



# Linear Classifiers

DCS310

Sun Yat-sen University

# Outline

---

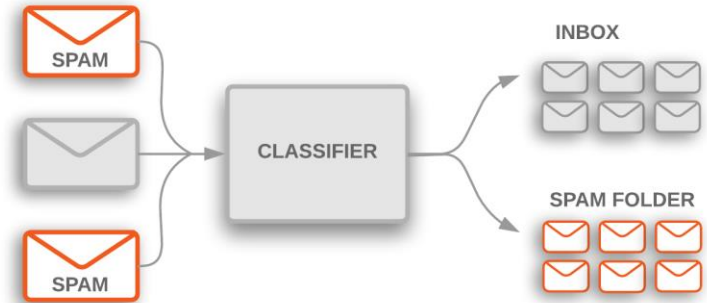
- Two-class Case
- Multi-class Case

# Examples

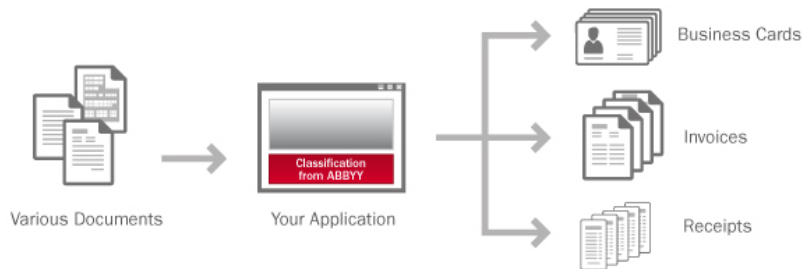
- Image category classification



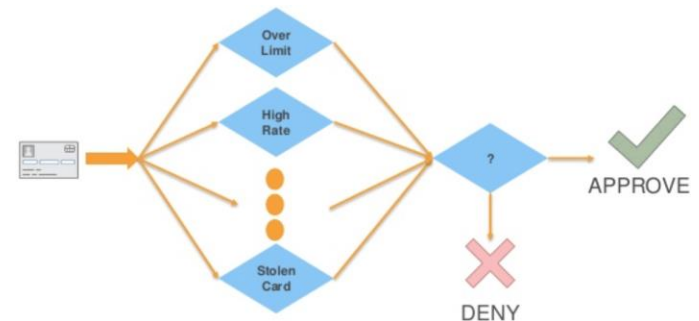
- Spam e-mails detection



- Document automatic categorization

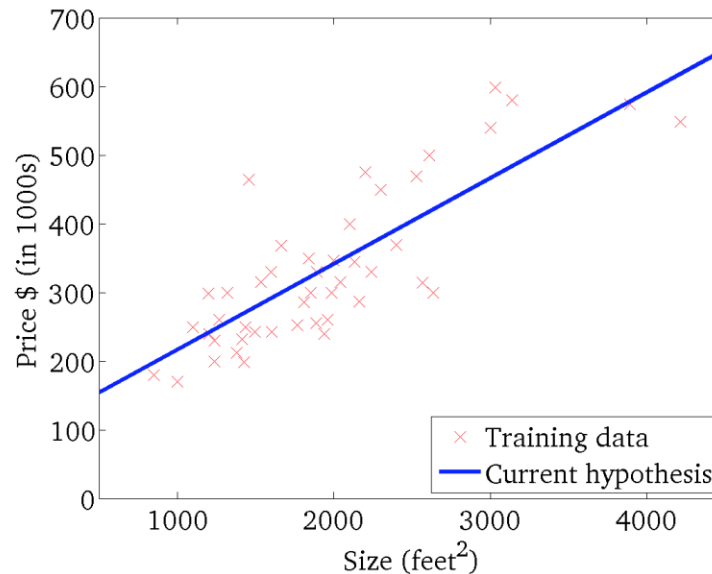


- Transaction fraud detection



# Logistic Regression

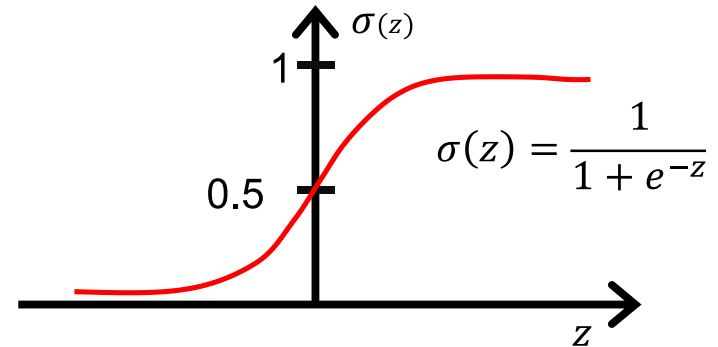
- In classification, the target variable  $y \in \{0, 1\}$
- In linear regression, the output  $f(x) = xw$  falls in the range  $[-\infty, +\infty]$



