

Softmax

1c) Prove that $\text{softmax}(x) = \text{softmax}(x+c)$ where $\text{softmax}(x)_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$

$$\text{softmax}(x+c)_i = \frac{e^{x_i+c}}{\sum_j e^{x_j+c}}$$

$$\Leftrightarrow \text{softmax}(x+c)_i = \frac{e^{x_i} \cdot e^c}{\sum_j e^{x_j} \cdot e^c}$$

$$\Leftrightarrow \text{softmax}(x+c)_i = \frac{e^{x_i} \cdot e^c}{e^c \sum_j e^{x_j}}$$

$$\Leftrightarrow \text{softmax}(x+c)_i = \frac{e^c}{e^c} \cdot \frac{e^{x_i}}{\sum_j e^{x_j}}$$

$$\Leftrightarrow \text{softmax}(x+c)_i = 1 \cdot \frac{e^{x_i}}{\sum_j e^{x_j}} = \text{softmax}(x)_i$$

Neural Network Basics

(2a) Derive the gradients of the sigmoid function and show that it can be rewritten as a function of the function value

The sigmoid function is given by

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

But we can also rewrite it as

$$\sigma(x) = v(z(x)) \quad \text{where} \quad v(z) = \frac{1}{z}$$
$$z(x) = 1 + e^{-x}$$

The derivative can thus be written as

$$\sigma'(x) = v'(z(x)) \cdot z'(x) \quad \text{where} \quad v'(z) = -\frac{1}{z^2}$$
$$z'(x) = -e^{-x}$$

Therefore, inserting $v'(z)$ and $z'(x)$ we get:

$$\sigma'(x) = -\frac{1}{(1 + e^{-x})^2} \cdot -e^{-x}$$

$$\Leftrightarrow \sigma'(x) = \frac{e^{-x}}{(1 + e^{-x})^2}$$

Keeping in mind that $\sigma(x) = \frac{1}{1 + e^{-x}}$ and considering that we can rewrite e^{-x} as $e^{-x} = \frac{1}{\sigma(x)} - 1$ we can rewrite $\sigma'(x)$ as a function of $\sigma(x)$

$$\sigma'(x) = \sigma(x)^2 \cdot \left(\frac{1}{\sigma(x)} - 1 \right)$$

$$\Leftrightarrow \sigma'(x) = \sigma(x) - \sigma(x)^2$$

$$\Leftrightarrow \sigma'(x) = \sigma(x) \cdot (1 - \sigma(x))$$

Neural Networks

- 2b) Derive the gradient with regards to the inputs of a softmax function when cross entropy loss is used for evaluation, i.e. find the gradients with respect to the softmax input vector θ , when the prediction is made by $\hat{y} = \text{softmax}(\theta)$

$$\text{Cross Entropy: } CE(y, \hat{y}) = - \sum_i y_i \log(\hat{y}_i)$$

$$\text{Softmax: } \text{softmax}(x) = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

Instead of writing θ for all the inputs we will take the derivative of input vector x

$$\frac{\partial CE}{\partial x_h} = - \sum_i y_i \frac{1}{\hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial x_h}$$

Therefore we first need the derivative of the softmax function with respect to x_h . For this we will differentiate between the two cases where $i=h$ and $i \neq h$

Case 1: $i = h$

$$\frac{\partial \hat{y}_i}{\partial x_h} = \frac{e^{x_h}}{\sum_j e^{x_j}} - \frac{e^{x_h}}{(\sum_j e^{x_j})^2} \cdot e^{x_h}$$

$$\frac{\partial \hat{y}_i}{\partial x_h} = \hat{y}_h - \hat{y}_h^2$$

$$\frac{\partial \hat{y}_i}{\partial x_h} = \hat{y}_h (1 - \hat{y}_h)$$

Case 2: $i \neq h$

$$\frac{\partial \hat{y}_i}{\partial x_h} = - \frac{e^{x_i}}{(\sum_j e^{x_j})^2} \cdot e^{x_h}$$

$$\frac{\partial \hat{y}_i}{\partial x_h} = - \frac{e^{x_i}}{\sum_j e^{x_j}} \cdot \frac{e^{x_h}}{\sum_j e^{x_j}}$$

$$\frac{\partial \hat{y}_i}{\partial x_h} = - \hat{y}_i \hat{y}_h$$

Plugging in those derivatives we can get the following

$$\frac{\partial CE}{\partial x_h} = - y_h \frac{1}{\hat{y}_h} \cdot (\hat{y}_h (1 - \hat{y}_h)) - \sum_{i \neq h} y_i \cdot \frac{1}{\hat{y}_i} \cdot (- \hat{y}_i \hat{y}_h)$$

