

# Notes On Machine Learning

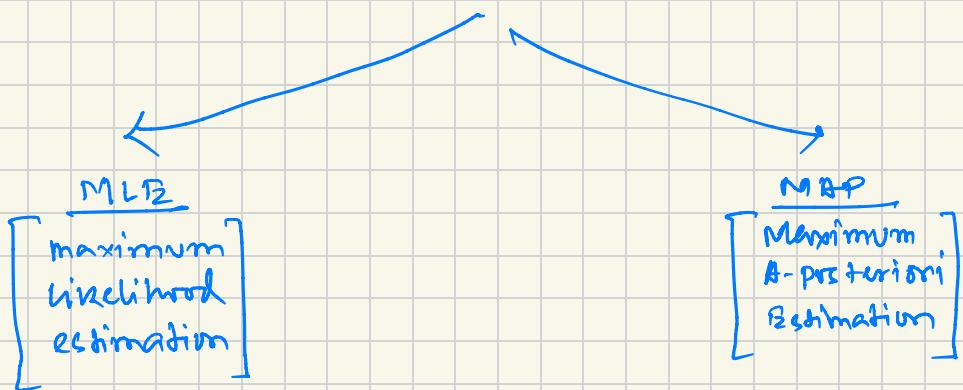
## Parameter Estimation

- \*\* given a parameter value of a distribution how the density/distribution function looks like.
- \*\* Sampling: If  $x$  follows a distribution, how can we obtain different values of  $x$  will take.
- \*\* Parameter estimation: Given different values of  $x$ , how can we obtain the parameters of the underlying function.
  - \* Given data, statistics/ML aim to learn parameters of the underlying distribution from data.
    - Gaussian Mixture Model / Probabilistic graphical model / linear regression / logistic regression.
  - \* Sampling is essential in probabilistic machine learning & Bayesian statistics.
    - Latent dirichlet allocation / Bayesian process / probabilistic graphical model.
    - Markov chain model (MCMC - Markov chain Monte Carlo).
- \* Any statistic used to estimate the value of an unknown parameter  $\theta$  is called an estimator of  $\theta$ .
  - Ex -  $\mu$  &  $\sigma$  for normal distribution  
 $\lambda$  for precision.

The strategy is to choose the parameter in such a way that maximize the likelihood between the prediction & the label or minimize the loss between the prediction & true label.

Once you define loss function only need to optimize the loss function with the parameter.

Two ways to defining loss function



## ① LSE (Least Square Error)

This is special form of MLE. we may define the loss function as:

$$\hat{w}_{LSE} = \underset{w}{\operatorname{argmin}} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

$$= \underset{w}{\operatorname{argmin}} \sum_{i=1}^N (x_i w_i - y_i)^2$$

$$= \underset{w}{\operatorname{argmin}} (x w - y)^2$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_p \end{bmatrix}$$

$$x = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & & & \\ x_m & \cdots & \cdots & x_{mp} \end{bmatrix}$$

This is closed form of above

Now we need to choose  $w$  that minimizes

$$(x w - y)^2$$

So, Loss function

$$L(w) = \gamma_2 \sum_{i=1}^N (x_i w_i - y_i)^2 \\ = \gamma_2 (xw - y)^T$$

$$\frac{\partial L(w)}{\partial w} \Rightarrow (xw - y)x^T = 0 \\ \Rightarrow x^T x w - x^T y = 0 \\ \Rightarrow w = (x^T x)^{-1} x^T y$$

Thus,  $\hat{w} = (x^T x)^{-1} x^T y$  gives least square estimation for our parameter  $w$ .

## ② Maximum Likelihood Estimation (MLE)

MLE to find the parameters that maximize the likelihood between  $y$  &  $\hat{y}$ .

We can define the loss function of MLE as

$$\hat{\theta}_{MLE} = \operatorname{argmax}_\theta P(D|\theta)$$

this represents probability on the dataset D, calculated by the parameter  $\theta$ .