- The output value of linear regression is **not compatible** with the target values in the classification tasks

- Sigmoid/logistic function

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



- Logistic regression

$$f(\mathbf{x}) = \sigma(\mathbf{x}\mathbf{w})$$

- Linear regression

$$f(\mathbf{x}) = \mathbf{x}\mathbf{w}$$

- The output range is transformed from  $[-\infty, +\infty]$  to  $[0, 1]$

# Cost Function

---

- Goal
  - If the true label  $y = 1$ , we want  $f(\mathbf{x}) = \sigma(\mathbf{x}\mathbf{w})$  to be close to 1
  - If the true label  $y = 0$ , we want  $f(\mathbf{x}) = \sigma(\mathbf{x}\mathbf{w})$  to be close to 0
- To achieve this goal, we can define a cost function similar to that in regression

$$L(\mathbf{w}) = (\sigma(\mathbf{x}\mathbf{w}) - y)^2$$

- **Alternatively**, we can also seek to minimize

$$-\log(\sigma(\mathbf{x}\mathbf{w})) \text{ if } y = 1 \quad \text{or} \quad -\log(1 - \sigma(\mathbf{x}\mathbf{w})) \text{ if } y = 0$$

- The objective above can be equivalently written as

$$L(\mathbf{w}) = -y \log(\sigma(\mathbf{x}\mathbf{w})) - (1 - y) \log(1 - \sigma(\mathbf{x}\mathbf{w}))$$

*If  $y = 1$ ,  $L(w)$  reduces to  $L(w) = -\log(\sigma(xw))$ ;*

*If  $y = 0$ ,  $L(w)$  reduces to  $L(w) = -\log(1 - \sigma(xw))$*

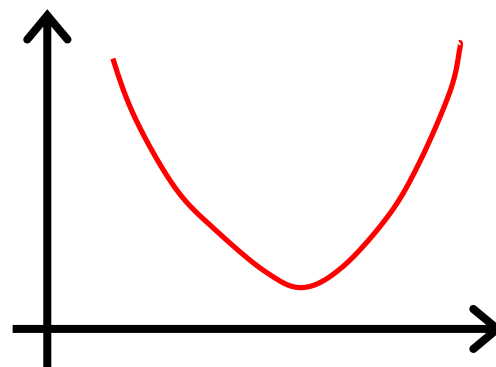
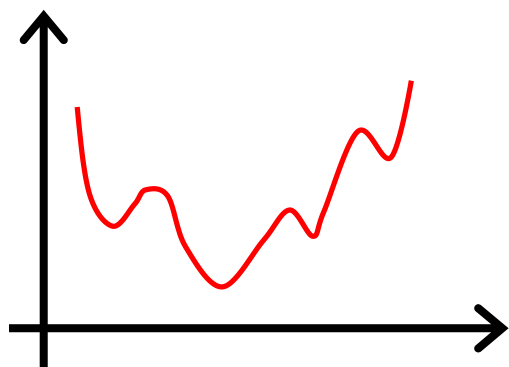
- The loss above is called the ***cross-entropy loss***

- Which cost function is better?

