

Supervised Learning

- Classification
- Regression

Generative

Discriminative.

you need

large amount of

Features $\rightarrow x^{(1)}, x^{(2)}, \dots \rightarrow \bar{x}$ data for this

Labels $\rightarrow \{0, 1\}$

classification problem

Sample = $\{\bar{x}_i, y_i\}$

Labels $\Rightarrow \in \mathbb{R}$

Regression problem.

Generative model

- probabilistic model

underlying assumption (x, y) is from some
 $(x, y) \sim P$ distribution P .

Simplest model

Classification problem.

Ex: Credit default problem.

$$y = \{0, 1\}$$

We will assign some prior probabilities

$$P_0 \triangleq P(y=0) = 0.8$$

$$P_1 \triangleq P(y=1) = 0.2$$

} estimated from historical data.

$$\begin{aligned}
 h_0 &\rightarrow \text{loan to everybody} \\
 h_1 &\rightarrow \text{loan to no one} \\
 h_2(z) &\rightarrow \begin{cases} 0 & \text{if } z=0 \\ 1 & \text{if } z=1 \end{cases} \quad p(z=1) = 0.2
 \end{aligned}$$

How to compare h_0, h_1, h_2 ?

$$\text{Loss} \Rightarrow L(h, y)$$

0-1 Loss,

$$L(1,1) = L(0,0) = 0$$

$$L(1,0) = L(0,1) = 1$$

$$E_y(L(h_0, y)) = \sum_y L(0, y) P_y$$

$$= L(0,0) P_0 + L(0,1) P_1$$

$$= \underline{\underline{0.2}}$$

$$E_y(L(h_1, y)) = \sum_y L(1, y) P_y$$

$$= \underline{\underline{0.8}}$$

$$E_{z,y}(L(h_2, y)) =$$

$$\begin{aligned}
 &L(0,0) P(z=0, y=0) + L(1,0) P(z=1, y=0) + \\
 &L(0,1) P(z=0, y=1) + L(1,1) P(z=1, y=1)
 \end{aligned}$$

$z \perp y$ (independant)

$$\begin{aligned}
 &= 1(P(z=1)) P(y=0) + P(z=0) P(y=1) \\
 &= (1-\beta) 0.2 + \beta (0.8)
 \end{aligned}$$

Best value here $\beta = 0$

$$E(h(z), \gamma) = 0.2$$

"Simply creating a new source of randomness might not help you?" :-)

$$X = \begin{cases} 1 & \rightarrow \text{late tax payer} \\ 0 & \rightarrow \text{not a late tax payer} \end{cases}$$

$$P_0 = P(Y=0) = 0.8$$

$$P_1 = P(Y=1) = 0.2$$

Class conditional density (likelihood).

$$f_1(x) = P(X=x | Y=1)$$

$$f_0(x) = P(X=x | Y=0)$$

$$f_1(1) = P(X=1 | Y=1) = 0.95$$

$$f_0(1) = P(X=1 | Y=0) = 0.1$$

Setup 2

$$X = \{0, 1\}$$

$$h: X \rightarrow \{0, 1\}$$

$$h(0) = ? \quad h(1) = ?$$

X	h ₁	h ₂	h ₃	h ₄
0	0	0	1	1 →
1	1	0	1	0

find the loss for everything? find minimum loss that would be h.

$$E(L(h_{\theta}(x), y))$$

$$= \sum_x \sum_y L(h_{\theta}(x), y) p(x, y)$$

$$= L(1, 0) p(1, 0) + L(1, 1) p(0, 1) \\ L(0, 0) p(1, 0) + L(1, 0) p(1, 1)$$

$$= P(1, 0) + P(1, 1)$$

$$= P(Y=0) P(X=0 | Y=0) + P(Y=1) P(X=1 | Y=1)$$

$$= (0.8)(0.9) + 0.2(0.95)$$

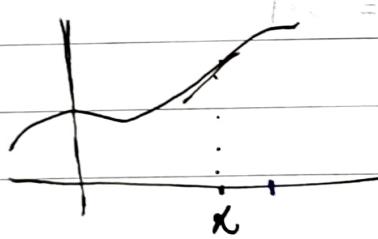
$$= 0.72 + 0.190$$

$$= 0.91$$

Huge Loss!! \rightarrow Bad classifier

Lecture 2 Derivative

10/1/23



The derivative is the slope.
 $f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$

$$f(x+h) \approx f(x) + f'(x)h \quad (\text{approximate a line})$$

Generalisation

$$\vec{f}(\vec{x} + \vec{h}) = f(\vec{x}) + T_f(\vec{h})$$

$\vec{f} \rightarrow$ is a vector valued fn
 $\mathbb{R}^n \rightarrow \mathbb{R}^m$

$$\vec{f} : (f_1, f_2, \dots, f_m) \quad f_i : \mathbb{R}^n \rightarrow \mathbb{R}$$

$T_f(h) \rightarrow$ linear transformation

special kind of transformation satisfying the following

$$1) T(x+y) = T(x) + T(y)$$

$$2) T(\alpha x) = \alpha T(x)$$

conditions:

consider a vector in \mathbb{R}^n

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \alpha_1 \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} + \alpha_2 \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} + \alpha_3 \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

$$T \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \alpha_1 T \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} + \alpha_2 T \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} + \alpha_3 T \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

$$T(e_1) = e_2 \quad T(e_2) = e_3 \quad T(e_3) = e_1$$

(Here T is a simple rotation.)

$$\begin{aligned} T \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} &= 1 T(e_1) + 2 T(e_2) + 3 T(e_3) \\ &= 1 \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} + 2 \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} + 3 \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \\ &= \begin{pmatrix} 3 \\ 1 \\ 2 \end{pmatrix} \end{aligned}$$

This is same as.

$$+ \checkmark \begin{pmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} = \begin{pmatrix} 3 \\ 1 \\ 2 \end{pmatrix}$$

$$\begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & & \frac{\partial f_m}{\partial x_n} \\ \hline & & & mxn \end{pmatrix}$$

T_f

$$f(x+h) = f(x) + T_f h$$

Rules of differentiation.

$$1. h = f + g$$

$\mathbb{R}^n \rightarrow \mathbb{R}^m \quad \mathbb{R}^n \rightarrow \mathbb{R}^m$

$$T(h(x)) = T(f(x)) + T(g(x))$$

$$2. \text{ chain rule: } f: \mathbb{R}^n \rightarrow \mathbb{R}^k$$

$$g: \mathbb{R}^k \rightarrow \mathbb{R}^m$$

$$h = g(f(x)) \quad h: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$3. T_h = T_g(f(x)) \downarrow f. T_f(x)$$

composition opn is
nothing but matrix
multiplication -

$$T_1 = \mathbb{R}^n \rightarrow \mathbb{R}^k \rightarrow (k \times n) \text{ matrix}$$

$$T_2 = \mathbb{R}^k \rightarrow \mathbb{R}^m \rightarrow (m \times k) \text{ matrix}$$

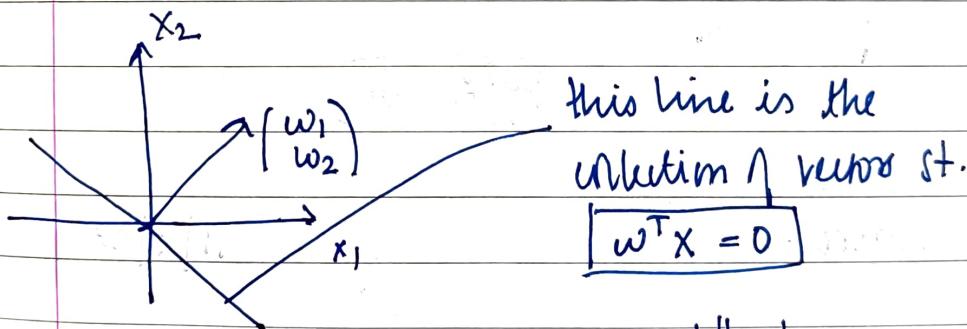
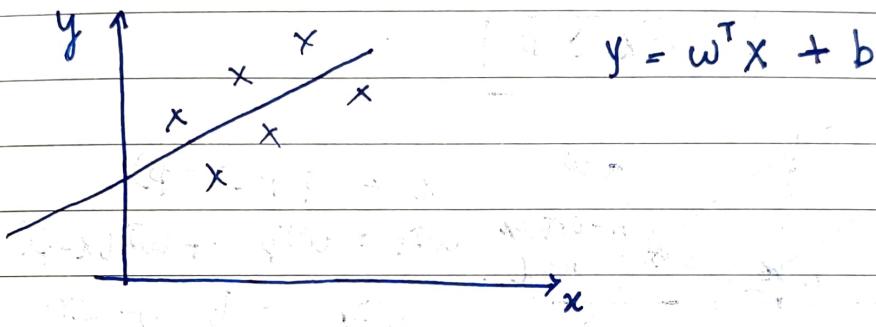
Ex: $f: \mathbb{R}^2 \rightarrow \mathbb{R}^3$

$$f(x, y) = (x^2 + y, xy, x + y)$$

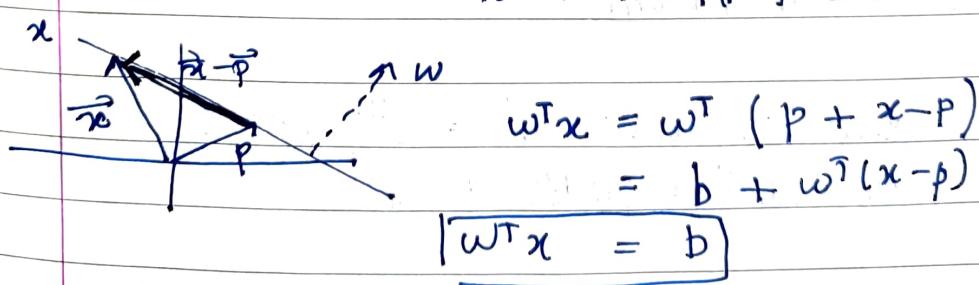
$$\quad \quad \quad f_1 \quad f_2 \quad f_3$$

