

# Linear Classification

Monday, 23 May 2022 10:47 AM

## Classification

Supervised learning

$$D = \{(x_n, t_n)\}_{n=1}^N$$

Regression  
 $t_n \in \mathbb{R}$

classification

$t_n$  is discrete

$\in \{\text{'dog', 'cat'}\}$

$\in \{\text{'red', 'blue'}\}$

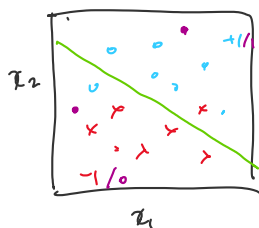
## 2-classes

$$D = \{(x_n, t_n)\}_{n=1}^N$$

features

$\{-1, +1\}$

(or  $\{0, 1\}$  or "one-hot"  
 $\begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ )



$$t_n = \begin{bmatrix} \text{class 0} & \text{class 1} \\ 1 & 0 \end{bmatrix}$$

$$\begin{bmatrix} \text{class 0} & \text{class 1} \\ 0 & 1 \end{bmatrix}$$

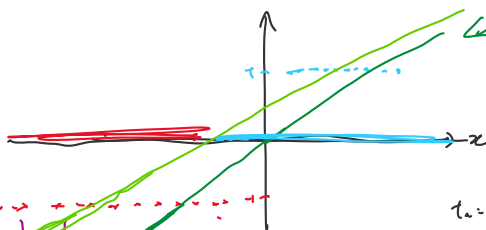
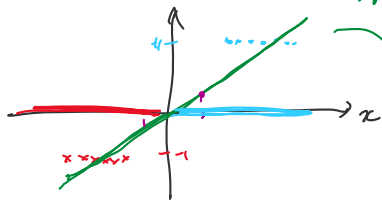
$$\begin{bmatrix} \text{class 0} & \text{class 1} & \text{class 2} \\ 0 & 1 & 0 \\ \vdots & \vdots & \vdots \end{bmatrix}$$

(SLA - Super Lazy Approach)

① use Linear Regression.

$$y = w^T x$$

$$\text{sign}(y)$$

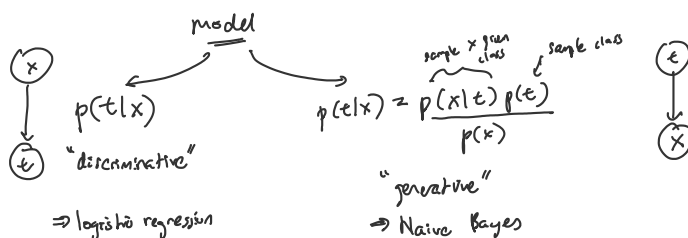


$$t_n = w^T x_n + \epsilon$$

$$\epsilon \sim \mathcal{N}(\epsilon | 0, \beta^{-1})$$

① Design linear model for classification.

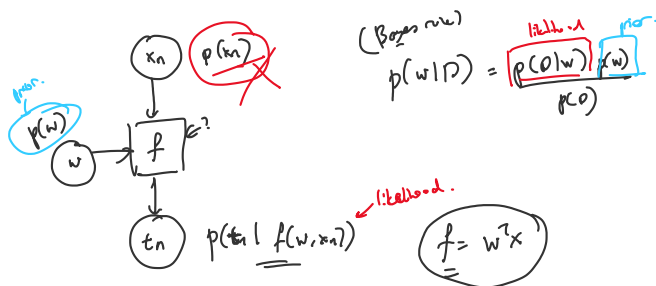
(clever lazy approach)



Discriminative Model

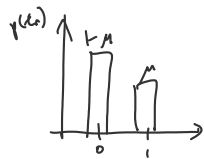
# Logistic Regression (Sec 4.3 of PRML)

$D = \{(x_n, t_n)\}$ .  $t_n \in \{0, 1\}$ . 2-class (binary classification)



prior:  $p(w) = N(0, \alpha^{-1}I)$

likelihood:  $p(D|w) = \prod_{n=1}^N p(t_n | x_n, w)$  (i.i.d)

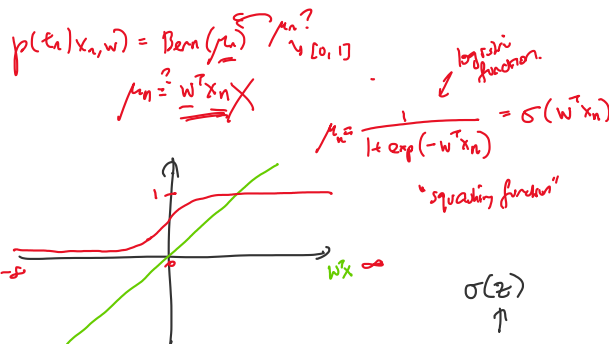


$$p(t_n | x_n, w) = \text{Bern}(\mu_n)$$

$$= \begin{cases} \mu_n & \text{if } t_n = 1 \\ 1 - \mu_n & \text{if } t_n = 0 \end{cases}$$