**Squared error loss:**  $L(\mathbf{w}) = (\sigma(\mathbf{xw}) - y)^2$

**Cross entropy:**  $L(\mathbf{w}) = -[y \log(\sigma(\mathbf{xw})) + (1 - y) \log(1 - \sigma(\mathbf{xw}))]$

- Squared loss is non-convex
- Cross entropy is convex



Convex function is easier to optimize

- In next lecture, another advantage of using cross-entropy loss will be manifested from the perspective of *more accurate modeling*



The Thirty-Third AAAI Conference on Artificial Intelligence (AAAI-19)

### DeepCF: A Unified Framework of Representation Learning and Matching Function Learning in Recommender System

Zhi-Hong Deng,<sup>1</sup> Ling Huang,<sup>1</sup> Chang-Dong Wang,<sup>1</sup> Jian-Huang Lai,<sup>1</sup> Philip S. Yu<sup>2,3</sup><sup>1</sup>School of Data and Computer Science, Sun Yat-sen University, Guangzhou, China<sup>2</sup>Department of Computer Science, University of Illinois at Chicago, Chicago, USA<sup>3</sup>Institute for Data Science, Tsinghua University, Beijing, China  
dengzh7@mail2.sysu.edu.cn, huangling1@hotmail.com, changdongwang@hotmail.com  
stsljh@mail.sysu.edu.cn, psyu@uic.edu

mentioned in **Problem Statement**, we assume  $y_{ui}$  obeys a Bernoulli distribution, i.e.,  $y_{ui} \sim \text{Bern}(p_{ui})$ . By replacing  $p_{ui}$  with  $\hat{y}_{ui}$  in Equation 2, we can define the likelihood function as

$$\begin{aligned} L(\Theta) &= \prod_{(u,i) \in \mathcal{Y}^+ \cup \mathcal{Y}^-} P(y_{ui} | \Theta) \\ &= \prod_{(u,i) \in \mathcal{Y}^+ \cup \mathcal{Y}^-} \hat{y}_{ui}^{y_{ui}} (1 - \hat{y}_{ui})^{1-y_{ui}}, \end{aligned} \quad (3)$$

where  $\mathcal{Y}^+$  denotes all the observed interactions in  $\mathbf{Y}$  and  $\mathcal{Y}^-$  denotes the sampled unobserved interactions, i.e., the negative instances. Furthermore, taking the negative logarithm of the likelihood (NLL), we obtain

$$\begin{aligned} \ell_{BCE} &= - \sum_{(u,i) \in \mathcal{Y}^+ \cup \mathcal{Y}^-} y_{ui} \log \hat{y}_{ui} \\ &\quad + (1 - y_{ui}) \log(1 - \hat{y}_{ui}). \end{aligned} \quad (4)$$

Based on all the above assumptions and formulations, we finally obtain an objective function which is suitable for learning from implicit feedback data, i.e., the binary cross-entropy loss function.

## Cross-Domain Explicit-Implicit-Mixed Collaborative Filtering Neural Network

Chang-Dong Wang<sup>✉</sup>, Member, IEEE, Yan-Hui Chen, Wu-Dong Xi<sup>✉</sup>, Ling Huang<sup>✉</sup>, Member, IEEE, and Guangqiang Xie<sup>✉</sup>, Member, IEEE

### B. Loss Function

To predict the interaction probability  $\hat{y}_{uid}^d \forall d = 1, 2$ , an appropriate loss function needs to be designed. The recommender system focusing on the explicit ratings usually formulates the recommendation problem as a rating prediction problem. Similarly, for the recommendation problem with the implicit interactions, it is often formulated as an interaction prediction problem. For any domain  $\mathbb{D}^d \forall d = 1, 2$ , assume that the implicit interaction  $y_{uid}^d$  obeys a Bernoulli distribution, the likelihood function in CEICFNet can be defined as follows:

$$\begin{aligned} L^1(\Theta) &= \prod_{(u,i^1) \in \mathcal{Y}_+^1 \cup \mathcal{Y}_-^1} (\hat{y}_{ui^1}^1)^{y_{ui^1}^1} (1 - \hat{y}_{ui^1}^1)^{1-y_{ui^1}^1} \\ L^2(\Theta) &= \prod_{(u,i^2) \in \mathcal{Y}_+^2 \cup \mathcal{Y}_-^2} (\hat{y}_{ui^2}^2)^{y_{ui^2}^2} (1 - \hat{y}_{ui^2}^2)^{1-y_{ui^2}^2} \end{aligned} \quad (2)$$

where  $\mathcal{Y}_+^d$  and  $\mathcal{Y}_-^d$ , respectively, denote the observed interactions and the randomly selected negative samples in  $\mathbb{D}^d \forall d = 1, 2$ .

