1 = XW_1 + b_1$, for convenience $z_1 = z^{(1)}$; $z^{(1)}: \mathbb{R}^n \rightarrow \mathbb{R}^m$

Jacobian $\frac{\partial z^{(1)}}{\partial X} = \begin{bmatrix} \frac{\partial z_1^{(1)}}{\partial x_1} & \dots & \frac{\partial z_1^{(1)}}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_m^{(1)}}{\partial x_1} & \dots & \frac{\partial z_m^{(1)}}{\partial x_n} \end{bmatrix} = W_1^T$ W_1^T for input to hidden neurons
 $b_1 = b^{(1)}$

$z^{(1)} = [x_1w_{11} + x_2w_{12} + \dots + x_nw_{1n} + b_1] \dots [x_1w_{m1} + x_2w_{m2} + \dots + x_nw_{mn} + b_m]$

$$\frac{\partial z^{(1)}}{\partial X} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1n} \\ w_{21} & w_{22} & \dots & w_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m1} & w_{m2} & \dots & w_{mn} \end{bmatrix} = W_1^T = W_1^T$$

$$\frac{\partial CE(y, \hat{y})}{\partial X} = \frac{\partial CE(y, \hat{y})}{\partial z_1} \cdot \frac{\partial z_1}{\partial x_1} \cdot \frac{\partial z_1}{\partial x_2} \dots \frac{\partial z_1}{\partial x_n}$$

$-(\hat{y} - y) W_2^T \sigma'(z_1)(1-\sigma'(z_1)) W_1^T$

(a) How many parameters are there in this neural network assuming the input is 2-dimensional, the output is 0, 1-dimensional and there are 11 hidden units?

#parameters: $D_X H + H + H D_Y + b_1$
 $= (D_X + 1)H + (H + 1)D_Y$

(b) W_1 for a neural network with a single hidden layer. Calculate the derivatives $\frac{\partial CE}{\partial w_1}, \frac{\partial CE}{\partial b_1}, \frac{\partial CE}{\partial w_2}, \frac{\partial CE}{\partial b_2}$

$z_1 = \text{sigmoid}(w_1x + b_1)$
 $\hat{y} = \text{softmax}(w_2z_1 + b_2)$
 $CE(y, \hat{y}) = -\sum_i y_i \log(\hat{y}_i)$

$$\frac{\partial CE(y, \hat{y})}{\partial w_2} = \frac{\partial CE(y, \hat{y})}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_2} = (\hat{y} - y) \frac{\partial z_1}{\partial w_2}$$

$z_1 = z^{(1)}: \mathbb{R}^n \rightarrow \mathbb{R}^m$

for convenience: $W_2 = W^{(2)}$; $W_3 = W^{(3)}$ weights of neurons hidden
 $b_1 = b^{(1)}$ for output
 $z_1 = \underbrace{[w_{11}x_1 + w_{12}x_2 + \dots + w_{1n}x_n + b_1]}_{z_1} \dots \underbrace{[w_{m1}x_1 + w_{m2}x_2 + \dots + w_{mn}x_n + b_m]}_{z_m}$

$W_1^{(1)} = \begin{bmatrix} w_{11}^{(1)} & w_{12}^{(1)} & \dots & w_{1n}^{(1)} \\ w_{21}^{(1)} & w_{22}^{(1)} & \dots & w_{2n}^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m1}^{(1)} & w_{m2}^{(1)} & \dots & w_{mn}^{(1)} \end{bmatrix}$ $z_1^{(1)}: \mathbb{R}^n \rightarrow \mathbb{R}^m$ $\frac{\partial z_1^{(1)}}{\partial w_1^{(1)}} \in \mathbb{R}^{m \times n}$
 $b_1 \in \mathbb{R}^m$

