

# 08-Linear Classification (Discriminative)

## Probabilistic assumption

- We would like to impose probabilistic assumptions to classification tasks, so that we could have a principled way of learning, e.g., MLE and MAP estimation.
- **Bernoulli distribution**

A Bernoulli distribution, parametrized by  $p$ , is a two-point distribution

$$(1 - p) < 0 > + p < 1 >$$

That is,

$$\Pr\{t = 1\} = p \quad \text{and} \quad \Pr\{t = 0\} = 1 - p$$

The probability of  $t$  is  $p^t \cdot (1 - p)^{1-t}$

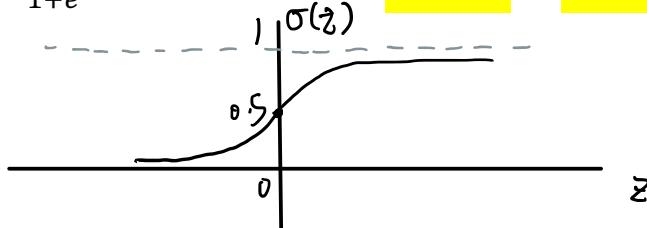
- If  $t = 1$ , the probability is  $p$ .
- If  $t = 0$ , the probability is  $1 - p$ .

- For a classification task, we can always assume  $t|x \sim \text{Bernoulli}(p)$
- And, we further assume  
 $t|x \sim \text{Bernoulli}(f_{\theta}(x))$   
for some unknown parameters  $\theta$ .
- In linear regression,  $f_{\theta}(x) = w^T x + b \in \mathbb{R}$ , which does not work for  $p \in [0,1]$ .
- For classification, we may squash the linear function into  $(0,1)$  as

$$f_{\theta}(x) = \sigma(w^T x + b)$$

where

$\sigma(z) = \frac{1}{1+e^{-z}}$  is known as the sigmoid or logistic function.



The resulting model is known as logistic regression

- Why sigmoid function?
  - Exponential family
  - Other squashing functions are also acceptable

- E.g., the probit function  $\Phi(x) = \int_{-\infty}^x N(t; 0, 1) dt$

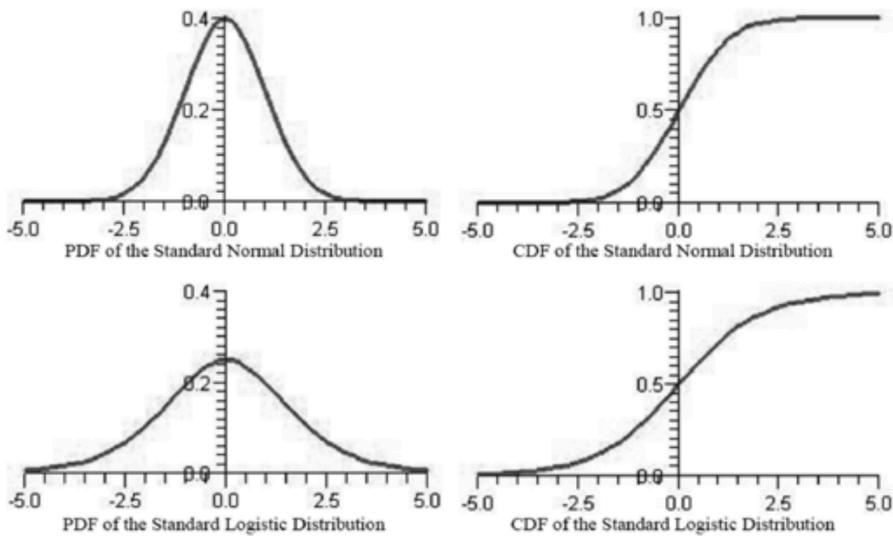


Figure 1. The Standard Normal and Standard Logistic Probability Distributions

Source : Park (2010)

Note: The difference between a probit function or a sigmoid function is not merely scaling.