Finally, by taking the negative logarithm of the likelihood (NLL), we obtain the loss function as follows:

$$\begin{aligned} l^1(\Theta) &= - \sum_{(u,i^1) \in \mathcal{Y}_+^1 \cup \mathcal{Y}_-^1} y_{ui^1}^1 \log \hat{y}_{ui^1}^1 + (1 - y_{ui^1}^1) \log(1 - \hat{y}_{ui^1}^1) \\ l^2(\Theta) &= - \sum_{(u,i^2) \in \mathcal{Y}_+^2 \cup \mathcal{Y}_-^2} y_{ui^2}^2 \log \hat{y}_{ui^2}^2 + (1 - y_{ui^2}^2) \log(1 - \hat{y}_{ui^2}^2) \\ l(\Theta) &= l^1(\Theta) + \lambda l^2(\Theta) \end{aligned} \quad (3)$$

where  $\lambda$  is the tradeoff parameter that is utilized to tune the importance of the two domains.

# Gradient Descent

---

- The gradient of the cross-entropy loss

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{n} \sum_{i=1}^n [\sigma(\mathbf{x}^{(i)} \mathbf{w}) - y^{(i)}] \mathbf{x}^{(i)T}$$

- The optimal  $\mathbf{w}^*$  can be obtained by solving  $\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = 0$ . But here, *the analytical solution does not exist*
- Thus, we can only resort to the numerical methods
  - Gradient descent
  - Newton methods
  - Coordinate descent
  - .....

- Since the cross-entropy loss is convex, the gradient descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - r \cdot \frac{\partial L(\mathbf{w})}{\partial \mathbf{w}}$$

is guaranteed to converge to the optimal value  $\mathbf{w}^*$

- By examining the gradient

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{n} \sum_{i=1}^n \left[ \underbrace{\sigma(\mathbf{x}^{(i)} \mathbf{w}) - y^{(i)}}_{\text{prediction error}} \right] \mathbf{x}^{(i)T}$$

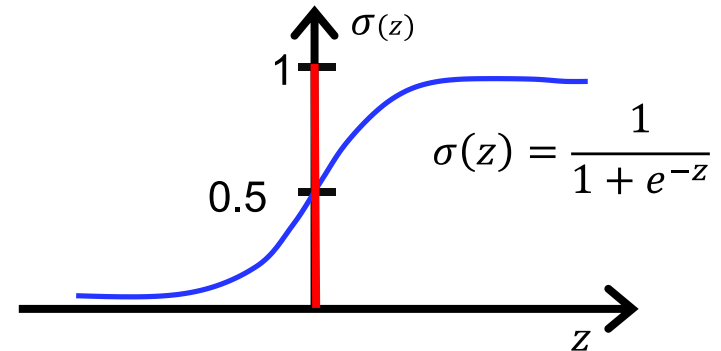
we see that the GD always seeks to reduce the prediction error

- If  $y^{(i)} = 1$ , the algorithm derives  $\sigma(\mathbf{x}^{(i)} \mathbf{w})$  towards 1
- If  $y^{(i)} = 0$ , the algorithm derives  $\sigma(\mathbf{x}^{(i)} \mathbf{w})$  towards 0

# Decision Boundary

- The sample is classified into 1 and 0 as

$$\hat{y} = \begin{cases} 1, & \text{if } \sigma(\mathbf{x}\mathbf{w}) \geq 0.5 \\ 0, & \text{if } \sigma(\mathbf{x}\mathbf{w}) < 0.5 \end{cases}$$