Parameter to estimate which is similar to  $w$

D is the dataset

This is similar to  $P(\hat{x}|y) = P(x,w|y)$

## Ex: MLE for Bernoulli parameter!

$$\text{prob} = P(x) = \begin{cases} 1-p & \text{for } x=0 \\ p & \text{for } x=1 \end{cases}$$

It is also written as  $p^n(1-p)^{n-1}$  as a closed form.

$$f(x_1, x_2, \dots, x_n | p) = p(x=x_1, x=x_2, x=x_3, \dots, x=x_n | p)$$

$$\Rightarrow p(x=x_1 | p) p(x=x_2 | p) \cdots p(x=x_n | p)$$

$$\Rightarrow p^{x_1} (1-p)^{1-x_1} \cdot p^{x_2} (1-p)^{1-x_2} \cdots p^{x_n} (1-p)^{1-x_n}$$

$$\Rightarrow p^{\sum_{i=1}^n x_i} (1-p)^{\sum_{i=1}^n 1 - n}$$

taking log both side to find log likelihood.

$$\log f(x_1, x_2, \dots, x_n | p) = \sum_{i=1}^n x_i \log p + \left( \sum_{i=1}^n 1 - n \right) \log (1-p)$$

$$\frac{\partial}{\partial p} (\log f(x_1, x_2, \dots, x_n | p)) = \frac{1}{p} \sum_{i=1}^n x_i + \frac{1}{1-p} \left( \sum_{i=1}^n 1 - n \right) = 0$$

$$\Rightarrow \frac{1}{p} \sum_{i=1}^n x_i + \frac{1}{1-p} \sum_{i=1}^n 1 - n = 0$$

$$\Rightarrow \sum_{i=1}^n x_i \left( \frac{1}{p} + \frac{1}{1-p} \right) = \frac{n}{1-p}$$

$$\Rightarrow \sum_{i=1}^n x_i \left( \frac{1-p+p}{p(1-p)} \right) = \frac{n}{1-p}$$

$$\Rightarrow p = \frac{1}{n} \sum_{i=1}^n x_i$$

## MLB for normal Distribution:

$$P(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

take log likelihood,

$$\log P(x | \mu, \sigma^2) = \log \left[ \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}} \right]$$

$$L = \sum_{i=1}^n \log \left( \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}} \right)$$

$$L = \sum_{i=1}^n \log \left( \frac{1}{\sqrt{2\pi}\sigma} \right) + \sum_{i=1}^n w_i e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$

$$L_2 = \sum_{i=1}^n w_i \left( \frac{1}{\sqrt{2\pi}\sigma} \right) + \sum_{i=1}^n -\frac{(x_i-\mu)^2}{2\sigma^2}$$

$$\frac{\partial L}{\partial \mu} = \sum_{i=1}^n -\frac{1}{2\sigma^2} 2(x_i-\mu) (-1)$$

$$= \sum_{i=1}^n \frac{x_i - \mu}{\sigma^2} = 0 \Rightarrow$$

$$\sum_{i=1}^n x_i - n\mu = 0$$

$$\Rightarrow \boxed{\mu = \frac{1}{n} \sum_{i=1}^n x_i}$$

$$\frac{\partial L}{\partial \sigma} = \frac{\partial}{\partial \sigma} n \log(2\pi\sigma^2)^{1/2} + -\frac{1}{2} \sum_{i=1}^n (\mu - \mu) \frac{\partial}{\partial \sigma} (\sigma^{-2})$$

$$= -\frac{n}{2} \frac{1}{2\pi\sigma^4} * 2\pi\sigma^2 + -\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2 - 2\sigma^{-3}$$

$$\Rightarrow -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^2 = 0$$

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 = \sigma^2$$

$$\Rightarrow \boxed{\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2}}$$

### ③ Maximum A-posteriori Estimation (MAP)

$$\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} P(\theta | D) = \underset{\theta}{\operatorname{argmax}} \frac{P(D|\theta) P(\theta)}{P(D)}$$

$$\propto \underset{\theta}{\operatorname{argmax}} P(D|\theta) P(\theta)$$

Since  $P(D)$  in the above expression doesn't depend on  $\theta$ , we can just maximize the numerator  $P(D|\theta) P(\theta)$ .

$$\begin{aligned}\hat{\theta}_{MAP} &= \underset{\theta}{\operatorname{argmax}} P(D|\theta) P(\theta) \\ &= \underset{\theta}{\operatorname{argmax}} \left\{ \log(P(D|\theta)) + \log P(\theta) \right\} \\ &= \underset{\theta}{\operatorname{argmax}} \left\{ \sum_{i=1}^n \log P(D_i|\theta) + \log P(\theta) \right\}\end{aligned}$$

This is same as MLE with extra log-prior-distribution term.

MAP allows incorporating our prior knowledge about  $\theta$  in its distribution.

\* prior belief of  $p$ :

→ This prior belief of  $p$  will be 0 outside  $[0,1]$

→ In most cases we want to choose a distribution for prior beliefs that peaks somewhere in the interval  $[0,1]$ .