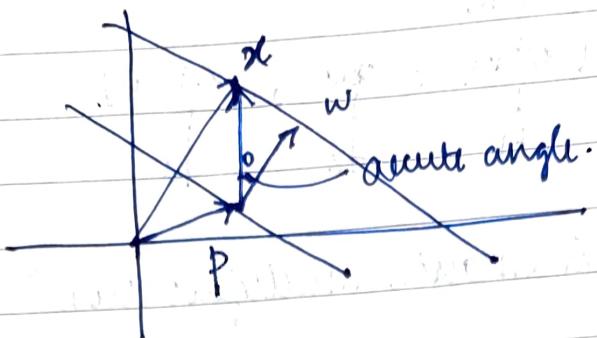
"Chain Rule is matrix multiplication"

A simple Neural Network - Linear Regression



$$w^T x = \|w\| \|x\| \cos \theta$$



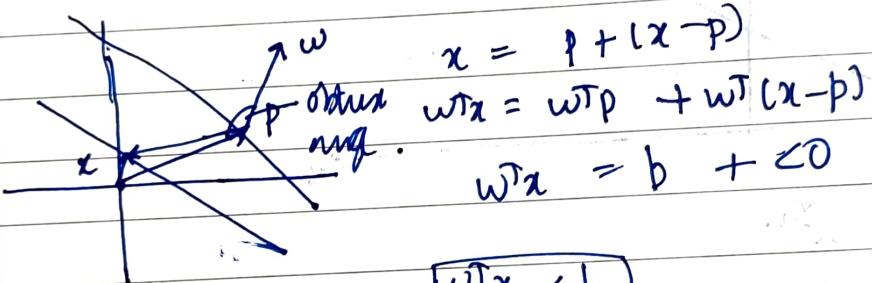


$$x = p + (x - p)$$

$$w^T x = w^T p + w^T (x - p)$$

$$w^T x \neq b + \downarrow > 0$$

$$\boxed{w^T x \geq b}$$



$$x = p + (x - p)$$

$$w^T x = w^T p + w^T (x - p)$$

$$w^T x = b + \downarrow < 0$$

$$\boxed{w^T x < b}$$

Lecture 3

11/1/23

Regression problem.

Given the feature $\bar{x} = x^{(1)}, x^{(2)}, \dots$
we need to predict $y \in \mathbb{R}$

Ex: $x \rightarrow$ features of a person
 $y \rightarrow$ family income.

Through linear regression, we aim to find a line, plane, or hyperplane that classifies the problem.

$$w^T x = b \rightarrow \text{eqn of the line}$$

$$y = w^T x + w_0 \rightarrow (\text{Rearrange})$$

$$\left(\begin{array}{l} \bar{w}_1 x + \bar{w}_2 y + b = 0 \\ y = \frac{-\bar{w}_1}{\bar{w}_2} x - \frac{b}{\bar{w}_2} \end{array} \right)$$

Goal:

Input $\{x_i, y_i\}$

$$x_i^{(1)}, \dots, x_i^{(d)}, y \in \mathbb{R}$$

Find $w = w_1, \dots, w_d, w_0$

$$w^T x_i + w_0 \approx y_i$$

How do we measure this?

Using the loss fn.

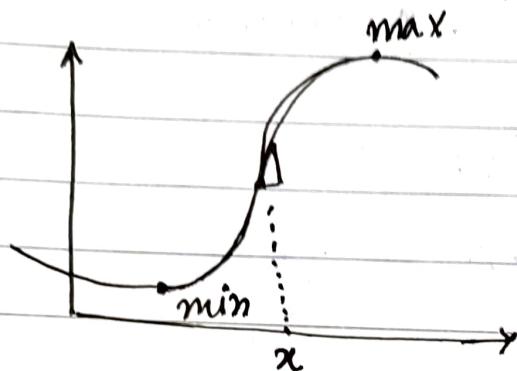
- the standard used to approximate the error is the root mean square error.

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

↳ mean squared error (MSE).

$$\min_w \frac{1}{n} \sum_i (w^T x_i + w_0 - y_i)^2$$

↳ unconstrained optimization problem

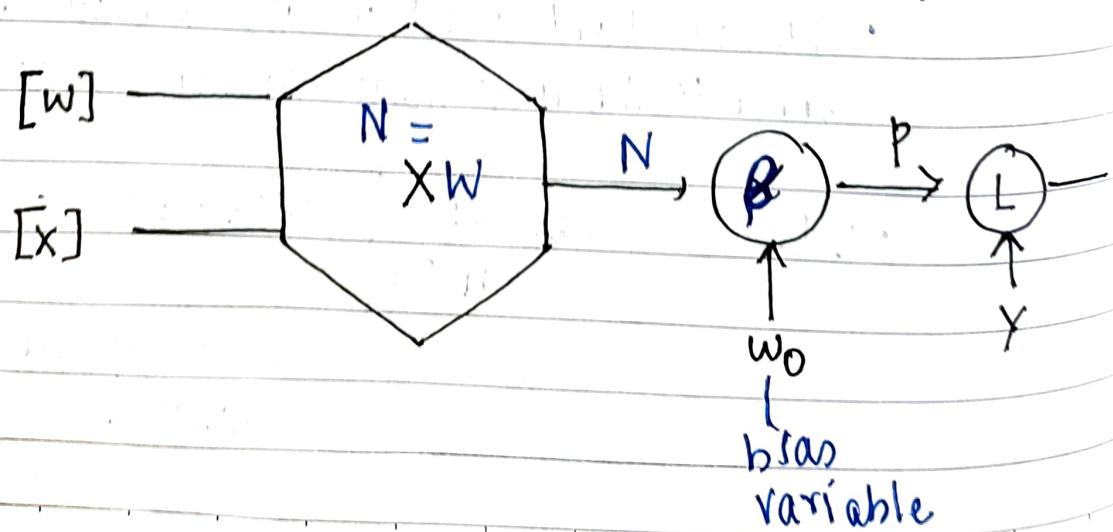


$$\begin{aligned} x + \eta f'(x) &\rightarrow \text{increase} \\ x - \eta f'(x) &\rightarrow \text{decrease} \end{aligned} \quad \left. \begin{array}{l} \text{in 2D space.} \end{array} \right\}$$

In general, going in the direction of increasing gradient will increase y.

So to minimize the loss, we will move in the direction of decreasing gradient.

A simple Neural Network model (linear regression)



$$X = \begin{bmatrix} X_{11} & X_{12} & X_{13} \\ \vdots & \vdots & \\ X_{n1} & \cdots & X_{n3} \end{bmatrix} \quad w = \begin{bmatrix} w_{11} \\ w_{21} \\ w_{31} \end{bmatrix}$$

$$N = Xw \quad P = Xw + w_0$$

$$L = \sum (P_i - Y_i)^2$$

$$L : \mathbb{R}^n \times \mathbb{R}^{10} \rightarrow \mathbb{R}$$

$$L(P, Y) = \frac{\partial L}{\partial P_1} \frac{\partial L}{\partial P_2} \dots \frac{\partial L}{\partial P_n}$$

$$\frac{\partial L}{\partial P_i} = 2 (P_i - Y_i)$$

$$P : \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R}^n$$

↓ |
 (N_1, \dots, N_n) w

$$P = (N_1, \dots, N_n), w_0 + (P_1, \dots, P_n)$$

$$P_i = N_i + w_0$$

$$\frac{\partial L}{\partial N_i} = \frac{\partial L}{\partial P_i} \cdot \underbrace{\frac{\partial P_i}{\partial N_i}}_1 = \frac{\partial L}{\partial P_i}$$

$$\frac{\partial L}{\partial w_0} = \frac{\partial L}{\partial P_1} \frac{\partial P_1}{\partial w_0} + \frac{\partial L}{\partial P_2} \frac{\partial P_2}{\partial w_0} + \dots + \frac{\partial L}{\partial P_n} \frac{\partial P_n}{\partial w_0}$$

$$\frac{\partial L}{\partial w_0} = \sum_{i=1}^m \frac{\partial L}{\partial p_i}$$

$$X : \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ \vdots & & \\ x_{n1} & x_{n2} & x_{n3} \end{bmatrix}$$

$$W : \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix}$$

Linear regression is a simple neural network, here W will be an $n \times 1$ matrix, but generally it will be an $n \times l$ matrix (l is the number of layers)

$$\begin{bmatrix} \frac{\partial L}{\partial w_{11}} & \frac{\partial L}{\partial w_{12}} \\ \vdots & \\ \frac{\partial L}{\partial w_{31}} \end{bmatrix}$$

$$\begin{bmatrix} \frac{\partial L}{\partial N_{11}} & \dots & \frac{\partial L}{\partial N_{12}} \\ \vdots & & \vdots \\ \frac{\partial L}{\partial N_{n1}} & \dots & \frac{\partial L}{\partial N_{n2}} \end{bmatrix}$$

$$\frac{\partial L}{\partial x_{11}} = \frac{\partial L}{\partial N_{11}} \cdot \frac{\partial N_{11}}{\partial x_{11}} + \frac{\partial L}{\partial N_{12}} \cdot \frac{\partial N_{12}}{\partial x_{11}}$$

$$\frac{\partial L}{\partial N_{11}} w_{11} + \frac{\partial L}{\partial N_{12}} w_{12}$$