$$\Leftrightarrow \frac{\partial CE}{\partial x_h} = - y_h + y_h \hat{y}_h + \sum_{i \neq h} y_i \hat{y}_h$$

$$\Leftrightarrow \frac{\partial CE}{\partial x_h} = - y_h + \sum_i y_i \hat{y}_h$$

cont'd \rightarrow

$$\frac{\partial CE}{\partial x_h} = -y_h + \hat{y}_h \underbrace{\sum_i y_i}$$

This equals 1 because there is only
one career class label

$$\frac{\partial CE}{\partial x_h} = \hat{y}_h - y_h$$

4/24/20

Something missing

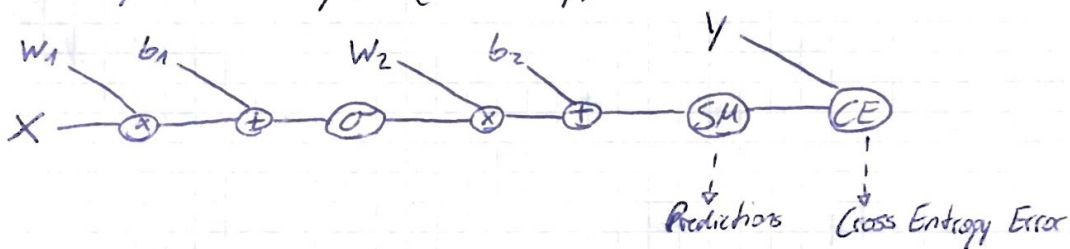
Neural Network Basics

- (2c) Derive the gradients with respect to the inputs x to a one-hidden layer neural network that is, find $\frac{\partial J}{\partial x}$ where $J = CE(y, \hat{y})$ is the cost function for the neural network. The neural network employs sigmoid activation function for the hidden layer, and softmax for the output layer. Assume the one-hot label vector is y , and cross entropy is used. (Feel free to use $\sigma(x)$ as the shorthand for the sigmoid gradient, and feel free to define any variables whenever you see fit.)

The given functions for the neural net are:

$$J = CE(y, \hat{y}) \quad \hat{y} = \text{softmax}(hW_2 + b_2) \quad h = \text{sigmoid}(xW_1 + b_1)$$

The Computation Graph (Visual help) looks like this:



For the purpose of the task I will rewrite the functions as follows:

$$J = CE(y, \text{softmax}(\hat{y})) \quad \hat{y} = hW_2 + b_2 \quad h = \text{sigmoid}(z) \quad z = xW_1 + b_1$$

The Gradients with respect to x can thus be defined as follows

$$\begin{aligned} \frac{\partial J}{\partial x} &= \frac{\partial CE}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial x} \\ &= \frac{\partial CE}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h} \frac{\partial h}{\partial x} \\ &= \frac{\partial CE}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h} \frac{\partial h}{\partial z} \frac{\partial z}{\partial x} \end{aligned} \quad \text{where} \quad \begin{aligned} \frac{\partial CE}{\partial \hat{y}} &= \hat{y} - y \\ \frac{\partial \hat{y}}{\partial h} &= W_2 \\ \frac{\partial h}{\partial z} &= \sigma'(z) \\ \frac{\partial z}{\partial x} &= W_1 \end{aligned}$$

The Gradients are thus given by

$$\begin{aligned} \delta_1 &= \hat{y} - y \\ \delta_2 &= \delta_1 \cdot W_2^T \\ \delta_3 &= \delta_2 \odot \sigma'(z) \\ \delta_4 &= \delta_3 \cdot W_1^T \end{aligned}$$

$$\Rightarrow \frac{\partial J}{\partial x} = \delta_4 \cdot W_1^T$$

Neural Network Basics

(2d) How many parameters are there in this neural network (2c), assuming the input is D_x -dimensional, the output is D_y dimensional and there are H hidden units?

The trainable parameters of this Network are given by W_1 , b_1 , W_2 , and b_2

The Dimensions of the input X and Output Y for a batch size of n are:

$$X \in \mathbb{R}^{n \times D_x}$$

$$Y \in \mathbb{R}^{D_y \times n}$$

The dimensions of our parameters must thus be as follows

$$W_1 \in \mathbb{R}^{D_x \times H}$$

$$W_2 \in \mathbb{R}^{H \times D_y}$$

$$b_1 \in \mathbb{R}^{1 \times H}$$

$$b_2 \in \mathbb{R}^{1 \times D_y}$$

The total amount of parameters P is thus:

$$P = D_x \cdot H + H + H \cdot D_y + D_y$$

$$P = H \cdot (1 + D_x) + D_y \cdot (1 + H)$$

Word2Vec

(3a) Q: Assume you are given a "predicted" word vector v_c corresponding to the center word c for Ship-Gren, and word prediction is made with the softmax function found in word2vec models.

$$\hat{y}_0 = p(u_0 | v_c) = \frac{\exp(u_0^T v_c)}{\sum_w^V \exp(u_w^T v_c)}$$

where u_w ($w=1, \dots, V$) are the "input" word vectors for all words in the vocabulary. Assuming cross entropy cost is applied to the prediction and word o is the predicted (the o -th element of the one-hot vector is one), derive the gradients with respect to v_c .