□ Beta distribution:

Parameterized by two shape constants  $\alpha \geq \beta$ .

Prob( $p | \alpha, \beta$ ) =

$$\frac{1}{B(\alpha, \beta)} p^{\alpha-1} (1-p)^{\beta-1}$$

$$B = \frac{\Gamma(\alpha) \Gamma(\beta)}{\Gamma(\alpha + \beta)}$$

where  $\Gamma_n = (n-1)!$

### 1.1 Definition

Beta distribution is a type of statistical distribution, which has two free parameters. It is used as a prior distribution in Bayesian inference, due to the fact that it is the conjugate prior distribution for the binomial distribution, which means that the posterior distribution and the prior distribution are in the same family.

The probability distribution function (pdf) of the beta distribution is defined as,

$$f(x; \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha) + \Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1} = \frac{x^{\alpha-1} (1-x)^{\beta-1}}{B(\alpha, \beta)} \quad (1)$$

where  $\alpha > 0, \beta > 0, x \in [0, 1]$  and  $\Gamma(\cdot)$  denotes the gamma function.

### 1.2 Application

Considering the classical Bernoulli problem (repeated coin flipping), after  $n$  trials, there are  $s$  successes (heads) and  $f$  failures (tails). Let a random variable  $x$  denote the success probability of each trial. The likelihood for parameters  $s$  and  $f$  given  $x = p$  is the following binomial distribution,

$$L(s, f | x = p) = \binom{n}{s} x^s (1-x)^{n-s} \quad (2)$$

If belief about prior probability information is reasonably well approximated by a beta distribution,

$$P(x = p; \alpha, \beta) = \frac{x^{\alpha-1} (1-x)^{\beta-1}}{B(\alpha, \beta)} \quad (3)$$

According to Bayes's theorem, the posterior probability is given by the product of the likelihood function and the prior probability normalised by the integral as follows,

$$\begin{aligned} P(x = p | s, f) &= \frac{L(s, f | x = p) P(x = p; \alpha, \beta)}{\int_0^1 L(s, f | x = p) P(x = p; \alpha, \beta) dx} \\ &= \frac{\binom{n}{s} x^s (1-x)^{n-s} x^{\alpha-1} (1-x)^{\beta-1}}{\int_0^1 \binom{n}{s} x^s (1-x)^{n-s} x^{\alpha-1} (1-x)^{\beta-1} dx} \\ &= \frac{x^{s+\alpha-1} (1-x)^{n-s+\beta-1}}{B(s+\alpha, n-s+\beta)} \end{aligned} \quad (4)$$

The posterior probability function is also a beta distribution (conjugate). It is convenient to compute. This is the main reason why we approximate the prior using a beta distribution.

\* Derive a-posteriori estimation for the parameter  $\theta$ :

$$\hat{P}_{MAP} = \left( \prod_{i=1}^n P_{\text{Bernoulli}}(x_i | \theta) \right) \cdot P_{\text{Beta}}(\theta | \alpha, \beta)$$

$$L = \sum_{i=1}^n \log P_{\text{Bernoulli}}(x_i | \theta) + \log P_{\text{Beta}}(\theta | \alpha, \beta)$$

$$= \sum_{i=1}^n \log \left[ \theta^{x_i} * (1-\theta)^{x_i-1} \right] + \log \frac{\alpha^{\alpha-1} (1-\theta)^{\beta-1}}{B(\alpha, \beta)}$$

$$L = \sum_{i=1}^n \log \theta^{x_i} + \sum_{i=1}^n \log (1-\theta)^{x_i-1} + \log \theta^{\alpha-1} + \log (1-\theta)^{\beta-1} - \log B(\alpha, \beta)$$

$$L = \sum_{i=1}^n x_i \log \theta + \sum_{i=1}^n (x_i-1) \log (1-\theta) + (\alpha-1) \log \theta + (\beta-1) \log (1-\theta) - \log B(\alpha, \beta)$$

$$\frac{\partial L}{\partial \theta} = \sum_{i=1}^n \frac{x_i}{\theta} - \sum_{i=1}^n \frac{(x_i-1)}{1-\theta} + \frac{\alpha-1}{\theta} - \frac{\beta-1}{1-\theta} = 0$$

$$\Rightarrow \frac{1}{\theta} \sum_{i=1}^n x_i + \frac{1}{1-\theta} \sum_{i=1}^n x_i - \frac{n}{1-\theta} + \frac{\alpha-1}{\theta} - \frac{\beta-1}{1-\theta} = 0$$