### Maximum Likelihood Estimation (MLE) for Logistic Regression

Suppose we have a dataset  $\{(x^{(m)}, t^{(m)})\}_{m=1}^M$ , and the above assumption says

$$t^{(m)} \sim \text{Bernoulli}(f_{w,b}(x))$$

where  $w, b$  are true, but unknown, parameters.

The log-likelihood is

$$\begin{aligned}
 & \log p(t^{(1)}, \dots, t^{(M)} | x^{(1)}, \dots, x^{(M)}) \\
 &= \log \prod_{m=1}^M p(t^{(m)} | x^{(m)}) \quad [\text{data iid, and } t^{(m)} \text{ only depends on } x^{(m)}] \\
 &= \sum_{m=1}^M \log p(t^{(m)} | x^{(m)}) \quad [\log \prod \Leftrightarrow \sum \log] \\
 &= \sum_{m=1}^M \log [f_{w,b}(x)]^{t^{(m)}} \cdot [1 - f_{w,b}(x)]^{1-t^{(m)}} \quad [\text{Bernoulli}] \\
 &= \sum_{m=1}^M \left[ t^{(m)} \log f_{w,b}(x) + (1-t^{(m)}) \cdot \log (1 - f_{w,b}(x)) \right]
 \end{aligned}$$

The training objective is to minimize

$$J(w, b) = \frac{1}{M} \sum_{m=1}^M \left[ -t^{(m)} \log f_{w,b}(x) - (1-t^{(m)}) \log (1-f_{w,b}(x)) \right]$$

$t \in \{0,1\}$  can be thought of as a distribution  $t$ .

$$\text{If } t = 0, \text{ then } t = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad \text{If } t = 1, \text{ then } t = \begin{pmatrix} 0 \\ 1 \end{pmatrix}.$$

$t$  is the index representation, whereas  $t$  is the one-hot representation  
They are equivalent.

It is noted that  $\begin{pmatrix} 1-f_{w,b}(x) \\ f_{w,b}(x) \end{pmatrix}$  is also a distribution

The above objective is the cross-entropy loss of  $f_{w,b}(x)$  against  $t$ .

The Kullback–Leibler (KL) divergence is an important measure of the dissimilarity between two distributions. The KL of two distributions  $p$  and  $q$ , in the discrete case, is defined as

$$KL(p \parallel q) = \sum_{i=1}^k p_i \log \frac{p_i}{q_i}$$

[HW] Show that minimizing the cross-entropy loss is equivalent to minimizing the KL divergence between the predicted distribution and the one-hot target distribution. Either  $KL(f \parallel t)$  or  $KL(t \parallel f)$ , but which and why?

## Gradient for cross-entropy loss

We introduce an intermediate variable

$$z^{(m)} = w^T x^{(m)} + b$$

$$y^{(m)} = f(z^{(m)}) = \frac{1}{1+e^{-z^{(m)}}}$$

We don't need the hat for  $y$  because we denote groundtruth label as  $t$ .

$$J^{(m)}(w, b) = [-t^{(m)} \log y^{(m)} - (1-t^{(m)}) \log (1-y^{(m)})]$$

$$\frac{\partial}{\partial z^{(m)}} J^{(m)}(w, b) = \frac{\partial}{\partial z^{(m)}} \left[ t^{(m)} \log (1+e^{-z^{(m)}}) - (1-t^{(m)}) \log \left( \frac{e^{-z^{(m)}}}{1+e^{-z^{(m)}}} \right) \right]$$

$$= t^{(m)} \underbrace{\frac{1}{1+e^{-z^{(m)}}} \cdot e^{-z^{(m)}} \cdot (-1)}_{1-y^{(m)}} - (1-t^{(m)}) \cdot (-1)$$

$$+ (1-t^{(m)}) \cdot \underbrace{\frac{1}{1+e^{-z^{(m)}}} \cdot e^{-z^{(m)}} \cdot (-1)}_{1-y^{(m)}}$$

