

## 0.1 Softmax [2 points]

- 1) [1 point] Prove that softmax is invariant to constant shifts in the input, i.e., for any input vector  $\mathbf{x}$  and a constant scalar  $c$ , the following holds:

$$\text{softmax}(\mathbf{x}) = \text{softmax}(\mathbf{x} + c),$$

where  $\text{softmax}(\mathbf{x})_i \triangleq \frac{e^{x_i}}{\sum_{i'} e^{x_{i'}}}$ , and  $\mathbf{x} + c$  means adding  $c$  to every dimension of  $\mathbf{x}$ .

- 2) [1 point] Let  $\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{c}$ , where  $\mathbf{W}$  and  $\mathbf{c}$  are some matrix and vector, respectively. Let

$$J = \sum_i \log \text{softmax}(\mathbf{z})_i.$$

Calculate the derivatives of  $J$  w.r.t.  $\mathbf{W}$  and  $\mathbf{c}$ , respectively, i.e., calculate  $\frac{\partial J}{\partial \mathbf{W}}$  and  $\frac{\partial J}{\partial \mathbf{c}}$ .

$$\begin{aligned} 1) \text{softmax}(\mathbf{x} + c)_i &= \frac{e^{x_i + c}}{\sum_{i'} e^{x_{i'} + c}} = \frac{e^{x_i} \cdot e^c}{\sum_{i'} e^{x_{i'}} \cdot e^c} = \frac{e^{x_i}}{\sum_{i'} e^{x_{i'}}} \\ &= \frac{e^{x_i}}{\sum_{i'} e^{x_{i'}}} = \text{softmax}(\mathbf{x})_i \\ \therefore \text{softmax}(\mathbf{x}) &= \text{softmax}(\mathbf{x} + c) \end{aligned}$$

$$\begin{aligned} 2) \quad \frac{\partial J}{\partial w} &= \frac{\partial J}{\partial z_i} \cdot \frac{\partial z}{\partial w} \\ \because J &= \sum_i \log \text{softmax}(\mathbf{z})_i \\ \therefore \frac{\partial J}{\partial z} &= \sum_i \frac{1}{\text{softmax}(\mathbf{z})_i} \cdot (\text{softmax}(\mathbf{z}))'_i \\ \because \text{softmax}(\mathbf{x}) &= \frac{e^{x_i}}{\sum_j e^{x_j}} \\ \frac{\partial J}{\partial z} &= \sum_i \frac{\log \text{softmax}(\mathbf{z})_i}{\partial z_i} \\ &= \sum_i \frac{1}{\text{softmax}(\mathbf{z})_i} \cdot \left( \frac{e^{z_i}}{\sum_j e^{z_j}} \right)' \\ &= \sum_i \frac{1}{\text{softmax}(\mathbf{z})_i} \cdot \frac{e^{z_i} \cdot z_i e^{z_j} - e^{z_i} \cdot e^{z_j}}{(\sum_j e^{z_j})^2} \\ &= \frac{1}{\text{softmax}(\mathbf{z})_i} \cdot \frac{e^{z_i}}{\sum_j e^{z_j}} \cdot \frac{\sum_j z_j e^{z_j} - e^{z_i}}{\sum_j e^{z_j}} \\ &\quad \downarrow \text{softmax}(\mathbf{z}) \\ &= \frac{\sum_j z_j e^{z_j} - e^{z_i}}{\sum_j e^{z_j}} \\ &= 1 - \text{softmax}(\mathbf{z})_i \end{aligned}$$

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial z_j} \cdot \frac{\partial z}{\partial w} = (1 - \text{softmax}(\mathbf{z})_i) x_i \quad i = (1 \dots N)$$

$$\frac{\partial J}{\partial C} = \frac{\partial J}{\partial z_i} = (1 - \text{softmax}(z)_i) \quad i = (1 \dots N)$$

## 0.2 Logistic Regression with Regularization [2 points]

- 1) [1 point] Let the data be  $(\mathbf{x}_i, y_i)_{i=1}^N$ , where  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \{0, 1\}$ . Logistic regression is a binary classification model, with the probability of  $y_i$  being 1 as:

$$p(y_i; \mathbf{x}_i, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^T \mathbf{x}_i) \triangleq \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}_i}},$$

where  $\boldsymbol{\theta}$  is the model parameter. Assume we impose an  $L_2$  regularization term on the parameter, defined as:

$$\mathcal{R}(\boldsymbol{\theta}) = \frac{\lambda}{2} \boldsymbol{\theta}^T \boldsymbol{\theta}$$

with a positive constant  $\lambda$ . Write out the final objective function for this logistic regression with regularization model.

- 2) [1 point] If we use gradient descent to solve the model parameter. Derive the updating rule for  $\boldsymbol{\theta}$ . Your answer should contain the derivation, not just the final answer.