$$\Rightarrow \frac{1}{\theta} \left[ \sum_{i=1}^n x_i + \alpha - 1 \right] + \frac{1}{1-\theta} \left[ \sum_{i=1}^n x_i - n - \beta + 1 \right] = 0$$

$$\Rightarrow (1-\theta) \left[ \sum_{i=1}^n x_i + \alpha - 1 \right] + \theta \sum_{i=1}^n x_i - n\theta - \beta\theta + \theta = 0$$

$$\Rightarrow \sum_{i=1}^n x_i + \alpha - 1 - \theta \sum_{i=1}^n x_i - \alpha\theta + \theta + \theta \sum_{i=1}^n x_i - n\theta - \beta\theta + \theta = 0$$

$$\Rightarrow \sum_{i=1}^n x_i + \alpha - 1 + \theta [\alpha - \beta - n + 2] = 0$$

$$\Rightarrow \theta [n + \alpha + \beta - 2] = \sum_{i=1}^n x_i + \alpha - 1$$

$$\Rightarrow \theta_{MAP} = \frac{\sum_{i=1}^n x_i + \alpha - 1}{n + \alpha + \beta - 2}$$

## \* Apply MAP on multinomial distribution

The multinomial distribution is a common distribution for characterizing categorical variables. Suppose a random variable  $Z$  has  $k$  categories, we can code each category as an integer, leading to  $Z \in \{1, 2, \dots, k\}$ . Suppose that  $P(Z = k) = p_k$ . The parameter  $\{p_1, \dots, p_k\}$  describes the entire distribution of  $k$  (with the constraint that  $\sum_j p_j = 1$ ). Suppose we generate  $Z_1, \dots, Z_n$  IID from the above distributions and let

$$X_j = \sum_{i=1}^n I(Z_i = j) = \# \text{ of observations in the category } j.$$

Then the random vector  $X = (X_1, \dots, X_k)$  is said to be from a multinomial distribution with parameter  $(n, p_1, \dots, p_k)$ . We often write

$$X \sim M_k(n; p_1, \dots, p_k)$$

to denote a multinomial distribution.

**Example (pet lovers).** The following is a hypothetical dataset about how many students prefer a particular animal as a pet. Each row (except the 'total') can be viewed as a random vector from a multinomial distribution. For instance, the first row  $(18, 20, 6, 4, 2)$  can be viewed as a random draw from a multinomial distribution  $M_5(n = 50; p_1, \dots, p_5)$ . The second and the third row can be viewed as other random draws from the same distribution.

	cat	dog	rabbit	hamster	fish	total
Class 1	18	20	6	4	2	50
Class 2	15	15	10	5	5	50
Class 3	17	18	8	4	3	50

## 7.1 Properties of multinomial distribution

The PMF of a multinomial distribution has a simple closed-form. If  $X \sim M_k(n; p_1, \dots, p_k)$ , then

$$p(X = x) = p(X_1 = x_1, \dots, X_k = x_k) = \frac{n!}{x_1! x_2! \dots x_k!} p_1^{x_1} \dots p_k^{x_k}.$$

The *multinomial coefficient*  $\frac{n!}{x_1! x_2! \dots x_k!} = \binom{n}{x_1, x_2, \dots, x_k}$  is the number of possible ways to put  $n$  balls into  $k$  boxes. The famous *multinomial expansion* is

$$(a_1 + a_2 + \dots + a_k)^n = \sum_{x_i \geq 0, \sum_i x_i = n} \frac{n!}{x_1! x_2! \dots x_k!} a_1^{x_1} a_2^{x_2} \dots a_k^{x_k}.$$

This implies that  $\sum_{x_i \geq 0, \sum_i x_i = n} p(X = x) = 1$ .

$$L(\theta | D) \propto p_{\text{multinomial}}(x_1, x_2, \dots, x_N | \alpha_1, \alpha_2, \dots, \alpha_N) \cdot p_{\text{Dirichlet}}(\alpha_1, \alpha_2, \dots, \alpha_N | \alpha_1, \alpha_2, \dots, \alpha_N)$$

$$= \ln \prod_{i=1}^N \frac{\alpha_i^{x_i}}{\Gamma(x_i)} \times \left[ \frac{\Gamma(\sum_{j=1}^N \alpha_j)}{\prod_{j=1}^N \Gamma(\alpha_j)} \right] \prod_{i=1}^N \alpha_i^{x_i - 1}$$

$$\begin{aligned} & \text{removing all constant} \\ & \prod_{i=2}^N \alpha_i^{x_i} \quad \prod_{i=2}^N \alpha_i^{x_i - 1} \end{aligned}$$

$$= \prod_{i=1}^n \theta_i^{x_i + \alpha_i - 1}$$

$$L = \log \left( \prod_{i=1}^n \theta_i^{x_i + \alpha_i - 1} \right) = \sum_{i=1}^n (x_i + \alpha_i - 1) \log \theta_i$$

we have constraint  $\sum_{i=1}^n \theta_i = 1$

To impose this constraint we will apply lagrange multiplier  $\lambda$  to log-likelihood.