$$= y^{(m)} - t^{(m)}$$

$$\frac{\partial J^{(m)}}{\partial w_i} = \frac{\partial J^{(m)}}{\partial z^{(m)}} \cdot \frac{\partial z^{(m)}}{\partial w_i} = (y^{(m)} - t^{(m)}) \cdot x_i$$

$$\frac{\partial J}{\partial b_i} = y^{(m)} - t^{(m)}$$

Since  $J = \sum_{m=1}^M J^{(m)}$ ,

$$\frac{\partial J}{\partial w_i} = \sum_{m=1}^M (y^{(m)} - t^{(m)}) x_i^{(m)} \quad \frac{\partial J}{\partial w} = \mathbf{x}^\top (\mathbf{y} - \mathbf{t})$$

$$\frac{\partial J}{\partial b} = \sum_{m=1}^M (y^{(m)} - t^{(m)})$$

## Multi-class classification

If  $t \in \{1, 2, \dots, K\}$ , it is safe to assume  $t|x \sim \text{Multinomial}(\mathbf{p})$ , where  $\mathbf{p}$  is a probability of a **multinomial distribution** with  $p_k = \Pr\{t = k|x\}$

### Note:

- A Bernoulli trial  $x \sim \text{Bernoulli}(p)$  is a two-point distribution  $(1-p)\langle 0 \rangle + p\langle 1 \rangle$
- A **binomial distribution**  $n \sim \text{Binomial}(p, N)$  is the sum of Bernoulli( $p$ ) repeated  $N$  times, i.e.,  

$$n = \sum_{i=1}^N x_i, \text{ where } x_i \sim \text{Bernoulli}(p)$$
- In a general setting, a multinomial distribution is parametrized by  $\text{Multinomial}(\mathbf{p}, N)$ , and  

$$n = \sum_{i=1}^N x_i, \text{ where } \Pr\{x_i = k\} = p_k$$

In classification problems, we usually have  $N = 1$ , which is omitted for simplicity. It is also called a **multinoulli distribution**, or a **categorical distribution**.

For classification problems, it suffices to model uncertainty by a multinomial distribution. However, we still need to design the machine learning model that computes  $\mathbf{p}$  from  $\mathbf{x}$ .

## Linear multi-class classification

We still hope to have a linear model like  $\mathbf{W}\mathbf{x} + \mathbf{b}$  and then normalize it as a probability. We make things work in a heuristic manner as follows.

- We define a scoring function  $z_i = \mathbf{w}_i^\top \mathbf{x} + b_i$   
i.e.,  $\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b}$ , where  $\mathbf{x} \in \mathbb{R}^d$  (number of input features)  
 $\mathbf{z} \in \mathbb{R}^K$ ,  $\mathbf{b} \in \mathbb{R}^K$  (number of categories)  
 $\mathbf{W} \in \mathbb{R}^{K \times d}$
- We make the scoring function positive  

$$\tilde{y}_i = \exp(z_i)$$

- We make the scoring function normalized

$$y_i = \frac{\tilde{y}_i}{\sum_j \tilde{y}_j} = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

Such computation is known as ( $k$ -way) softmax:

$$\mathbf{y} = \text{softmax}(\mathbf{z})$$

The terminology is said in comparison to (hard) argmax, which gives a one-hot index of the maximum element in the vector  $\mathbf{z}$ . Instead, softmax gives a soft distribution over all categories, and the maximum entry in  $\mathbf{z}$  is also the one in  $\mathbf{y}$ .