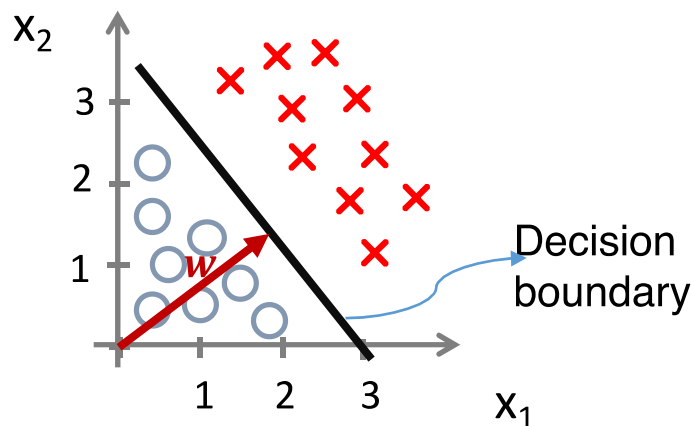


- This is equivalent to classify the samples as

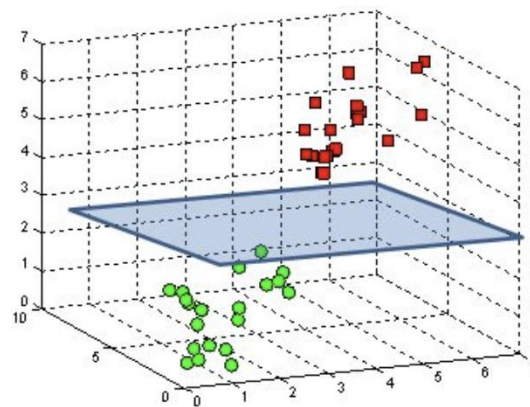
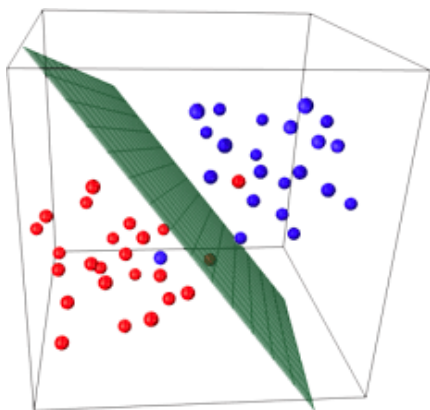
$$\hat{y} = \begin{cases} 1, & \text{if } \mathbf{x}\mathbf{w} \geq 0 \\ 0, & \text{if } \mathbf{x}\mathbf{w} < 0 \end{cases}$$

- The decision boundary consists of  *$\mathbf{x}$  that satisfy  $\mathbf{x}\mathbf{w} = 0$*

- Since  $w$  is a vector, all  $x$  that satisfies  $xw = 0$  constitute *a space that is orthogonal to  $w$*



- In the *two*-dimensional case, the space is a straight *line*
- In the *three*-dimensional case, the space is a *plane*

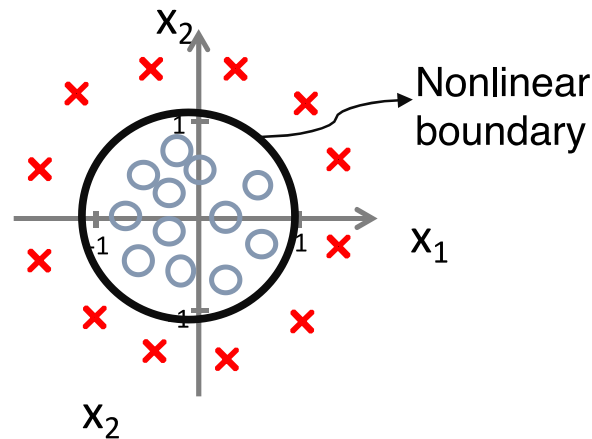


- For a fixed vector  $\mathbf{w} \in \mathbb{R}^{d \times 1}$ , the set of points

$$\mathbf{x} \in \{\mathbf{x} | \mathbf{x}\mathbf{w} = 0\}$$

constitute a  $(d - 1)$ -dimensional *hyper-plane*

- The hyper-planes can *never represent a nonlinear decision boundary*, e.g.,



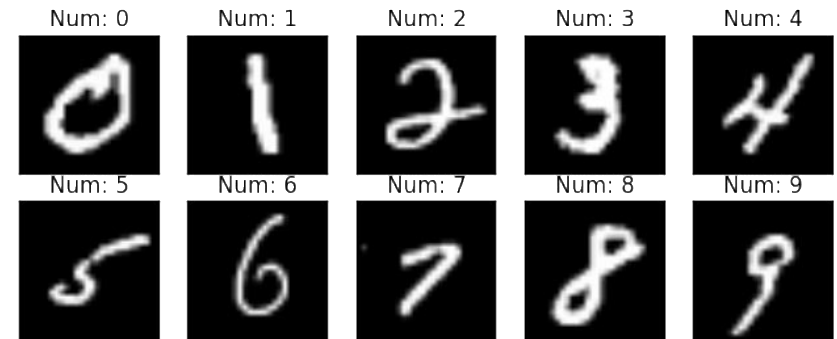
- That's why we call the logistic regression a linear classifier