$$L = \sum_{i=1}^n (x_i + \alpha_i - 1) \log \theta_i - \lambda \left( \sum_{i=1}^n \theta_i - 1 \right) = 0$$

taking derivative on single  $\theta_i$  cancels out other terms

$$\frac{\partial L}{\partial \theta_i} = \frac{x_i + \alpha_i - 1}{\theta_i} - \lambda = 0 \Rightarrow \theta_i = \frac{x_i + \alpha_i - 1}{\lambda}$$

$$\sum_{i=1}^n \theta_i = \frac{1}{\lambda} \sum_{i=1}^n x_i + \alpha_i - 1 = 1$$

$$\Rightarrow \lambda = \sum_{i=1}^n x_i + \alpha_i - 1$$

so, replacing  $\lambda$ 's value

$$\hat{\theta}_{i \text{ MAP}} = \frac{x_i + \alpha_i - 1}{\sum_{i=1}^n x_i + \alpha_i - 1}$$

# Linear Regression

Learn a function which maps input to output

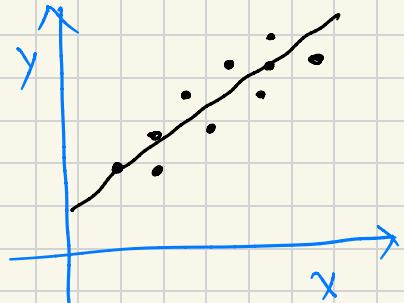
$$f: X \rightarrow Y$$

$$\hat{Y}_i = w_0 + w_1 x_i$$

$$= w_0 + w^T x$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_p \end{bmatrix} \quad p \times 1$$

$$X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & & & \\ x_m & \cdots & \cdots & x_{mp} \end{bmatrix} \quad m \times p$$



$$\text{Loss } L(w) = \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$$L(w) = \frac{1}{2N} \sum_{i=1}^N \{ y_i - (w^T x_i + w_0) \}^2$$

The goal is to minimize this loss function.

We can also write it as

$$\frac{\partial L(w)}{\partial w} = \frac{\partial}{\partial w} \left[ \frac{1}{2N} (xw - y)^T \right]$$

$$\nabla E(w) \Rightarrow \frac{1}{2N} 2 (xw - y)^T = 0$$

$$2) x^T x w = x^T y$$

$$\Rightarrow w = (x^T x)^{-1} x^T y$$

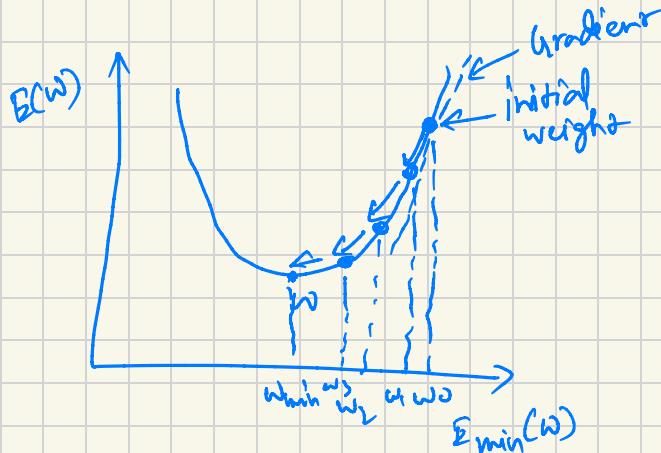
$\underbrace{x^T x}_{P \times N}$     $\underbrace{(x^T x)^{-1}}_{N \times P}$     $\underbrace{x^T y}_{P \times N}$     $\underbrace{y}_{N \times 1}$   
 $\underbrace{P \times P}_{P \times 1}$     $\underbrace{P \times 1}_{P \times 1}$

~~\* \*~~ update weight

$$\omega^{t+1} = \omega^t - \eta \nabla E \rightarrow (x\omega - y)x^T$$

$t$  = iteration number

$\eta$  = learning rate parameter



$$w_1 = w_0 - \eta \nabla E$$

$$w_2 = w_1 - \eta \nabla E$$

$$w_3 = w_2 - \eta \nabla E$$

$$W_{min} = w_3 - \eta \nabla E$$

~~\* \*~~ To apply gradient descent, the function needs to be a convex function.