Input to the backward pass:

$$\begin{bmatrix} \frac{\partial L}{\partial N_{11}} & \frac{\partial L}{\partial N_{12}} \\ \vdots & \vdots \\ \frac{\partial L}{\partial N_{n1}} \end{bmatrix}$$

$n \times 2$

$$\begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \end{bmatrix}$$

2×3

 w^T

$$= \begin{bmatrix} \frac{\partial L}{\partial x_{11}} & \frac{\partial L}{\partial x_{12}} & \frac{\partial L}{\partial x_{13}} \end{bmatrix}^{n \times 3}$$

Lecture 4

features.

16/1/23

$$x = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{1a} \end{bmatrix} \} \text{first sample.}$$

why $\frac{\partial L}{\partial x}$? Because we can have many layers later.

$$w_1 x_1 + w_2$$

$$\frac{dL}{dw}$$

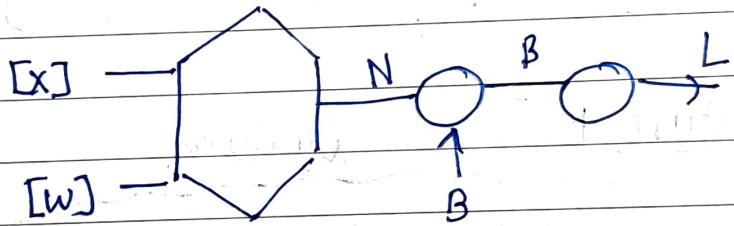
$$N_{11} = x_{11} w_1 + x_{12} w_2 + \dots$$

$$N_{21} = x_{21} w_{11} + x_{22} w_{21} + \dots$$

$$W = \begin{bmatrix} w_{11} & w_{22} \\ w_{21} & \vdots \\ \vdots & \ddots \\ w_{11} & w_{22} \end{bmatrix} \quad N = \begin{bmatrix} N_{11} & N_{12} \\ N_{21} & N_{22} \\ \vdots & \vdots \end{bmatrix}$$

$$\frac{dL}{dw_{11}} = \frac{dL}{dN_{11}} \times \frac{dN_{11}}{dw_{11}} + \frac{dL}{dN_{21}} \frac{dN_{21}}{dw_{11}} + \dots$$

$$= \frac{dL}{dN_{11}} x_{11} + \frac{dL}{dN_{21}} x_{21} + \dots \frac{dL}{dN_{m1}} \frac{dN_{m1}}{dw_{11}}$$



$$\begin{bmatrix} \frac{dL}{dN_{11}} & \frac{dL}{dN_{12}} \\ \frac{dL}{dN_{21}} & \frac{dL}{dN_{22}} \\ \vdots & \vdots \\ \frac{dL}{dN_{m1}} & \frac{dL}{dN_{m2}} \end{bmatrix}$$

$$\begin{bmatrix} \frac{dL}{dW_{11}} & \frac{dL}{dW_{12}} \\ \frac{dL}{dW_{12}} & \vdots \\ \vdots & \end{bmatrix} = \begin{bmatrix} x_{11} & x_{21} & \dots & x_{n1} \\ \vdots & & & \\ x_{13} & & & x_{n3} \\ \vdots & & & \\ x_{1n} & & & x_{nn} \end{bmatrix} \begin{bmatrix} \frac{dL}{dN_{11}} & \frac{dL}{dN_{12}} \\ \vdots & \\ \frac{dL}{dN_{n1}} & \frac{dL}{dN_{n2}} \\ \vdots & \\ n \times 2 \end{bmatrix}$$

"It is a very modular approach, no matter how many layers we have."

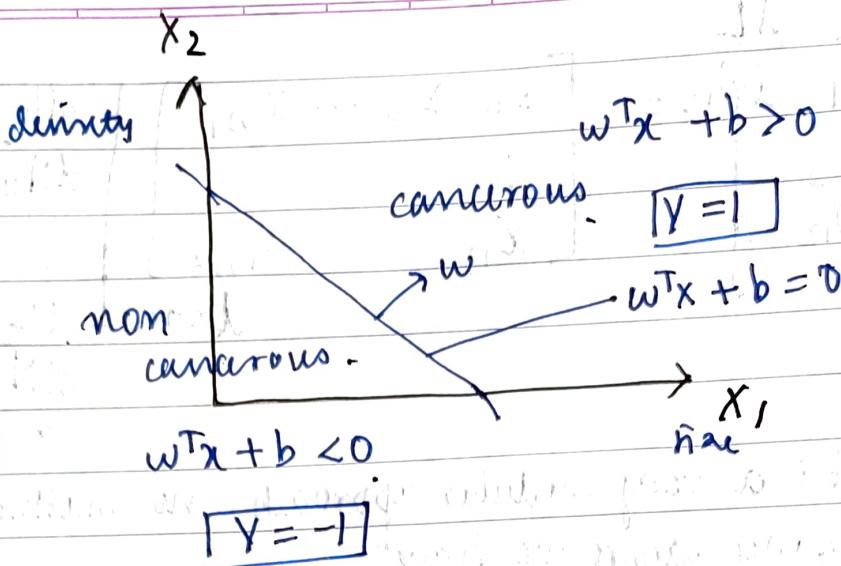
$$S_n = \{(x_i, y_i)\}_{i=1}^n$$

$$FP(S_n) \quad BP(S_n) \quad w^{new} = w^{old} - \alpha \frac{dL}{dw}$$

- If you use the entire training sample, it is called standard gradient descent.
- If you take random batches of training sample, it is called stochastic gradient descent.
- Practically, it has been shown that stochastic GD works better.

Classification:

We can reuse the same model. → used in regression model to solve classification problem. , you will only have to tweak the loss fn.



We want to find $(Y)(w^T x + b) \geq 0$

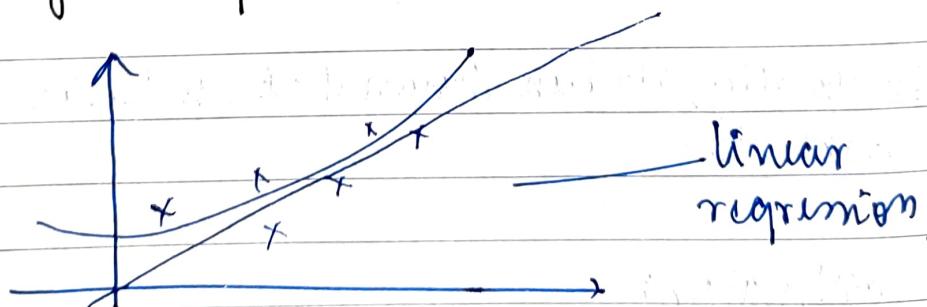
Linear model for classification is called as perceptron.

Loss is called as cross-entropy loss.

Lecture - 5

17/11/23.

Logistic Regression



$$w^T x + b \rightarrow \text{range: } (-\infty, \infty)$$

$$w^T x + b \in [0, 1]$$

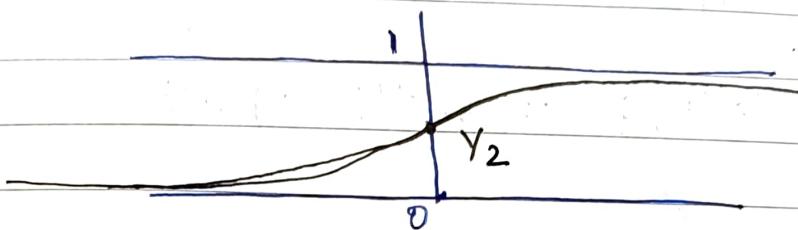
we can treat the num as a probability

sigmoid function :