# Outline

---

- Two-class Case
- Multi-class Case

- Many applications have more than 2 classes



- Two methods to deal with multiclass classification

- One-vs-All

Transform the multi-class problem into multiple binary problem

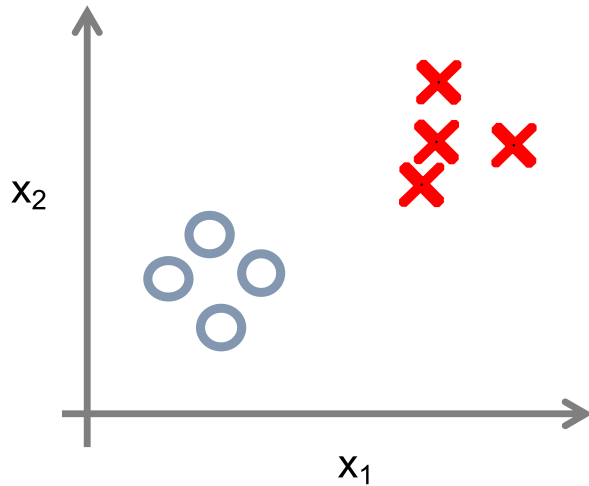
- Softmax function

Classifying the sample into one of the classes directly

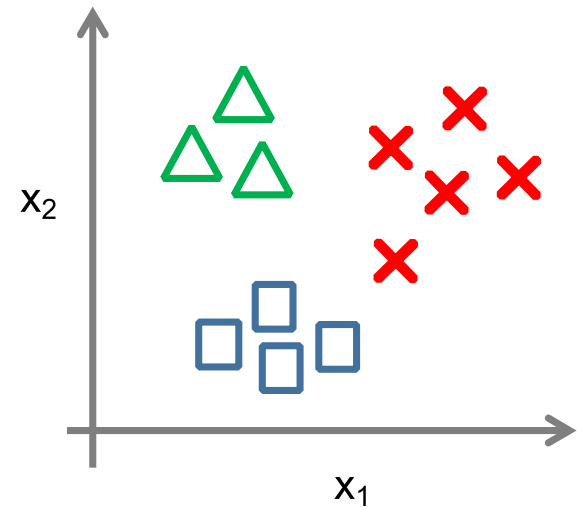


# One-vs-All

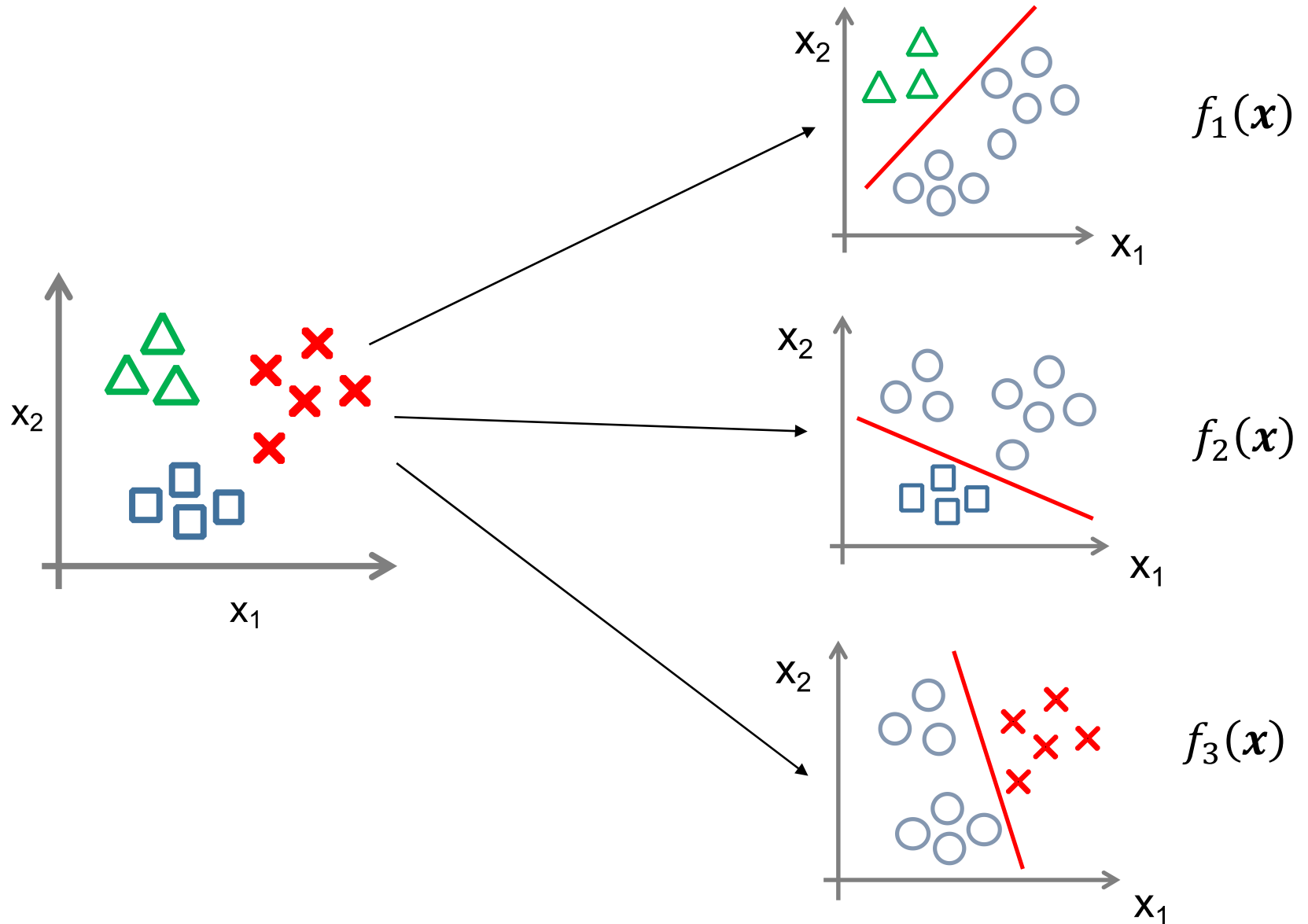
Binary classification



Multiclass classification

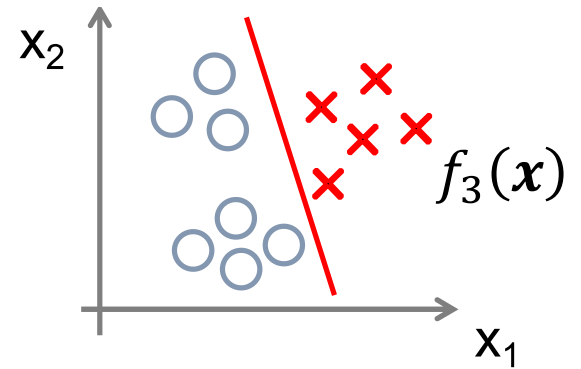
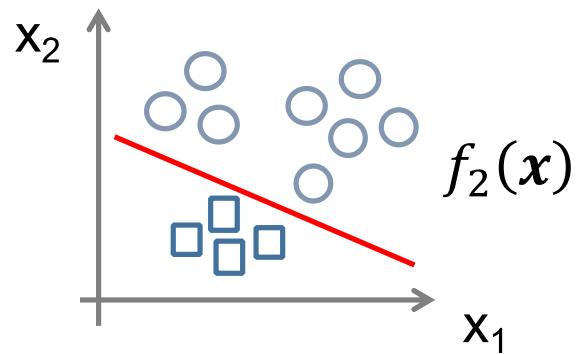
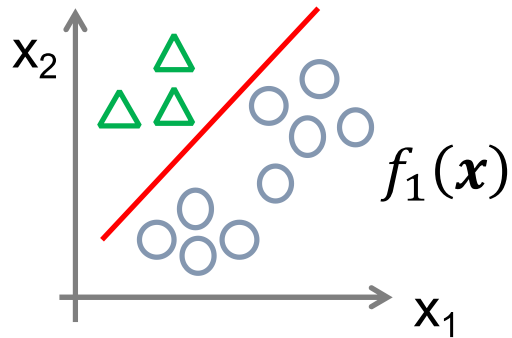