→ This function will have only one global minimum.

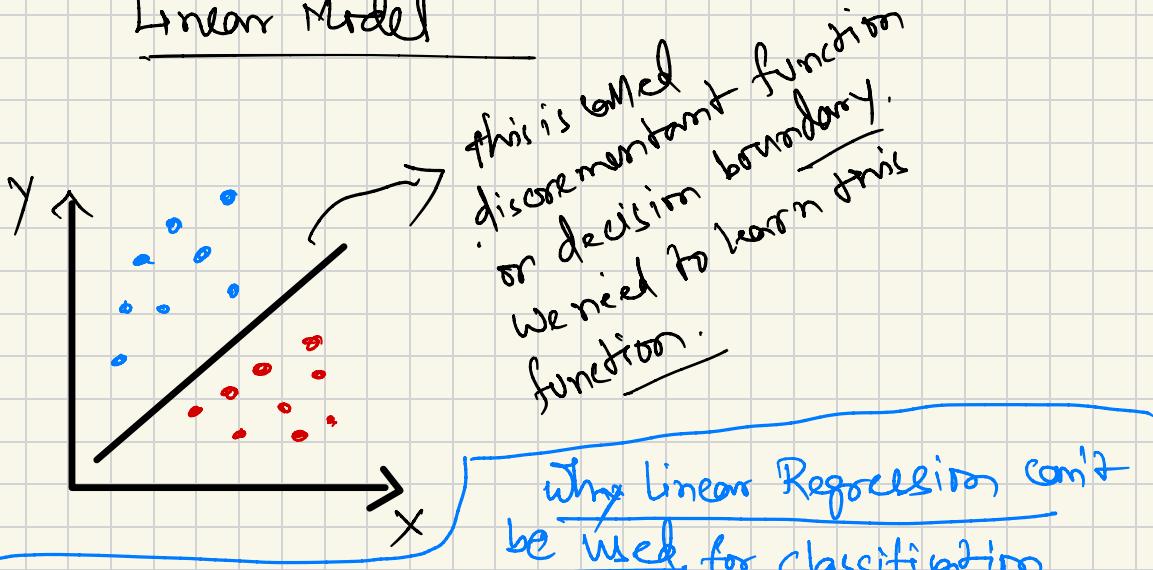
→ It follows a property called Jensen's Inequality.



# Logistic Regression:

## Classification Problem

### Linear Model



→ For classification problem using linear regression, the probability is a Normal distribution. So, output  $P(y=0|x) + P(y=1|x) \neq 1$ .

→ even if you provide binary output to linear regression, it may provide output in real number like 1, 2, -2 etc which doesn't make sense in case of binary classification.

→ To make proper decision we need probabilistic output which sum upto 1 for all classes.

So, use of linear regression for classification problem is not a good option. We need good model which can distribute our probability value within 0 & 1 and for K-class classification within  $\{1, K\}$ .

our task is to,

$$P(y=1|x) \Rightarrow \mathbb{R}^D \xrightarrow{w^T x} \text{sigmoid function} \rightarrow \{0, 1\}$$

also called probit function.

$$\text{If } P(x) = N(x|0, 1)$$

$$\text{CDF}(x) = \int_{-\infty}^x p(x) dx$$

$$\text{CDF}(-\alpha) = 0$$

$$\begin{aligned} \text{CDF}(0) &= \int_{-\infty}^0 p(x) dx \\ &= 0.5 \end{aligned}$$