A: We rewrite the functions $CE(y, \hat{y})$ and \hat{y}_0 as follows:

$$CE(y, \hat{y})_0 = -\sum y \cdot \hat{y}_0 \quad \hat{y}_0 = \log \left(\frac{\exp(u_0^T v_c)}{\sum_w^V \exp(u_w^T v_c)} \right) = u_0^T v_c - \log \sum_w^V \exp(u_w^T v_c)$$

$$\text{Then: } \frac{\partial CE}{\partial v_c} = -\sum y \frac{\partial \hat{y}_0}{\partial v_c}$$

$$\frac{\partial \hat{y}_0}{\partial v_c} = u_0 - \frac{1}{\sum_w^V \exp(u_w^T v_c)} \cdot \sum_x^V \exp(u_x^T v_c) \cdot u_x$$

$$\frac{\partial \hat{y}_0}{\partial v_c} = u_0 - \sum_x^V \frac{\exp(u_x^T v_c)}{\sum_w^V \exp(u_w^T v_c)} \cdot u_x$$

$$\frac{\partial \hat{y}_0}{\partial v_c} = u_0 - \sum_x^V \hat{y}_x \cdot u_x$$

Inserting this into $\frac{\partial CE}{\partial v_c}$ we get our result

$$\frac{\partial CE}{\partial v_c} = -\sum y \cdot u_0 - \sum_x^V \hat{y}_x u_x$$

This is always 0 except in our case where o is the correct expected word

$$\Leftrightarrow \frac{\partial CE}{\partial v_c} = -u_0 + \sum_x^V \hat{y}_x u_x$$

Word2Vec

(3b) Q: As in the previous part, derive gradients for the "output" word vectors u_h 's (including u_0)

A: We again rewrite the functions $CE(y, \hat{y})$ and \hat{y} as follows

$$CE(y, \hat{y}) = - \sum_i y_i \hat{y}_i \quad \hat{y}_i = \log \frac{\exp(u_i^T v_c)}{\sum_w \exp(u_w^T v_c)} = u_i^T v_c - \log \sum_w \exp(u_w^T v_c)$$

The derivative wrt u_h is given by

$$\frac{\partial CE}{\partial u_h} = - \sum_i y_i \frac{\partial \hat{y}_i}{\partial u_h}$$

To calculate the derivative of \hat{y} we have to differentiate between the two cases where $k=i$ and $k \neq i$

For $k=i$:

$$\begin{aligned} \frac{\partial \hat{y}_i}{\partial u_h} &= v_c - \frac{1}{\sum_w \exp(u_w^T v_c)} \cdot \exp(u_h^T v_c) \cdot v_c \\ &= v_c - \frac{\exp(u_h^T v_c)}{\sum_w \exp(u_w^T v_c)} \cdot v_c \\ &= v_c \cdot (1 - \hat{y}_h) \end{aligned}$$

For $k \neq i$:

$$\begin{aligned} \frac{\partial \hat{y}_i}{\partial u_h} &= - \frac{1}{\sum_w \exp(u_w^T v_c)} \cdot \exp(u_h^T v_c) \cdot v_c \\ &= - v_c \hat{y}_h \end{aligned}$$

Cont'd \rightarrow

Inserting these derivatives into $\frac{\partial CE}{\partial u_h}$ we get the following:

$$\frac{\partial CE}{\partial u_h} = -y_h \cdot v_L (1 - \hat{y}_h) + \sum_{i \neq h} y_i \hat{y}_h v_L$$

$$= -y_h v_L + y_h v_L \hat{y}_h + \sum_{i \neq h} y_i \hat{y}_h v_L$$

$$= -y_h v_L + \sum_i y_i \hat{y}_h v_L$$

$$= -y_h v_L + \hat{y}_h v_L \cdot \underbrace{\sum_i y_i}_{=1}$$

$$\boxed{\frac{\partial CE}{\partial u_h} = v_L (\hat{y} - y)^T}$$