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

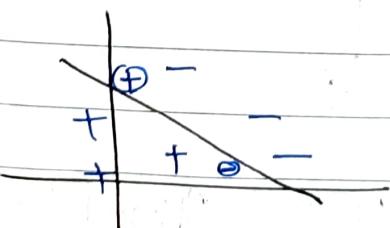
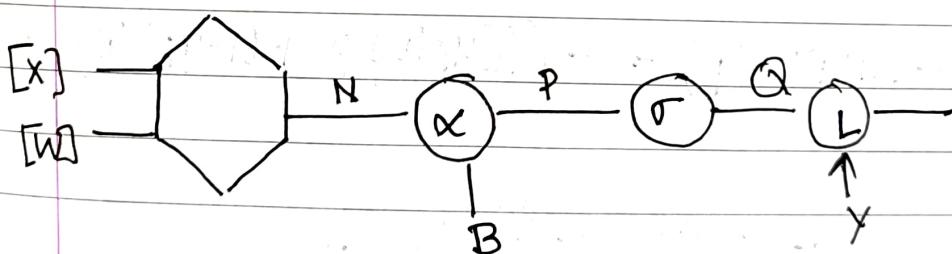
$$z \rightarrow \infty \quad \sigma(z) \rightarrow 1$$

$$z \rightarrow -\infty \quad \sigma(z) \rightarrow 0$$



$$\text{Idea } P(Y=1 | X) = \sigma(w^T x + b)$$

Basically if the value is $>$ than some probability (threshold), then we predict as 1



$$q_i = \sigma(p_i)$$

$$\theta = H \ H \ H \ T \ T \ H \ H$$

$$\mathbb{P}(\theta | p) = p^5 (1-p)^2$$

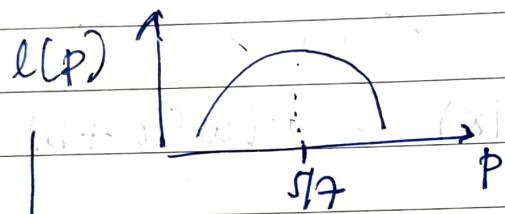
$$\theta = \begin{smallmatrix} 1 & 1 & 1 \\ / & \backslash & \backslash \\ y_1 & y_2 & y_3 \end{smallmatrix}$$

the i^{th} value is $y_i \Rightarrow \{0, 1\}$

$$P(\theta | p) = \prod_{i=1}^n (p)^{y_i} (1-p)^{1-y_i}$$

$$\ln P(\theta | p) =$$

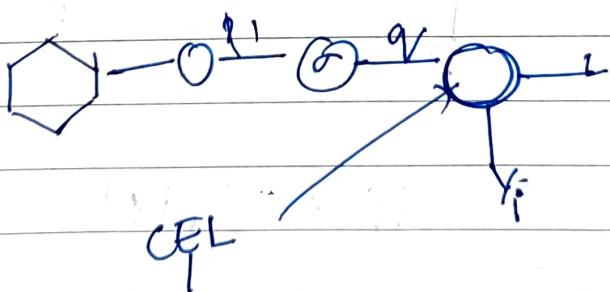
$$\sum y_i \log p + (1-y_i) \log (1-p)$$



maximum likelihood estimator.

cross entropy loss :

$$-\ln P(\theta | p) = -(\sum y_i \log p + (1-y_i) \log (1-p))$$



$$= - \sum y_i \log q_i + (1-y_i) \log (1-q_i)$$

$y_i \rightarrow q_i$. In the optimized var.

$\log q_i \rightarrow q_i \rightarrow 0$

minimize the loss .

$$\frac{dL}{dq_i} = - \left(\frac{u_i}{q_i} + \frac{(1-u_i)}{(1-q_i)} (-1) \right)$$

$$= \frac{1-u_i}{1-q_i} - \frac{u_i}{q_i}$$

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

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

$$\frac{dL}{dP_i} = \frac{dL}{dq_i} \frac{dq_i}{dP_i}$$

$$= \frac{dL}{dq_i} (q_i) (1-q_i)$$

whole training sample

- Batch gradient descent.

for \leftarrow lab
small portions

- minibatch SGD

one training

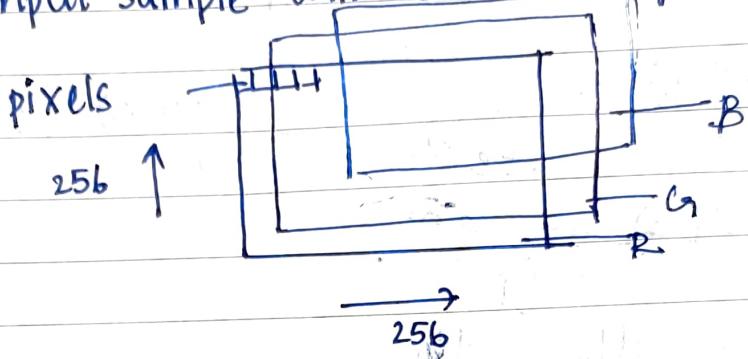
- stochastic gradient descent

~~Σ~~

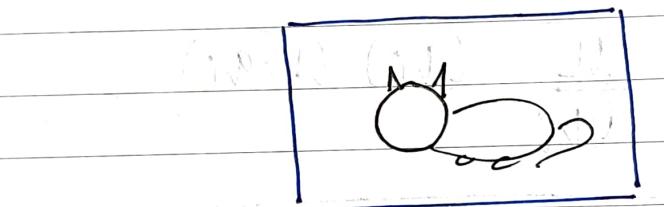
$\sum_k \frac{dL}{dP_i}$

Convolutional Neural Network (CNN)

- used for computer vision problem.
- input sample will be an image



Let us assume there is only a single channel.



filters.

-1	-1	-1
-1	8	-1
-1	-1	-1

9	4	9
4	4	4
4	4	4

0

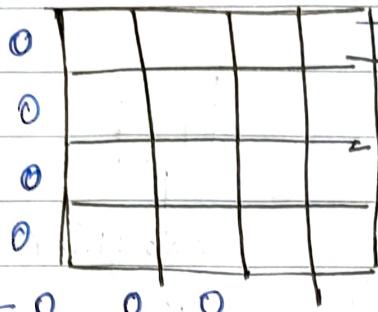
using this filter we will be able to detect the edges.

convolution will give 0 if underlying pixels are all the same i.e. we are trying to learn these values

filter of dimension

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix}$$

design choices? filter size, number of filters to use



→ output image will have dimensions 2×2 with a 3×3 filter.

padding → 1 padding for 3×3 , 2 padding for 5×5 etc...

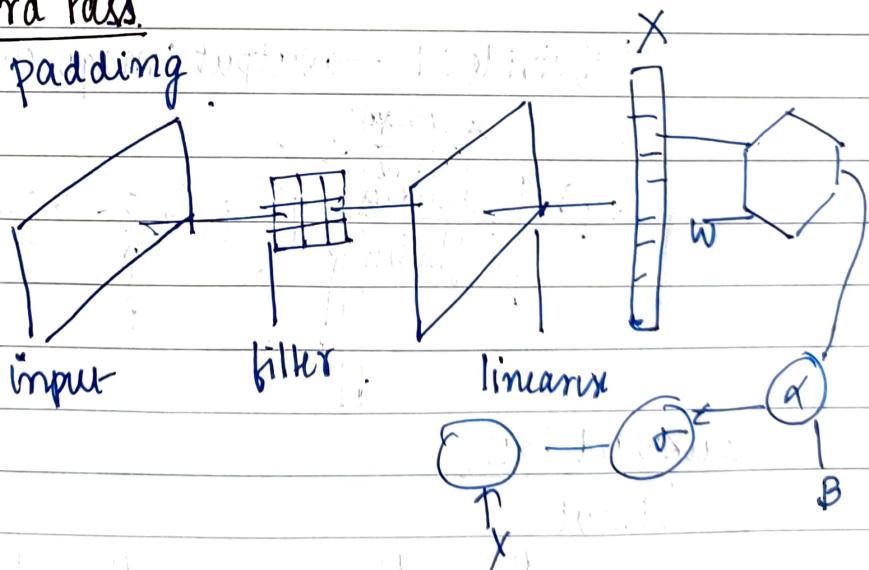
→ now output will be of dimension 4×4 .

→ stride → by how many units do we move the filter.

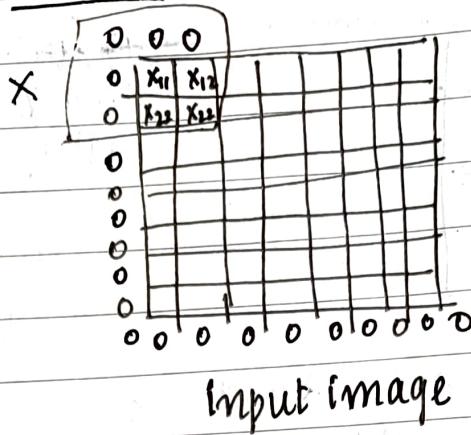
→ If stride is higher, you will get a smaller image.

Forward Pass.