1) Loss function of logistic regression.

$$J(\boldsymbol{\theta}) = \frac{1}{N} \left( -\log \prod_{i=1}^N p(y_i; \mathbf{x}_i, \boldsymbol{\theta}) + \frac{\lambda}{2} \boldsymbol{\theta}^T \boldsymbol{\theta} \right)$$

$$= \frac{1}{N} \left[ \sum_{i=1}^N \log p(y_i; \mathbf{x}_i, \boldsymbol{\theta}) + \frac{\lambda}{2} \boldsymbol{\theta}^T \boldsymbol{\theta} \right]$$

$$\because p(y_i; \mathbf{x}_i, \boldsymbol{\theta}) = \frac{e^{\boldsymbol{\theta}^T \mathbf{x}_i}}{1 + e^{\boldsymbol{\theta}^T \mathbf{x}_i}}$$

$$\therefore J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \log \left( \frac{e^{\boldsymbol{\theta}^T \mathbf{x}_i}}{1 + e^{\boldsymbol{\theta}^T \mathbf{x}_i}} \right) + \frac{\lambda}{2N} \boldsymbol{\theta}^T \boldsymbol{\theta}.$$

$$2) \frac{\partial J(\boldsymbol{\theta})}{\partial \theta_i} = \frac{1}{N} \left[ \sum_{i=1}^N (\sigma(\boldsymbol{\theta}^T \mathbf{x}_i) - y_i) \mathbf{x}_i \right] + \frac{\lambda}{N} \theta_i$$

$$\theta_j = \theta_j - \eta \left[ \frac{1}{N} \sum_{i=1}^N (\sigma(\boldsymbol{\theta}^T \mathbf{x}_i) - y_i) \mathbf{x}_j + \frac{\lambda}{N} \theta_j \right] \quad (j=0, 1, 2, \dots)$$

## 0.3 Derivative of the Softmax Function [3 points]

- 1) [1 point] Define the loss function as

$$J(\mathbf{z}) = - \sum_{k=1}^K y_k \log \tilde{y}_k,$$

where  $\tilde{y}_k = \frac{e^{z_k}}{\sum_{k'} e^{z_{k'}}}$ , and  $(y_1, \dots, y_K)$  is a known probability vector. Derive the  $\frac{\partial J(\mathbf{z})}{\partial \mathbf{z}}$ .

Note  $\mathbf{z} = (z_1, \dots, z_K)$  is a vector so  $\frac{\partial J(\mathbf{z})}{\partial \mathbf{z}}$  is in the form of a vector. Your answer should contain the derivation, not just the final answer.

- 2) [1 point] Assume the above softmax is the output layer of an FNN. Briefly explain how the derivative is used in the backpropagation algorithm.
- 3) [1 points] Let  $\mathbf{z} = \mathbf{W}^T \mathbf{h} + \mathbf{b}$ , where  $\mathbf{W}$  is a matrix,  $\mathbf{b}$  and  $\mathbf{h}$  are vectors. Use the chain rule to calculate the gradient of  $\mathbf{W}$  and  $\mathbf{b}$ , i.e.,  $\frac{\partial J}{\partial \mathbf{W}}$  and  $\frac{\partial J}{\partial \mathbf{b}}$ , respectively.

$$\begin{aligned}
 1) \quad \frac{\partial J(z)}{\partial z} &= -\sum_{k=1}^K y_k (\log \tilde{y}_k)' \\
 &= -\sum_{k=1}^K y_k \cdot \frac{1}{\tilde{y}_k} \cdot \tilde{y}_k' \\
 &= -\sum_{k=1}^K y_k \cdot \frac{1}{\tilde{y}_k} \cdot \frac{e^{z_k} e^{z_k'} - e^{z_k} e^{z_k}}{(\sum_{k'} e^{z_k'})^2} \\
 &= -y_k \cdot \frac{1}{\tilde{y}_k} \cdot \tilde{y}_k \cdot \frac{\sum_{k'} e^{z_k'} - e^{z_k}}{\sum_{k'} e^{z_k'}} \\
 &= -y_k \cdot \frac{\sum_{k'} e^{z_k'} - e^{z_k}}{\sum_{k'} e^{z_k'}} \\
 &= -y_k \cdot (1 - \tilde{y}_k) \\
 &= y_k \tilde{y}_k - y_k
 \end{aligned}$$

2) use the result of question (1),

$$\text{update the } \theta = \theta - \eta \frac{\partial J(z)}{\partial z} = \theta - \eta (y_k \tilde{y}_k - y_k)$$

$$\begin{aligned}
 3) \quad \frac{\partial J}{\partial w} &= \frac{\partial J}{\partial z_k} \cdot \frac{\partial z_k}{\partial w} \\
 &= (y_k \tilde{y}_k - y_k) h. \quad k=(1 \dots K)
 \end{aligned}$$

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial z_k} \cdot \frac{\partial z_k}{\partial b} = \frac{\partial J}{\partial z_k} \cdot 1 = y_k \tilde{y}_k - y_k \quad k=(1 \dots K)$$