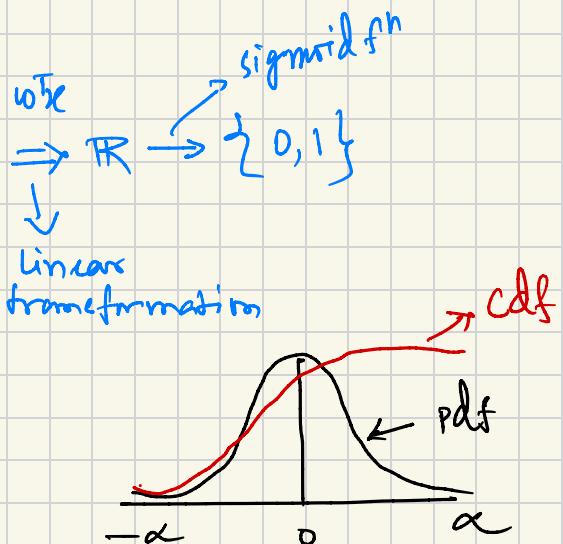
$$\text{CDF}(\alpha) = 1.$$

$$\sigma(-\alpha) = 0$$

$$\sigma(0) = 0.5$$

$$\sigma(\alpha) = 1$$

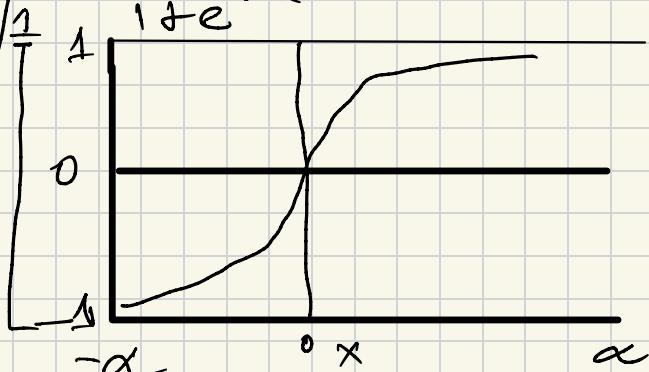
so, sigmoid function we can use to convert  $w^T x$  to  $\{0, 1\}$



Logistic Sigmoid:

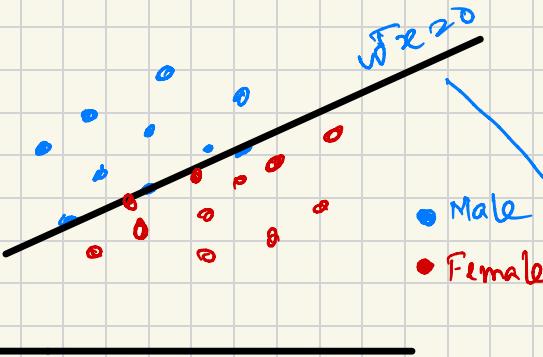
$$P(y=1|x) = \sigma(w^T x)$$

$$= \frac{1}{1 + e^{-w^T x}}$$



$$P(y=1|X) = \sigma(w^T x) = \frac{1}{1+e^{-w^T x}}$$

$$P(y=0|X) = 1 - \sigma(w^T x) = 1 - \frac{1}{1+e^{-w^T x}}$$



$$= \frac{e^{-w^T x}}{1+e^{-w^T x}}$$

$$= \frac{1}{1+e^{w^T x}}$$

→ this is decision boundary with equation  $w^T x = 0$

$$P(y=1|X) = P(y=0|X) = 0.5$$

points on lies on the line

$w^T x = 0$  have same probability

value = 0.5.

points above  $w^T x = 0$  line have high probability towards class 1.

or Male class & points lies below

$w^T x = 0$  have high probability towards class 0 or Female class.

This probability distribution over the output is good for decision making. But sometime we don't need probability to find confidence towards a positive class. In such

scenario we can go with approaches which just provide decision boundary. Ex - LDA, FLD, SVM..

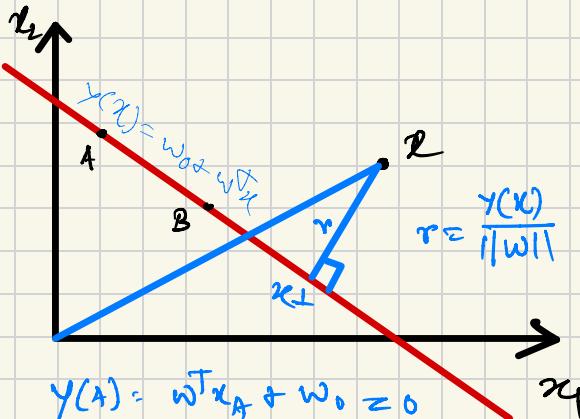
→ All they learn the decision boundary which is called Discriminant function.

## Discriminant function:

$$y(x) = w_0 + w^T x$$

$w_0$  denotes how far your decision boundary from origin. It's a perpendicular direction towards decision boundary from origin.