## Maximum likelihood estimation for softmax multi-class classification

$$\begin{aligned}
 & \log p(t^{(1)}, \dots, t^{(M)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}) \\
 &= \log \prod_{m=1}^M p(t^{(m)} | \mathbf{x}^{(m)}) \quad [\text{Data iid, } t^{(m)} \text{ depends only on } \mathbf{x}^{(m)}] \\
 &= \sum_{m=1}^M \log p(t^{(m)} | \mathbf{x}^{(m)}) \quad [\log \prod \Leftrightarrow \sum \log] \\
 &= \sum_{m=1}^M \log y_{t^{(m)}}^{(m)} \quad [\text{Index representation of } t^{(m)}] \\
 &= \sum_{m=1}^M \sum_{k=1}^K t_k^{(m)} \log y_k^{(m)} \quad [\text{One-hot representation } t^{(m)}]
 \end{aligned}$$

## Logistic regression vs. softmax

- Logistic regression is equivalent to a two-way softmax

For a two-way softmax with  $t \in \{0, 1\}$ .

$$y_0 = \hat{P}_r[t=0 | \mathbf{x}] = \frac{\exp\{-(\mathbf{w}_0^T \mathbf{x} + b_0)\}}{\exp\{-(\mathbf{w}_0^T \mathbf{x} + b_0)\} + \exp\{-(\mathbf{w}_1^T \mathbf{x} + b_1)\}}$$

$$\begin{aligned}
 y_1 &= \hat{P}_r[t=1 | \mathbf{x}] = \frac{\exp\{-(\mathbf{w}_1^T \mathbf{x} + b_1)\}}{\exp\{-(\mathbf{w}_0^T \mathbf{x} + b_0)\} + \exp\{-(\mathbf{w}_1^T \mathbf{x} + b_1)\}} \\
 &= \frac{1}{\exp\{-(\underbrace{(\mathbf{w}_0 - \mathbf{w}_1)^T \mathbf{x}}_{\mathbf{w}_0} + \underbrace{(b_0 - b_1))\})} + 1 \\
 &= \sigma(\mathbf{w}_0^T \mathbf{x} + b_0)
 \end{aligned}$$

- Most results for logistic regression still holds for softmax

[HW: derive the gradient of the cross-entropy loss for softmax]

Note that MLE for logistic regression is also equivalent to MLE for 2-

way softmax.

Also, MAP estimation/regularization can also be applied as in linear regression

### Inference criteria (decision rule)

For a  $k$ -category classification problem, the hypothesis function is

$$h : \mathcal{X} \rightarrow \{1, \dots, k\}$$

Suppose we have a probabilistic modeling of  $\Pr\{t = k | \mathbf{x}\}$ . Then,

**Max a posteriori (MAP) inference**  $\Leftrightarrow$  Maximizing accuracy

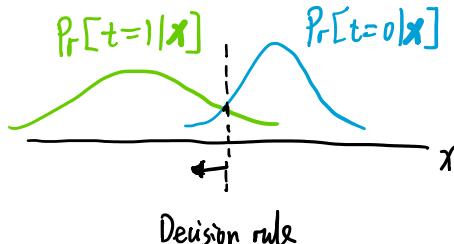
$$h(\mathbf{x}) = \operatorname{argmax}_{t=1}^k \Pr[t=k | \mathbf{x}] \quad \underset{t}{\operatorname{maximize}} \quad \mathbb{E}[1_{\{t=h(\mathbf{x})\}}]$$

Proof:  $\mathbb{E}[1_{\{t=h(\mathbf{x})\}}]$

$$= \sum_{t=1}^k p(t | \mathbf{x}) 1_{\{t=h(\mathbf{x})\}}$$

$$= p(h(\mathbf{x}) | \mathbf{x})$$

$$\underset{h(\mathbf{x})}{\operatorname{maximize}} \quad p(h(\mathbf{x}) | \mathbf{x}) \Leftrightarrow h(\mathbf{x}) = \operatorname{argmax}_{t=1}^k p(t | \mathbf{x})$$



This explains:

- For logistic regression, we predict
- For softmax, we predict

$$\hat{t} = \begin{cases} 1 & \text{if } y \geq 0.5 \\ 0 & \text{if } y < 0.5 \end{cases}$$

$$\hat{t} = \operatorname{argmax} y$$