→ padding



17/1/23

Lecture 5

3x3 filter

↓
layer 0
padding

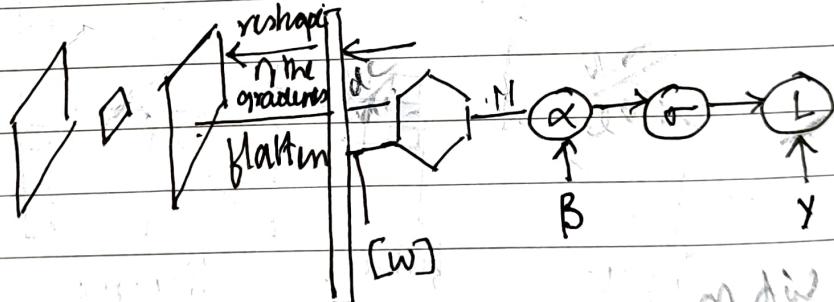
$$\begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix}$$

to maintain the dimension of input image.

$$\text{out}_{11} = x_{11}w_{22} + x_{12}w_{23} + x_{21}w_{32} + x_{22}w_{33}$$

stride = 1

↳ if stride > 1 → output image shrinks



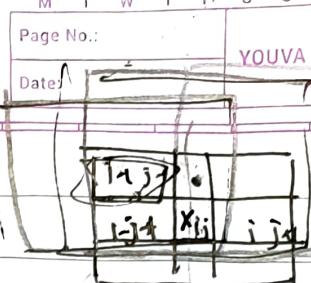
w + m div

Backward Pass

Input

$\frac{dL}{d\theta_{11}}$	$\frac{dL}{d\theta_{12}}$...	$\frac{dL}{d\theta_{m1}}$
$\frac{dL}{d\theta_{21}}$	$\frac{dL}{d\theta_{22}}$		

$$\frac{dL}{dx_{ij}}$$



If we change x_{ij} , what and all outputs change?

$$\hookrightarrow o_{(i-1)(j-1)}, o_{(i-1)(j)}, o_{(i-1)(j+1)} \dots$$

$$\frac{dL}{dx_{ij}}$$

$\frac{dL}{dO_{(i-1)(j-1)}} w_{33}$	$\frac{dL}{dO_{(i-1)j}} w_{32}$	$\frac{dL}{dO_{(i-1)(j+1)}} w_{31}$
$\frac{dL}{dO_{i(j-1)}} w_{23}$	$\frac{dL}{dO_{ij}} w_{22}$	$\frac{dL}{dO_{(i+1)(j+1)}} w_{21}$
w_{13}	w_{12}	w_{11}

$$\frac{dL}{dx_{ij}} = \frac{dL}{dO_{(i-1)(j-1)}} w_{33} + \dots + \frac{dL}{dO_{(i+1)(j+1)}} w_{11}$$

$$\frac{dL}{dw_{11}} = ?$$

Does w_{11} affect O_{11} ? No

~~What is the~~

What does w_{11} affect?

$$O_{22} = x_{11} w_{11} + x_{12} w_{12} + x_{13} w_{13} + \dots$$

$$\frac{dL}{dw_{11}} = \frac{dL}{dO_{22}} x_{11} + \frac{dL}{dO_{23}} x_{12} + \frac{dL}{dO_{24}} x_{13}$$

Claim

0	0	0	0
0	x_{11}	x_{12}	x_{13}
0	x_{21}	x_{22}	x_{23}

$$\begin{array}{cccc}
 & & & 0 \ 0 \\
 & & & x_{14} \ 0 \\
 0 \ x_{11} \ x_{12} \ x_{13} & & & x_{24} \ 0 \\
 & & & x_{34} \ 0 \\
 0 \ x_{21} \ x_{22} \ x_{23} & & & x_{34} \ 0 \\
 0 \ x_{31} \ x_{32} \ x_{33} & & & 0 \ 0 \\
 0 \ 0 \ 0 \ 0 & & & 0 \ 0
 \end{array}$$

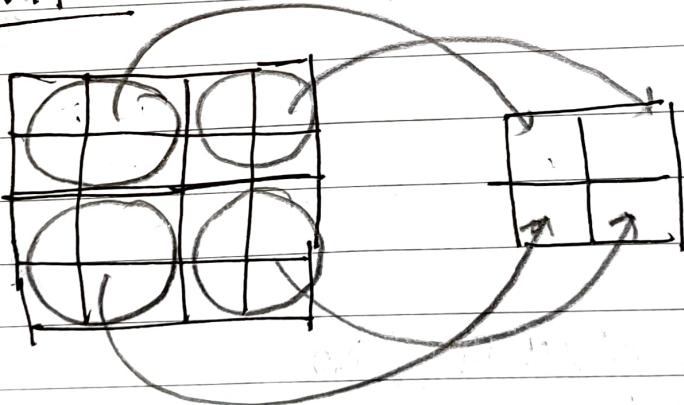
Input image - 3×4

filter - 3×3

Output image - 3×4

convolve \rightarrow Output Image with Input Img.

Maxpool -



Avg pool:

SAtmax \rightarrow Read about it!!

Lecture - 6

18/11/23

Venative Models: \rightarrow we are trying to model
Ex: Credit Default Problem: (x, y)

$$Y \in \{0, 1\}$$

No default \hookrightarrow Default

priors:

$$P_1 = P(Y=1) = 0.2 \quad P_0 = P(Y=0) = 0.8$$

I $h_1 = P(Y=1) = 0.2$ (Always predict 1)
 $P(Y=0) = 0.8$

II $h_2 = 0$ (Always predict 0)

$$\text{Perr}(h_2) = P(h_2 \neq Y) = P(Y=1) = 0.2$$

III $h_3(z) = \begin{cases} 1 & \text{if } z=1 \\ 0 & \text{if } z=0 \end{cases}$

$$\text{Perr}(h_3) = P(h_3 \neq Y)$$

$$\Rightarrow P(h_3)$$

feature space

$X \in \{0, 1\}$

↗ late tax payer.
 ↗ not a late tax payer

priors:

$$P_1 = P(Y=1) = \cancel{0.2} \quad 0.2$$

$$P_0 = P(Y=0) = 0.8$$

$$f_{Y|X}(x) = P(X=x | Y=y)$$

$$f_1(1) = P(X=1 | Y=1) = 0.95$$

$$f_0(1) = P(X=1 | Y=0) = 0.1$$

Goal: We need to find a fn h

$$h : X \rightarrow \{0, 1\}$$

Bayes Rule

$$q_{Y|X}(x) = P(Y=y | X=x)$$

$$= \frac{P_1 f_1(x)}{P_0 f_0(x) + P_1 f_1(x)}$$

$$\therefore q_{Y|X}(x) = \frac{P_1 f_1(x)}{P_0 f_0(x) + P_1 f_1(x)}$$

$$= \frac{P(Y=1) P(X=x | Y=1)}{P(Y=0) P(X=x | Y=0) +}$$

$$h_0(0) = 0 \quad h_0(1) = 1$$

$$q_1(0) = \frac{0.2 \times 0.95}{0.8 \times 0.1 + 0.2 \times 0.95} = \underline{\underline{0.4}}$$

$$q_0(1) = \frac{0.8 \times 0.1}{0.8 \times 0.1 + 0.2 \times 0.95} = 1 - q_1(0) = \underline{\underline{0.3}}$$

$$q_1(0) = \frac{0.2 \times 0.05}{0.8 \times 0.9 + 0.2 \times 0.05} = \underline{\underline{0.013}}$$

$$q_0(0) = 1 - q_1(0) = \underline{\underline{0.987}}$$

Type I error : False Alarm

There is no alarm, but my classifier is predicting 1.

$$\text{TI error} = P(h(x) = 1 | y=0)$$

$$= P(h(x) = 1 | x=1, y=0) P(x=1 | y=0) + P(h(x) = 1 | x=0, y=0)$$

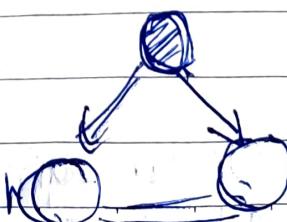
$$P(h(x) = 1 | y=0)$$

$$= P(h(x) = 1 | x=1, y=0) P(x=1 | y=0) + P(h(x) = 1 | x=0, y=0) P(x=0 | y=0)$$

$$= P(h(x) = 1 | x=1, y=0) P(x=1 | y=0)$$

$$= P(h(x) = 1 | x=1) P(x=1 | y=0)$$

$$= \underline{\underline{0.1}}$$



HW
2)

Type II error (missed detection)

$$\text{Perr}(h) = P_{\text{I}} \text{Type I} + P_{\text{II}} \text{Type II}$$

Bayes classifier

$$h_B(x) = \begin{cases} 1 & \{q_1(x) > q_0(x)\} \\ 0 & \{q_0(x) \geq q_1(x)\} \end{cases}$$

$$\text{Perr}(h_B) = P(h_B(x) \neq y)$$

$$= \sum_x P(h_B(x) \neq y | x=x) P(x=x)$$

(We are trying to prove Bayes classifier, is the most optimal classifier: proof of optimality).

$$= \sum_x \left[\sum_y P(h_B(x) \neq y | x=x, y=y) P(y=y | x=x) \right] P(x=x)$$

$$= \sum_x \left[\sum_y L(h_B(x), y) q_{y|x}(x) \right] P(x=x)$$

$$\rightarrow q_{y^*|x}(x) = \max_y q_{y|x}(x)$$

$$= \sum_x [1 - q_{y^*|x}(x)] P(x=x)$$

all the terms will be there
except $\max_y q_{y|x}(x)$

we are trying to minimax this loss, so
and in Bayes Class. we maximise $q_{y^*|x}(x)$

Here the classifier is optimal.

