

Assignment 4

Fudi(Fred) Wang

November 2019

Part 1

(a) Since $\mathbf{y}^{(t)}$ is a one hot vector, without loss of generality, let the index such that $y_j^{(t)} = 1$ be i , we have:

$$\log PP^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = \log \frac{1}{\sum_{j=1}^{|V|} y_j^{(t)} \cdot \hat{y}_j^{(t)}} = -\log(\hat{y}_i^{(t)}) = -\sum_{j=1}^{|V|} y_j^{(t)} \cdot \log \hat{y}_j^{(t)} = J^{(t)}(\theta)$$

Hence, minimizing the mean cross-entropy loss will also minimize the mean perplexity across the training set.

Now, since model predictions were completely random (chosen uniformly from the vocabulary), and the sum of the $\hat{y}_i^{(t)}$ s is equal to one, we have:

$$E(\hat{y}_i^{(t)}) = \frac{1}{|V|} \quad \text{for } i \in (1, |V|)$$

Therefore,

$$E(PP^{(t)}) = |V|$$

So the corresponding cross-entropy loss when $|V| = 10000$ is just

$$E(J^{(t)}(\theta)) = \log 10000 = 4$$

Of course, if log is base 2, then we get the number 13.29.

(b) First, let $\mathbf{x}^{(t)} = \mathbf{h}^{(t)}\mathbf{U} + \mathbf{b}_2$.

The softmax function is defined as:

$$\text{softmax}(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}} = \frac{1}{1 + e^{-x_i} \sum_{j \neq i} e^{x_j}}$$

where the second equality holds by dividing both the numerator and the denominator by e^{x_i} . Note that the derivative of the softmax function can be discussed as two cases:

$$\frac{\partial \text{softmax}(x_i)}{\partial x_i} = \frac{e^{-x_i} \sum_{j \neq i} e^{x_j}}{(1 + e^{-x_i} \sum_{j \neq i} e^{x_j})^2} = \text{softmax}(x_i)(1 - \text{softmax}(x_i))$$

and

$$\frac{\partial \text{softmax}(x_i)}{\partial x_j} = -\frac{e^{x_i} e^{x_j}}{(\sum_j e^{x_j})^2} = -\text{softmax}(x_i) \text{softmax}(x_j)$$

And recall that from previous assignments, $\sigma'(x) = \sigma(x)(1 - \sigma(x))$. Hence we have the following partial derivatives:

$$\begin{aligned}\frac{\partial J^{(t)}}{\partial \mathbf{b}_2} &= \frac{\partial J^{(t)}}{\partial \hat{\mathbf{y}}^{(t)}} \frac{\partial \hat{\mathbf{y}}^{(t)}}{\partial \mathbf{x}^{(t)}} \frac{\partial \mathbf{x}^{(t)}}{\partial \mathbf{b}_2} = -\frac{1}{\hat{\mathbf{y}}^{(t)}} \hat{\mathbf{y}}^{(t)} (1 - \hat{\mathbf{y}}^{(t)}) = \hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)} \\ \frac{\partial J^{(t)}}{\partial \mathbf{L}_{x^{(t)}}} &= \frac{J^{(t)}}{\partial \hat{\mathbf{y}}^{(t)}} \frac{\hat{\mathbf{y}}^{(t)}}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial e^{(t)}} \frac{\partial e^{(t)}}{\partial L_{x^{(t)}}} = (\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)}) \mathbf{U}^T \mathbf{h}^{(t)} (1 - \mathbf{h}^{(t)}) \mathbf{I}^T \\ \frac{\partial J^{(t)}}{\partial \mathbf{I}}|_{(t)} &= (\mathbf{e}^{(t)})^T (\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)}) \mathbf{U}^T \mathbf{h}^{(t)} (1 - \mathbf{h}^{(t)}) \\ \frac{\partial J^{(t)}}{\partial \mathbf{H}}|_{(t)} &= (\mathbf{h}^{(t-1)})^T (\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)}) \mathbf{U}^T \mathbf{h}^{(t)} (1 - \mathbf{h}^{(t)}) \\ \frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t-1)}} &= (\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)}) \mathbf{U}^T \mathbf{h}^{(t)} (1 - \mathbf{h}^{(t)}) \mathbf{H}^T\end{aligned}$$