

# Multi-class Logistic Regression

- Assume that there are  $s$  different classes.

Data set:  $\{(x_{n,1}, \dots, x_{n,k}; t_n)\}$ ,  $t_n = 1, 2, \dots, s$ ,  $n = 1, \dots, N$

- How to generalize logistic regression?

- $t_n = 1 \Leftrightarrow [1, 0, \dots, 0]$ ,  $t_n = 2 \Leftrightarrow [0, 1, 0, \dots, 0]$ ,  $\dots$ ,  
 $t_n = s \Leftrightarrow [0, \dots, 0, 1]$

- Obtain an  $N \times s$  matrix  $\mathbf{t} = [t_{n,m}]$  such that

$$t_{n,m} = \begin{cases} 1 & \text{if } t_n = m, \\ 0 & \text{if } t_n \neq m. \end{cases}$$

- Define  $\sigma : \mathbb{R}^s \rightarrow (0, 1)^s$  by

$$\sigma(\mathbf{a}) = \left( \frac{e^{a_1}}{\sum_{i=1}^s e^{a_i}}, \dots, \frac{e^{a_s}}{\sum_{i=1}^s e^{a_i}} \right),$$

where  $\mathbf{a} = (a_1, a_2, \dots, a_s)$ .

The function  $\sigma$  is called the **softmax** function.

- When  $s = 2$ , we have

$$\sigma(\mathbf{a}) = \left( \frac{e^{a_1}}{e^{a_1} + e^{a_2}}, \frac{e^{a_2}}{e^{a_1} + e^{a_2}} \right) = \left( \frac{1}{1 + e^{a_2 - a_1}}, \frac{e^{a_2 - a_1}}{1 + e^{a_2 - a_1}} \right).$$

- The softmax function is a generalization of the sigmoid function.

- Define  $\mathbf{y} = [y_1, \dots, y_s] = \boldsymbol{\sigma}(\mathbf{a})$ .
- For  $m, j = 1, \dots, s$ ,

$$\frac{\partial y_m}{\partial a_j} = y_m(\delta_{j,m} - y_j),$$

where  $\delta_{j,m}$  is the Kronecker's delta, i.e.

$$\delta_{j,m} = \begin{cases} 1 & \text{if } j = m, \\ 0 & \text{otherwise.} \end{cases}$$

- Consider a  $(k + 1) \times s$  matrix  $\mathbf{w} = [w_{p,q}]$ .

Define  $\mathbf{y} = \sigma(X\mathbf{w}) = [y_{n,m}]$ ,

where  $X$  is as before and  $\sigma$  is applied to the rows of  $X\mathbf{w}$ .

Each row of  $\mathbf{y}$  consists of probabilities for classes 1 through  $m$ .

- The likelihood function is given by

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N \prod_{m=1}^s y_{n,m}^{t_{n,m}}.$$

- The cross-entropy is

$$E(\mathbf{w}) = - \sum_{n=1}^N \sum_{m=1}^s t_{n,m} \ln y_{n,m}.$$

- One can check

$$\nabla E(\mathbf{w}) = \left[ \frac{\partial E}{\partial \mathbf{w}_{p,q}} \right] = \left[ \sum_{n=1}^N (y_{n,q} - t_{n,q}) x_{n,p} \right] = \mathbf{X}^T (\mathbf{y} - \mathbf{t}).$$

(Use

$$\sum_{m=1}^S t_{n,m} (\delta_{q,m} - y_{n,q}) = \sum_{m=1}^S (t_{n,m} \delta_{q,m} - t_{n,m} y_{n,q}) = t_{n,q} - y_{n,q}.)$$

- Gradient Descent

$$\mathbf{w}_{i+1} = \mathbf{w}_i - \eta \mathbf{X}^T (\mathbf{y} - \mathbf{t}).$$

- Let  $\mathbf{w}_i \rightarrow \mathbf{w}_*$  as  $i \rightarrow \infty$ .

Given  $\mathbf{x} = [x_1, \dots, x_k, 1]$ , the coordinates of the vector

$$\mathbf{y} = \sigma(\mathbf{x}\mathbf{w}_*)$$

represent the probabilities for the classes.

- The (multi-class) logistic regression is the simplest **neural network**.