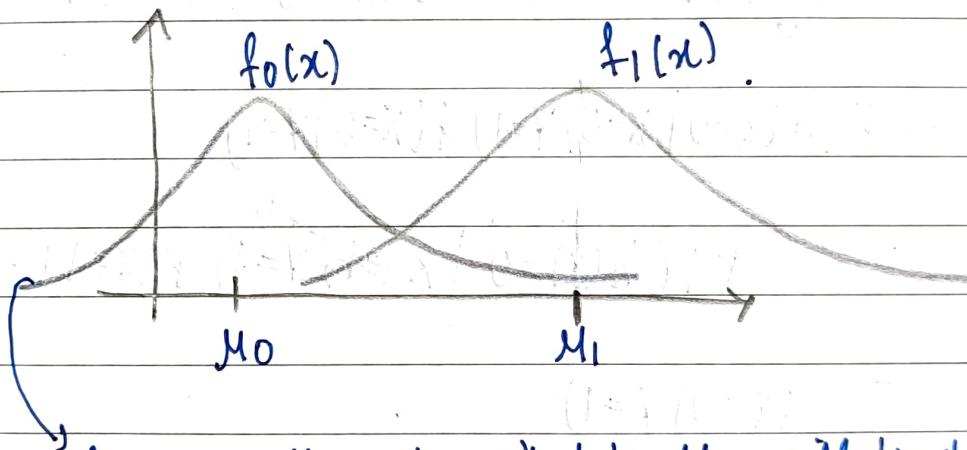
Continuum of Features:

$X \in \mathbb{R}$. (Distribution of time when taxes are paid)

$$p_0 = 0.8, p_1 = 0.2$$

$$f_1(x) = N(\mu_1, \sigma_1^2)$$

$$f_0(x) = N(\mu_0, \sigma_0^2)$$



A person that doesn't default will tend to pay taxes early.

$$q_1(x) = P(Y=1 | X=x) - x \text{ cuz } x \text{ is continuous}$$

$$q_1(x) = P(Y=1 | x \in (x-\delta/2, x+\delta/2))$$

$$q_1(x) = \frac{p_1 f_1(x)}{p_0 f_0(x) + p_1 f_1(x)}$$

$$h_B(x) = 1 \{ q_1(x) \geq q_0(x) \}$$

$$h: h(0) = 0$$

$$h(1) = 1$$

Haw Qn.

Type 2 error (missed detection)

$$h(x) = 0$$

$$y = 1$$

Type 2 error

$$= P(h(x) = 0 \mid y = 1)$$

$$= \sum_x P(h(x) = 0, x \mid y = 1)$$

$$= \sum_x P(h(x) = 0, x = 0 \mid y = 1) + P(h(x) = 0, x = 1 \mid y = 1)$$

$$= P(h(x) = 0 \mid x = 0, y = 1) P(x = 0 \mid y = 1)$$

$$+ P(h(x) = 0 \mid x = 1, y = 1) P(x = 1 \mid y = 1)$$

$$= P(y = 0 \mid y = 1)$$

$$= \underline{0.05}$$

$$P_{\text{err}}(h) = P_1 \text{ type 1 err} + P_2 \text{ type 2 err.}$$

$$\downarrow P_{\text{err}}(h) = P(h(x) \neq y)$$

$$= P(h(x) = 0, y = 1) + P(h(x) = 1, y = 0)$$

$$= P(h(x) = 0 \mid x = 1) P(y = 1) + P(h(x) = 1 \mid y = 0) P(y = 0)$$

$$= (\text{Type 1 error}) P_1 + (\text{Type 2 error}) P_2$$

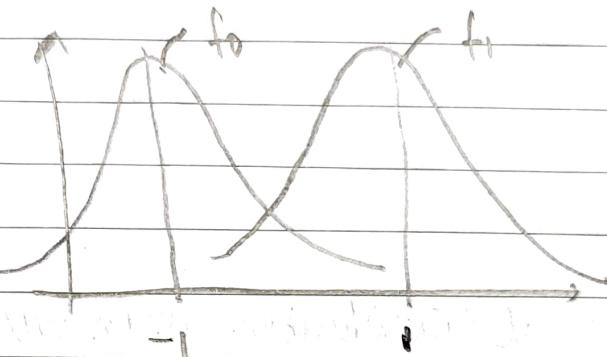
Lecture 7

23/1/23

$X = IR \sim \text{times when taxes are paid.}$

$$f_0(x) \sim N(-1, \sigma^2) \quad p_0 = 1/2$$

$$f_1(x) \sim N(1, \sigma^2) \quad p_1 = 1/2$$



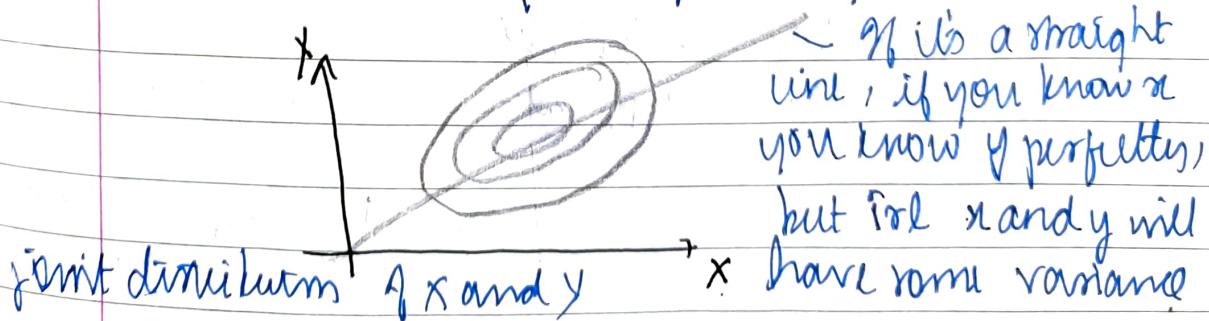
Conditional densities (recap...)

$$X \sim \text{height} \quad Y \sim \text{weight}$$

$$\sim N(55, \sigma_1^2) \quad \sim N(65, \sigma_2^2)$$

↳ Marginals.

$$X, Y \sim N \left(\begin{pmatrix} 55 \\ 65 \end{pmatrix}, \Sigma \right)$$



$$= \frac{1}{(2\pi)^{d/2}} |\det \Sigma|^{-1/2} \exp \left\{ -\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right\}$$

$$\propto \exp \left\{ -\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right\}$$

$$f(x|y) = \frac{P(x=y, Y=70)}{P(Y=70)}$$

$$\approx P[X \in [x - \delta/2, x + \delta/2], Y \in [70 - \delta/2, 70 + \delta/2]]$$

for normal distribution when we fix y ,
it becomes a normal distribution

$$P(Y=1 | X=x) \triangleq q_1(x) = \frac{p_1 f_1(x)}{Z}$$

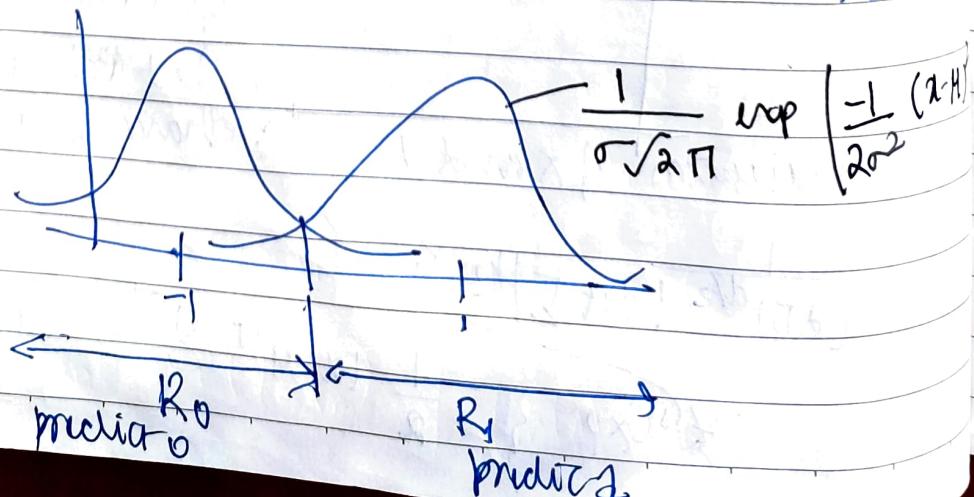
$$Z = p_0 f_0(x) + p_1 f_1(x)$$

$$q_0(x) > q_1(x)$$

$$p_0 f_0(x) > p_1 f_1(x)$$

$$f_0(x) > f_1(x)$$

$$p_0 = p_1 = 1/2$$



$$-(x+1)^2 > -(x-1)^2.$$

$$\boxed{x < 0}$$

$R_0 \rightarrow$ region in my feature space where we report 0.

$R_1 \rightarrow$ report 1.

$$R_0 = \{x \in X \mid h(x) = 0\}$$

PLT:

Type I error: $h(x) = 1$ but $y = 0$.

$$P_E = P(h(x) = 1 \mid y=0)$$

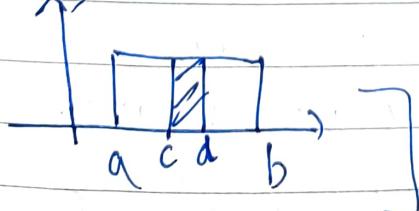
source of randomness
= x .