$$= \mu_n^{t_n} (1 - \mu_n)^{(1-t_n)}$$

$$\begin{cases} E[t_n] = \mu \\ V[t_n] = \mu(1-\mu) \end{cases}$$



$$p(t_n = 1 | w, x_n) = \sigma(w^T x_n)$$

$$p(t_n = 0 | w, x_n) = 1 - \sigma(w^T x_n)$$

$$p(t_n | w, x_n) = \sigma(w^T x_n)^{t_n} (1 - \sigma(w^T x_n))^{(1-t_n)}$$

① Est / inference

- MLE  $\max_w \log p(D|w)$

↑ - MAP  $\max_w \log p(w|D) = \log p(D|w) + \log p(w)$

- Bayes:  $p(w|D)$

MAP

$$\max_w \log \left[ \prod_{n=1}^N p(t_n | x_n, w) \right] + \log p(w)$$

$$= \min_w \left[ \sum_{n=1}^N \log p(t_n | x_n, w) \right] - \log p(w)$$

$$\log p(t_n | x_n, w)$$

$$\log p(w) = \log N(w | 0, \alpha^{-1}I)$$

$$= -\frac{\alpha}{2} w^T w - \frac{1}{2} \log \left[ \frac{1}{(2\pi)^D} \exp \left( -\frac{(w-0)^T}{2\alpha^{-1}I} \right) \right]$$

$$= \lg(\underbrace{\sigma_n}_{\sigma_n = \sigma(w^T x_n)} (1 - \sigma_n))$$

$$= \lg \sigma_n^{t_n} + \lg (1 - \sigma_n)^{1-t_n}$$

$$= t_n \lg \sigma_n + (1-t_n) \lg (1 - \sigma_n) \quad (\text{Cross entropy loss/error})$$

What I want

$$\min_w \mathcal{L}(w)$$

$$\mathcal{L}(w) = - \sum_{n=1}^N \left[ t_n \lg \sigma(w^T x_n) + (1-t_n) \lg (1 - \sigma(w^T x_n)) \right] + \frac{\lambda}{2} w^T w$$

regularizer

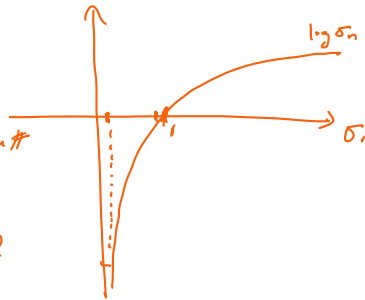
$$- [t_n \lg \sigma_n + (1-t_n) \lg (1 - \sigma_n)]$$

case  $t_n = 1$ , would like  $\sigma_n = ?$

if  $\sigma_n \approx 1$  then  $\lg \sigma_n \approx 0$

$\sigma_n \approx 0$  then  $\lg \sigma_n = \text{negation} \neq$   
 $\Rightarrow \uparrow \text{loss.}$

case  $t_n = 0$ , would like  $\sigma_n = ?$



Solve!

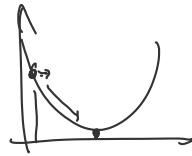
Attempt ①:  $\frac{\partial \mathcal{L}}{\partial w} = 0$ , Solve!

$\Rightarrow$  does not work (no closed form solution)

$$\frac{\partial \mathcal{L}}{\partial w} = 0$$

$$w = (X^T X + \lambda I)^{-1} X^T t$$

Attempt ②: gradient descent



$$\nabla_w \mathcal{L}(w)$$

$$= \nabla_w \left[ - \sum_{n=1}^N t_n \lg \sigma_n + (1-t_n) \lg (1 - \sigma_n) \right] + \nabla_w \left( \frac{\lambda}{2} w^T w \right)$$

exercise

$$= - \sum_{n=1}^N t_n \nabla_w \lg \sigma_n + (1-t_n) \nabla_w \lg (1 - \sigma_n)$$

$$\nabla_w \lg \sigma_n = \nabla_w \lg \sigma(w^T x_n)$$

$$= \frac{1}{\sigma(w^T x_n)} \cdot \sigma(w^T x_n) (1 - \sigma(w^T x_n)) \cdot x_n$$

$$= \frac{(1 - \sigma(w^T x_n)) \cdot x_n}{t_n = 1}$$

$$\nabla_w \lg (1 - \sigma_n) = \nabla_w \lg (1 - \sigma(w^T x_n))$$

$$= \frac{1}{1 - \sigma(w^T x_n)} (-\sigma(w^T x_n) (1 - \sigma(w^T x_n)) \cdot x_n)$$

$$= \frac{(0 - \sigma(w^T x_n)) \cdot x_n}{t_n = 0}$$

$$\Rightarrow (t_n - \sigma(w^T x_n)) x_n$$

FACT:

$$\sigma(z)$$

$$\frac{\partial \sigma}{\partial z} = \sigma(z) (1 - \sigma(z))$$

$$\nabla_w \mathcal{L}(w) = \sum_{n=1}^N \underbrace{(\sigma(w^T x_n) - t_n)}_{\text{error}} \cdot \underline{x_n} + (\text{homework})$$

alg:

$$w_{\text{EHL}} = w_c - \frac{1}{N} \sum_{n=1}^N (t_n - \sigma_n) \cdot x_n$$

Exkurs : ) nonlinear.  $\phi_n$   
 2) k-classes.

$$p(w|x) \propto \boxed{p(y|x)} p(w)$$

Categorical Distribution.



$t_n = [0, 0, 1, 0]$   
 "one-hot" rep / encoding

$$p(t_n | x) = \text{Cat}[\mu] = \prod_{k=1}^K \mu_k^{t_{nk}} = \prod_{k=1}^K \mu_k$$

$$\begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_K \end{bmatrix}$$

$$\underline{\underline{\mu_k = \frac{\exp(w_k^T x_n)}{\sum_k \exp(w_k^T x_n)}}} \quad \text{"softmax"}$$