

## Homework\_4

Shiqiang

(a)

Suppose  $y_i^t$  is the only nonzero element of  $y^t$ , then

$$J^{(t)}(\theta) = CE(\mathbf{y}^{(t)}, \mathbf{y}^{(t)}) = -\log \hat{y}_i^t = \log \frac{1}{\hat{y}_i^t}$$
$$PP(y^t, \hat{y}^t) = \frac{1}{\hat{y}_i^t}$$

Then,  $J^{(t)}(\theta) = \log PP(y^t, \hat{y}^t)$

If the model predictions are completely random,  $E[\hat{y}_i^t] = \frac{1}{|V|}$

The expected perplexity are  $E(PP(y^t, \hat{y}^t)) = E(\frac{1}{\hat{y}_i^t}) = 10000$

the expected cross-entropies are  $E(J^{(t)}(\theta)) = \log |V| = \log_2 10000 \approx 13.29$

(b)

As known, the soft max function is :

$$\hat{y}_i = \frac{\exp(o_i)}{\sum_j \exp(o_j)}$$

if  $i = j$ :

$$\frac{\partial \hat{y}_j}{\partial o_i} = \frac{\exp(o_i) \bullet \sum_j \exp(o_j) - (\exp(o_i))^2}{(\sum_j \exp(o_j))^2}$$
$$= \frac{\exp(o_i)}{\sum_j \exp(o_j)} (1 - \frac{\exp(o_i)}{\sum_j \exp(o_j)})$$
$$= \hat{y}_i (1 - \hat{y}_i)$$

if  $i \neq j$ :

$$\begin{aligned}\frac{\partial \hat{y}_j}{\partial o_i} &= \frac{-\exp(o_i) \bullet \exp(o_j)}{(\sum_j \exp(o_j))^2} \\ &= -\frac{\exp(o_i)}{\sum_j \exp(o_j)} \bullet \frac{\exp(o_j)}{\sum_j \exp(o_j)} \\ &= -\hat{y}_i \hat{y}_j\end{aligned}$$

The loss function is  $J^{(t)}(\theta) = CE(\mathbf{y}^{(t)}, \mathbf{\hat{y}}^{(t)}) = -\sum_{j=1}^{|V|} y_j^{(t)} \log \hat{y}_j^{(t)}$

$$\begin{aligned}\frac{\partial J^{(t)}}{\partial o_i} &= \frac{\partial(-\sum_{j=1}^{|V|} y_j^{(t)} \log \hat{y}_j^{(t)})}{\partial o_i} \\ &= -\sum_{j=1}^{|V|} y_j^{(t)} \frac{1}{\hat{y}_j^{(t)}} \frac{\partial \hat{y}_j^{(t)}}{\partial o_i} \\ &= -\sum_{i=j} \frac{y_i^{(t)}}{\hat{y}_j^{(t)}} \hat{y}_i^{(t)} (1 - \hat{y}_i^{(t)}) + \sum_{i \neq j} \frac{y_i^{(t)}}{\hat{y}_j^{(t)}} \hat{y}_i^{(t)} \hat{y}_j^{(t)} \\ &= -\sum_{i=j} y_i^{(t)} (1 - \hat{y}_i^{(t)}) + \sum_{i \neq j} y_i^{(t)} \hat{y}_i^{(t)} \\ &= \sum_{i \neq j} y_i^{(t)} \hat{y}_i^{(t)} - y_i^{(t)} + y_i^{(t)} \hat{y}_i^{(t)} \\ &= \hat{y}_i^{(t)} \sum_j y_i^{(t)} - y_i^{(t)} \\ &= \hat{y}_i^{(t)} - y_i^{(t)}\end{aligned}$$

Assume,  $\mathbf{g}^{(t)} = \mathbf{h}^{(t-1)} \mathbf{H} + \mathbf{e}^{(t)} \mathbf{I} + \mathbf{b}_1$   
 $\mathbf{o}^{(t)} = \mathbf{h}^{(t)} \mathbf{U} + \mathbf{b}_2$

Then,  $\frac{\partial J^{(t)}}{\partial o^{(t)}} = \hat{y}^{(t)} - y^{(t)}$

$$\frac{\partial J^{(t)}}{\partial \mathbf{g}^{(t)}} = \frac{\partial J^{(t)}}{\partial o^{(t)}} \frac{\partial o^{(t)}}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{g}^{(t)}} = (\hat{y}^{(t)} - y^{(t)}) \mathbf{U}^T \mathbf{h}^{(t)} (1 - h^{(t)})$$

$$\frac{\partial J^{(t)}}{\partial b_2} = \frac{\partial J^{(t)}}{\partial o^{(t)}} \frac{\partial o^{(t)}}{\partial b_2} = \hat{y}^{(t)} - y^{(t)}$$

$$\frac{\partial J^{(t)}}{\partial \mathbf{L}_{x^{(t)}}} = \frac{\partial J^{(t)}}{\partial \mathbf{g}^{(t)}} \frac{\partial \mathbf{g}^{(t)}}{\partial \mathbf{e}^{(t)}} \frac{\partial \mathbf{e}^{(t)}}{\partial \mathbf{L}_{x^{(t)}}} = (\hat{y}^{(t)} - y^{(t)}) \mathbf{U}^T \mathbf{h}^{(t)} (1 - h^{(t)}) \mathbf{I}^T$$

$$\left. \frac{\partial J^{(t)}}{\partial I} \right|_{(t)} = \frac{\partial J^{(t)}}{\partial g^{(t)}} \frac{\partial g^{(t)}}{\partial I} = (e^{(t)})^T (\hat{y}^{(t)} - y^{(t)}) \mathbf{U}^T h^{(t)} (1 - h^{(t)})$$

$$\left. \frac{\partial J^{(t)}}{\partial H} \right|_{(t)} = \frac{\partial J^{(t)}}{\partial g^{(t)}} \frac{\partial g^{(t)}}{\partial H} = (h^{(t-1)})^T (\hat{y}^{(t)} - y^{(t)}) \mathbf{U}^T h^{(t)} (1 - h^{(t)})$$

$$\frac{\partial J^{(t)}}{\partial h^{(t-1)}} = \frac{\partial J^{(t)}}{\partial g^{(t)}} \frac{\partial g^{(t)}}{\partial h^{(t-1)}} = (\hat{y}^{(t)} - y^{(t)}) \mathbf{U}^T h^{(t)} (1 - h^{(t)}) H^T$$