If x is fixed then $h(x)$.

Diagram

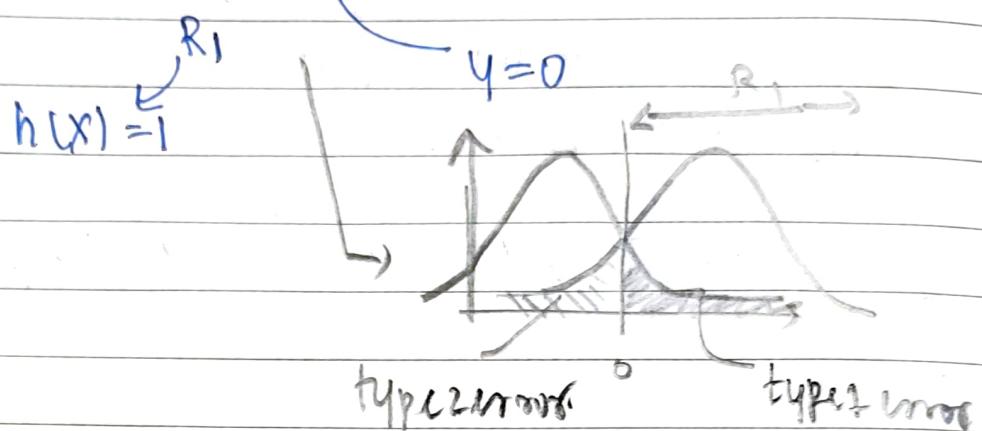
$X \sim \text{Unit } [a, b]$

$$P(X \in [c, d]) = \int_{x \in [c, d]} f_X(x) dx$$



$$P_I = P(h(x) = 1 \mid X=x)$$

$$= \int f_0(x) dx.$$



2) Neyman-Pearson.

$$\text{Power} = \underbrace{\beta_0 \text{ type I}}_{\text{fix}} + \underbrace{\beta_1 \text{ type II}}_{\text{power acc. fixed.}}$$

ROC curve, AUC

Bayesian Decision Theory

$$Y = \{c_1, c_2, \dots, c_m\}$$

there are m classes

$$p_1 = P(Y=c_1) \quad | \quad f_1(x) = P(X=x \mid Y=c_1)$$

:

$$p_m = P(Y=c_m) \quad | \quad f_m(x) = P(X=x \mid Y=c_m)$$

Instead of trying to report one of the m classes.
we will take an action.

$$A = \{\alpha_1, \alpha_2, \dots, \alpha_k\}$$

$$L: A \times Y \rightarrow \mathbb{R}$$

loss function

classifier:

$$h: X \rightarrow A$$

what I want to minimize:

Risk of a classifier.

$$R(h) = E [L(h(x), y)]$$

$$= E_x [E_{y|x} [L(h(x), y) | x]]$$

$$[E[x] = E_y [E_{x|y} [x|y]]]$$

$$E_{x|y} [x|y] = q(*)$$

$$= \int E(L(h(x), y) | x) f(x) dx$$

$$R(h(x) | x) = E(L(h(x), y) | x)$$

$$= \sum_{j=1}^m L(h(x), c_j) q_j(x)$$

Bayes classifier

It chooses.

$$\alpha_{\text{it}}^* = \operatorname{argmin}_{\alpha} \sum_k L(\alpha_j, c_k) q_k(x)$$

Lecture 8

24/1/23

$X \rightarrow \text{feature space}$

$Y \rightarrow \{0, 1\}$

$$\left. \begin{array}{l} p_0 \\ p_1 \\ \vdots \\ p_m \end{array} \right\} \text{priors} \quad f_0(x) \sim N(-1, \sigma^2) \quad f_1(x) \sim N(+1, \sigma^2)$$

$$q_i(x) = P(Y=i | X=x)$$

$$L: A \times Y \rightarrow \mathbb{R}$$

$$R(h(x) | X=x)$$

$$\operatorname{argmin}_{\alpha_j} \sum_k L(\alpha_j, c_k) q_k(x)$$



Bayes classifier.

Ex - (0-1) loss

P_{err} is minimum.

$$h_B(x) = 1 \{ q_1(x) \geq q_0(x) \}$$

Data $\xrightarrow{\text{procedure}}$ Estimate parameters.

given : $S_n = \{(x_i, y_i)\}_{i=1}^n$

$$\hat{p}_0 = \frac{\sum 1 \{ Y_i = 0 \}}{n} \rightarrow \text{Estimator for } p_0.$$

Estimators are some function of your data.

$$\hat{\theta}(x_1, x_2, \dots, x_n) \rightarrow \theta$$

Ex: Suppose there is a coin whose bias is p .

$$S_n = \{(x_i)\}_{i=1}^n \rightarrow \text{draw } n \text{ iid samples.}$$

$$\hat{p} = \frac{\sum x_i}{n} \rightarrow \text{unbiased estimate of } p.$$

- We want to find a good estimator and understand the characteristics of a good estimator.

Desirable properties of an estimator:

i) Unbiasedness

$$E[\theta] = \theta$$

p (fixed but unknown)

$$S_n = \{X_i\}_{i=1}^n$$

$$M_n = \frac{\sum_{i=1}^n X_i}{n}$$

$$M_{\text{odd}} = \frac{X_1 + X_3 + X_5 + \dots}{n_{\text{odd}}}$$

[EW - Expected value of M_n]

$$E[M_n] = p$$

$$E[M_{\text{odd}}] = p$$

$$\therefore E[M_n] = E\left[\frac{\sum X_i}{n}\right]$$

$$= \frac{\sum E(X_i)}{n} = \frac{np}{n} = p$$

X_i s are iid,

$$E(X_i) = \theta(1-p) + p$$

$$= p$$

$$\boxed{E[M_{\text{odd}}] = p}$$

Both have the same unbiasedness, but what is the diff?

$$\text{Var}(\bar{M}_n) = \text{Var}\left(\frac{\sum X_i}{n}\right)$$

$$= \frac{1}{n^2} \text{Var} \sum X_i \quad (X_i \text{ are independent})$$

$$= \frac{1}{n^2} \sum \text{Var} X_i$$

$$= \frac{(p)(1-p)}{n}$$

$$\text{Var}(M_{\text{odd}}) = \frac{p(1-p)}{n_{\text{odd}}}$$

$$\text{Var}(M_{\text{odd}}) > \text{Var}(\bar{M}_n)$$

Variance is big means the errors will be more.

- 2) The estimator should have low variance.
- 3) consistency.

$\hat{\theta}_n \xrightarrow{IP} \theta$ as $n \rightarrow \infty$
 must happen
 convergence in probability

$$\lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta| > \varepsilon) = 0 \quad \forall \varepsilon > 0$$

weak law of large numbers.

Probability of taking 100 samples and P of taking 1000 samples must be closer to P .

1) unbiasedness

2) low variance

3) consistency \rightarrow we are only looking at consistency

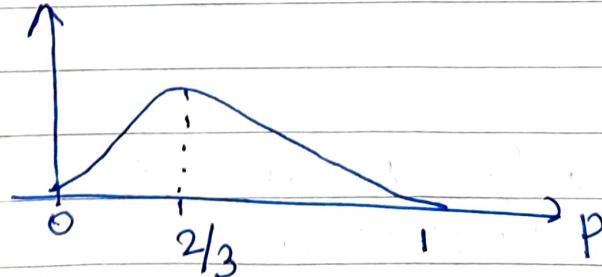
Maximum likelihood Estimation.

$$\mathcal{D} = \{H, H, T\}$$

$$l(p | \mathcal{D}) = P(\mathcal{D} | p)$$

$$\text{if } p=0, l(0 | \mathcal{D})=0$$

$$l(1 | \mathcal{D})=0$$



$$\hat{p}_{MLE} = \operatorname{argmax}_p l(p | \mathcal{D})$$

Lecture 9

25/1/23

Estimator

$$\hat{\theta}(x_1, x_2, \dots, x_n) \rightarrow \theta \text{ (fixed, unknown)}$$

Desirable properties:

$\hat{\theta}_n \Rightarrow$ random variable.

- 1) unbiasedness $E[\hat{\theta}_n] = \theta$
- 2) low variance $V(\hat{\theta}_n)$
- 3) consistency $\hat{\theta}_n \rightarrow \theta$

$$\forall \epsilon > 0 \quad \lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta| > \epsilon) = 0$$

$$\hat{M}_n = \frac{\sum_{i=1}^n x_i}{n} \quad \{x_i\}_{i=1}^n \text{ iid Bern}(p)$$

fix $\epsilon > 0$

$$P(|\hat{M}_n - p| > \epsilon)$$

Markov Inequality

$X > 0$ (non negative rr)

$$E(X) = \sum_{0 \leq x < a} x p_x(x) + \sum_{x \geq a} x p_x(x)$$

positive num.