- Train a classifier for each class



- To predict the class for a new sample  $\mathbf{x}$ , pick the class such that

$$k = \arg \max_i f_i(\mathbf{x})$$



# Softmax Function

---

- Softmax function

$$\text{softmax}_i(\mathbf{z}) = \frac{e^{z_i}}{\sum_{k=1}^K e^{z_k}}$$

It can be seen that  $\sum_{i=1}^K \text{softmax}_i(\mathbf{z}) = 1$

- The probability that a data  $\mathbf{x} \in \mathbb{R}^{1 \times d}$  is classified to the  $i$ -th class is

$$f_i(\mathbf{x}) = \text{softmax}_i(\mathbf{x}\mathbf{W}) = \frac{e^{\mathbf{x}\mathbf{w}_i}}{\sum_{k=1}^K e^{\mathbf{x}\mathbf{w}_k}}$$

where  $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K] \in \mathbb{R}^{d \times K}$

- If  $\mathbf{x}$  belongs to the  $i$ -th class, the model should encourage  $f_i(\mathbf{x})$  as large as possible

- The softmax function with  $K = 2$  is equivalent to logistic function

➤ Under the two-class case, we have

$$\begin{aligned} \text{softmax}_1(\mathbf{x}W) &= \frac{e^{xw_1}}{e^{xw_1} + e^{xw_2}} & \text{softmax}_2(\mathbf{x}W) &= \frac{e^{xw_2}}{e^{xw_1} + e^{xw_2}} \\ &= \frac{1}{1 + e^{-x(w_1 - w_2)}} & &= \frac{e^{-x(w_1 - w_2)}}{1 + e^{-x(w_1 - w_2)}} \end{aligned}$$

➤ It can be seen that

$$\text{softmax}_1(\mathbf{x}W) = \sigma(x(w_1 - w_2))$$

$$\text{softmax}_2(\mathbf{x}W) = 1 - \sigma(x(w_1 - w_2))$$

The two-class softmax classification is equivalent to the logistic regression, with the parameter being  $w_1 - w_2$

# Cost Function

- For a training dataset with  $K$  classes, its label  $\mathbf{y}$  is represented by a **one-hot vector**, which is illustrated as follows

$$\begin{aligned} &[1, 0, 0, \dots, 0], \\ &[0, 1, 0, \dots, 0], \\ &\vdots \\ &[0, 0, 0, \dots, 1] \end{aligned}$$

- The objective is to maximize the corresponding probability  $f_i(\mathbf{x})$ . Thus, the cost function can be written as

$$L(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K) = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K y_k^{(i)} \log \text{softmax}_k(\mathbf{x}^{(i)} \mathbf{W})$$

- $y_k^{(i)}$  is the  $k$ -th element of  $\mathbf{y}^{(i)} \in \mathbb{R}^{1 \times K}$  **Cross-entropy loss**

# Gradient Descent

- The gradient *w.r.t.*  $\mathbf{w}_j \in \mathbb{R}^{d \times 1}$  is

$$\frac{\partial L(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K)}{\partial \mathbf{w}_j} = \frac{1}{n} \sum_{i=1}^n \left( \text{softmax}_j(\mathbf{x}^{(i)} \mathbf{W}) - y_j^{(i)} \right) \mathbf{x}^{(i)T}$$

Note that all  $\mathbf{w}_j$  for  $j = 1, \dots, K$  should be updated simultaneously

- By representing  $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K] \in \mathbb{R}^{d \times K}$ , we have

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)T} \left( \text{softmax}(\mathbf{x}^{(i)} \mathbf{W}) - \mathbf{y}^{(i)} \right)$$

- $\text{softmax}(\mathbf{x}^{(i)} \mathbf{W}) = [\text{softmax}_1(\mathbf{x}^{(i)} \mathbf{W}), \dots, \text{softmax}_K(\mathbf{x}^{(i)} \mathbf{W})] \in \mathbb{R}^{1 \times K}$   
*is a row vector*
- $\mathbf{x}^{(i)T} \in \mathbb{R}^{d \times 1}$  *is a column vector*

- Updating:  $\mathbf{W}_{t+1} = \mathbf{W}_t - r \cdot \left. \frac{\partial L(\mathbf{W})}{\partial \mathbf{W}} \right|_{\mathbf{W}=\mathbf{W}_t}$