$w$  is orthogonal to every vector lying within the decision surface, so  $w$  determines the orientation of the decision surface.



$$y(A) = w^T x_A + w_0 = 0$$

$$y(B) = w^T x_B + w_0 = 0$$

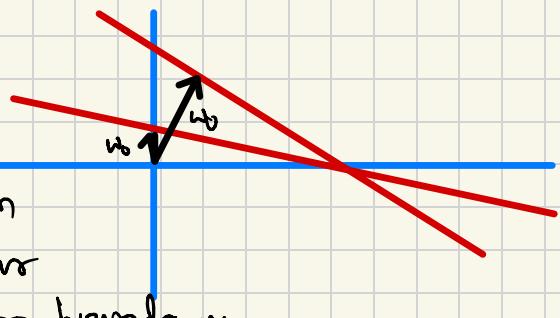
Both are on the line.

so,

$$w^T (x_A - x_B) = 0$$

product of two vectors is 0

so they are perpendicular to each other. so,  $w$  is orthogonal to every vector lying within the decision surface.



### Notes:

Simplistic discriminant function  $y \geq w_0 + w^T x$  represent a hyperplane in  $(m-1)$  dimensional space. where  $m$  is # of features.

# features	Discriminant function
1	point
2	line
3	plane
4	Hyperplane in 3D space
...	...
$m$	Hyperplane in $(m-1)$ D space

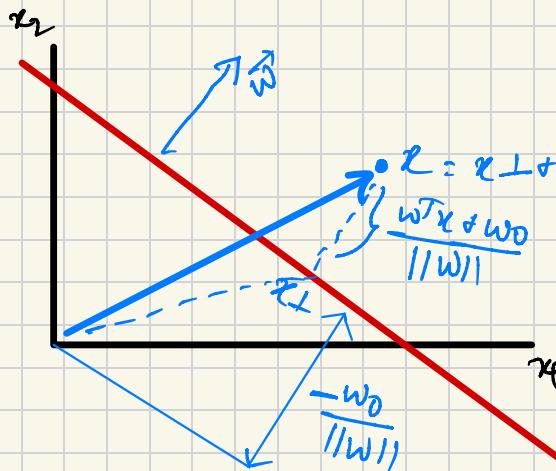
We have  $w_0 + w^T x = 0$

$$w^T x = -w_0$$

Normalizing both sides with the length of the vector  $\|w\|$ , we get normal distance from the origin to the decision surface.

$$d = \frac{w^T x}{\|w\|} = -\frac{w_0}{\|w\|}$$

So,  $w_0$  shifts the boundary away from origin.



this is unit vector for direction.

$$y(x) = w^T x + w_0$$

replacing  $x$  with  $x + r \frac{w}{\|w\|}$

$$y(x) \geq \boxed{w^T x +} + r \frac{w^T w}{\|w\|} + w_0$$

both addition is 0 bcz they are on same line/plane.

$$y(x) = r \frac{w^T w}{\|w\|} = r \frac{\|w\|^2}{\|w\|}$$

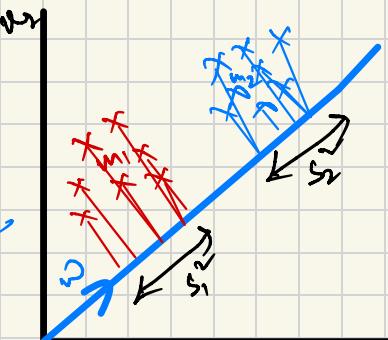
$$y(x) = r \|w\|$$

$$\text{So, } \boxed{r \geq \frac{y(x)}{\|w\|}}$$

$y(x)$  determines the distance of a point to surface.

# Supervised Linear Discriminant Analysis (LDA)

- $w$  as the direction to project  $x$ .
- Find  $w$  such that when  $x$  is projected, classes are well separated.
- Simplest measure of the separation of the classes, when projected onto  $w$ , is the separation of the projected class means.



$$m_1 = \frac{1}{N_1} \sum_{n \in C_1} x_n \quad m_2 = \frac{1}{N_2} \sum_{n \in C_2} x_n$$

we need to maximize  $m_2 - m_1 = w^T(m_2 - m_1)$

this expression can be made arbitrarily large simply by increasing the magnitude of  $w$ . To solve this problem we could constrain  $w$  to have unit length, so that  $\sum w_i^2 = 1$

$$L = \underset{w}{\operatorname{argmax}} w^T(m_2 - m_1) - \lambda \|w\|^2$$

$$\nabla L = (m_2 - m_1) - 2\lambda w = 0$$

$$\Rightarrow w = \frac{1}{2\lambda} (m_1 - m_2)$$

so,  $w \propto (m_1 - m_2)$

value if  $w$  is still depending on mean of two classes.

therefore, we can see there are still considerable overlap when projected onto the line joining their means. To eliminate this, Fisher proposes a function which also considers variance of both class points.