$\frac{\partial \hat{y}}{\partial w_2^{(2)}} = h^T: \mathbb{R}^m \rightarrow \mathbb{R}^n$ $\frac{\partial \hat{y}}{\partial w_2^{(2)}} = \begin{bmatrix} \frac{\partial \hat{y}_1}{\partial w_{21}^{(2)}} & \dots & \frac{\partial \hat{y}_n}{\partial w_{2n}^{(2)}} \end{bmatrix}$
 $\frac{\partial CE(y, \hat{y})}{\partial w_2} = h^T(\hat{y} - y)$

$\frac{\partial CE(y, \hat{y})}{\partial b_2} = \frac{\partial CE(y, \hat{y})}{\partial z_1} \cdot \frac{\partial z_1}{\partial b_2} = (\hat{y} - y)$

$\frac{\partial z_1}{\partial b_1}$, $z_1: \mathbb{R}^n \rightarrow \mathbb{R}^m$ $\hat{y} = \frac{\partial z_1}{\partial b_2} \in \mathbb{R}^{m \times m}$; $b_2 = b^{(2)}$; $z_2 = z^{(2)}$

$z_1 = \underbrace{[w_{11}x_1 + w_{12}x_2 + \dots + w_{1n}x_n + b_1]}_{z_1} \dots \underbrace{[w_{m1}x_1 + w_{m2}x_2 + \dots + w_{mn}x_n + b_m]}_{z_m}$

$J = \begin{bmatrix} \frac{\partial z_1^{(1)}}{\partial w_{11}^{(1)}} & \frac{\partial z_1^{(1)}}{\partial w_{12}^{(1)}} & \dots & \frac{\partial z_1^{(1)}}{\partial w_{1n}^{(1)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_m^{(1)}}{\partial w_{m1}^{(1)}} & \frac{\partial z_m^{(1)}}{\partial w_{m2}^{(1)}} & \dots & \frac{\partial z_m^{(1)}}{\partial w_{mn}^{(1)}} \end{bmatrix} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix} = I$

$\frac{\partial CE(y, \hat{y})}{\partial b_2} = \frac{\partial CE(y, \hat{y})}{\partial z_1} \cdot \frac{\partial z_1}{\partial b_2} = (\hat{y} - y) I = \hat{y} - y$

• $\frac{\partial CE(y, \hat{y})}{\partial w_1}$, $W_1 = W^{(1)}$

$\frac{\partial CE(y, \hat{y})}{\partial w_1} = \frac{\partial CE(y, \hat{y})}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1} = (\hat{y} - y) W_2^T \sigma'(z_1)(1-\sigma'(z_1)) \frac{\partial z_1}{\partial w_1}$

$z_1: \mathbb{R}^n \rightarrow \mathbb{R}^m$, $\frac{\partial z_1}{\partial w_1} = h^T \Rightarrow \frac{\partial z_1}{\partial w_1} = X^T$ $\begin{pmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{pmatrix}$

Jacobian $\frac{\partial CE(y, \hat{y})}{\partial w_1}$:
 $\frac{\partial CE(y, \hat{y})}{\partial w_1} = X^T ((\hat{y} - y) W_2^T \sigma'(z_1)(1-\sigma'(z_1)))$ $W_1 \in \mathbb{R}^{m \times n}$

$= X^T ((\hat{y} - y) W_2^T \sigma'(z_1)(1-\sigma'(z_1)))$
 $= X^T ((\hat{y} - y) W_2^T \sigma'(z_1)(1-\sigma'(z_1)))$
 $= X^T ((\hat{y} - y) W_2^T \sigma'(z_1)(1-\sigma'(z_1)))$

• $\frac{\partial CE(y, \hat{y})}{\partial b_1} = \frac{\partial CE(y, \hat{y})}{\partial z_1} \cdot \frac{\partial z_1}{\partial b_1} = (\hat{y} - y)$
 $z_1: \mathbb{R}^n \rightarrow \mathbb{R}^m$ $\frac{\partial z_1}{\partial b_1} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix} = I$

• $\frac{\partial CE(y, \hat{y})}{\partial b_1} = (\hat{y} - y) W_2^T \sigma'(z_1)(1-\sigma'(z_1))$
 $= (\hat{y} - y) W_2^T \sigma'(z_1)(1-\sigma'(z_1))$