↑

$$\begin{aligned} E(X) &\geq \sum_{x \geq a} x p_x(x) \\ &\geq \sum a p_x(x) \end{aligned}$$

$$P(X \geq a) \leq \frac{E(X)}{a} \quad \boxed{\text{markov inequality}}$$

union bound

$$P\left(\bigcup_{i=1}^n A_i\right) \leq \sum P(A_i)$$

chebyshov inequality (no restriction $X \in \mathbb{R}$)

$$P(|X| \geq a) = P(X^2 \geq a^2)$$

$$\Rightarrow P(|X - \mu| \geq a)$$

$$\Rightarrow P((X - \mu)^2 \geq a^2)$$

$$\Rightarrow \leq \frac{E(X - \mu)^2}{a^2} \quad (\text{markov inequality})$$

$$\Rightarrow \leq \frac{\text{Var}(X)}{a^2}$$

$$\Rightarrow \hat{M}_n = \frac{1}{n} \sum X_i$$

fix $\epsilon > 0$

$$P(|\hat{M}_n - p| > \epsilon)$$

$$= P(|\hat{M}_n - p|^2 > \epsilon^2)$$

$$\leq E \frac{(\hat{M}_n - p)^2}{\epsilon^2} = \frac{p(1-p)}{n\epsilon^2}$$

$$\therefore \lim_{n \rightarrow \infty} P(|M_n - p| > \epsilon) = 0.$$

$$\tilde{M}_n = \frac{1}{n-1} \sum_{i=1}^n x_i$$

unbiased?

$$E[\tilde{M}_n] = \frac{n}{n-1} \mu \neq \mu$$

$\boxed{\mu = p}$

consistency? - Yes.

$$P(|\tilde{M}_n - p| > \epsilon) \leq \frac{\text{Var}(\tilde{M}_n)}{\epsilon^2}$$

$$\leq \frac{n \sigma^2}{(n-1)^2 \epsilon^2}$$

$$\sigma^2 = (p)(1-p)$$

$$\text{As } n \rightarrow \infty \quad P(|\tilde{M}_n - p| > \epsilon) \rightarrow 0.$$

Maximum Likelihood Estimator.

$$L(\theta | \mathcal{D}) \triangleq p(\mathcal{D} | \theta) \quad (\mathcal{D} = \{x_i\}_{i=1}^n \text{ iid})$$

$$= \prod p(x_i | \theta)$$

MLE

$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} \quad l(\theta | \mathcal{D})$$

MLE for Bernoulli

Assume p ($0 < p < 1$) fixed but unknown

$$L \{x_i\}_{i=1}^n \quad \text{iid} \quad P(x_i = 1) = p$$

$$l(p|\theta) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i}$$

Ex: $\{H, H, T, T, T\} \quad \{1, 1, 0, 0\}$

log likelihood

$$\hat{p}_{MLE} = \underset{p}{\operatorname{argmax}} \sum x_i \log p + (1-x_i) \log(1-p)$$

$$\frac{d}{dp} L(p) = \cancel{\sum x_i} - \cancel{\sum (1-x_i)}$$

$$= \sum \frac{x_i}{p} - \frac{(1-x_i)}{(1-p)} = 0$$

$$= \sum \frac{x_i(1-p) - (1-x_i)p}{p(1-p)}$$

=

$$= \frac{\sum x_i}{n}$$

$$\hat{P}_{MLE} = \frac{\sum x_i}{n}$$

MLE for multinomial.

$$x_i \underbrace{\quad}_{\begin{matrix} 1 \\ 2 \\ \vdots \\ K \end{matrix}}$$

$$p(x_i=1) = p_1$$

$$p(x_i=2) = p_2$$

$$\sum_{i=1}^K p_i = 1$$

$x_i \rightarrow$ one hot encoding

$$x_i = 1, 1, 2, 2, 3, 1, \dots$$

$$\left(\begin{array}{c} 1 \\ 0 \\ 0 \end{array} \right), \left(\begin{array}{c} 0 \\ 1 \\ 0 \end{array} \right), \left(\begin{array}{c} 0 \\ 0 \\ 1 \end{array} \right), \left(\begin{array}{c} 0 \\ 0 \\ 0 \end{array} \right), \left(\begin{array}{c} 0 \\ 0 \\ 0 \end{array} \right)$$

$$x_i \rightarrow \tilde{x}_i = \begin{pmatrix} x_i^{(1)} \\ \vdots \\ x_i^{(K)} \end{pmatrix}$$

MLE

$$l(p_1, p_2, \dots, p_K | D) = \prod_{i=1}^n \prod_{j=1}^K p_j^{x_i^{(j)}}$$

$$L(l(p_1, p_2, \dots, p_K)) = \sum_i \sum_j x_i^{(j)} \log p_j$$

$$p_1, p_2, \dots = \arg \max_{p_1, p_2, \dots} = \sum_i \sum_j x_i^{(j)} \log p_j$$

$$\text{st } \sum p_j = 1 \quad p_j > 0$$

Page No.:
Date: YOUVA

(Bishop - Appendix)
Optimization of Unconstrained

$$\max \sum_{i=1}^n \sum_{j=1}^k x_i^{(j)} \log_e(p_j)$$

subject to $\sum p_i = 1$

$$L(x, \lambda) = f(x) + \lambda g(x)$$

$$\max_{x, \lambda} L(x, \lambda)$$

$$L(p_1, p_2, \dots, p_k, \lambda)$$

$$= \sum_{i=1}^n \sum_{j=1}^k x_i^{(j)} \ln(p_j) + \lambda \left(\sum_{i=1}^n p_i - 1 \right)$$

$$\frac{\partial L}{\partial p_j} = \sum_{i=1}^n \frac{x_i^{(j)}}{p_j} + \lambda = 0$$

$$\boxed{\frac{\partial L}{\partial \lambda} = 0}$$

$$= \sum_{i=1}^n x_i^{(j)} + \lambda p_j = 0$$

Summing over $j = \{1, 2, \dots, k\}$

$$= \sum_{j=1}^k \sum_{i=1}^n x_i^{(j)} + \lambda \sum_{j=1}^k p_j = 0$$

$$= n + \lambda = 0$$

$$\boxed{\lambda = -n}$$

$$\frac{\partial L}{\partial p_j} = \sum x_i^{(j)} + \lambda p_j = 0$$

=

$$p_{j, MLE} = \frac{\sum x_i^{(j)}}{n}$$

30/1/23

MLE for Normal distribution.

$$\mathcal{D} = \{x_i\}_{i=1}^n \quad x_i \sim N(\mu, \sigma^2) \text{ iid.}$$

$$L(\theta | \mathcal{D}) = P(\mathcal{D} | \theta)$$

$$= \prod_{i=1}^n f(x_i | \mu, \sigma^2)$$

$$f(x_i | \mu, \sigma^2) = \frac{1}{\sigma \sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2} \right\}$$

$$\ell(\mu, \sigma^2 | \mathcal{D})$$

$$\ell(\theta | \mathcal{D}) = \sum_{i=1}^n \ln \left(\frac{1}{\sigma \sqrt{2\pi}} \right) - \frac{1}{2\sigma^2} (x_i - \mu)^2$$

$$\frac{\partial L}{\partial \mu} = 0$$

$$\mu_{MLE} = \frac{\sum x_i}{n}$$

$$\frac{\partial L}{\partial \theta} = 0$$

$$\sigma_{MLE}^2 = \frac{1}{n} \sum (x_i - \mu_{MLE})^2$$

Bias variance tradeoff

MSE \rightarrow mean squared error

$$MSE = E(\hat{\theta} - \theta)^2$$

$$\hat{\theta}(x_1, x_2, \dots, x_n)$$

$$= E[\hat{\theta} - E(\hat{\theta}) + E(\hat{\theta}) - \theta]^2$$

$$= E[(\hat{\theta} - E(\hat{\theta}))^2 + (E(\hat{\theta}) - \theta)^2 + 2(\hat{\theta} - E(\hat{\theta}))(E(\hat{\theta}) - \theta)]$$

$$= E(\hat{\theta} - E(\hat{\theta}))^2 + (E(\hat{\theta}) - \theta)^2 + 2(E(\hat{\theta}) - \theta) E[\hat{\theta}]$$

$$= E(\hat{\theta} - E[\hat{\theta}])^2 + (E[(\hat{\theta})] - \theta)^2.$$

$$= \text{Var}(\hat{\theta}) + \text{Bias}^2$$

$$\text{MSE}(\hat{\theta}) = \underbrace{\text{Var}(\hat{\theta})}_{\text{How big the bag is?}} + \underbrace{\text{Bias}^2}_{\text{to cope with the variance you have to reduce the amount of bias in the "bag", this will give us a bias.}}$$

Bayesian Estimators:

Ex: $P(\theta) = 0.1$

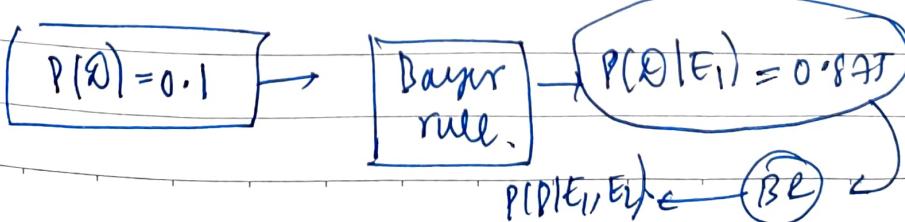
↳ probability of disease.

$$P(+|\theta) = 0.7 \quad P(+|\theta^c) = 0.01$$

$$P(\theta|+) = \frac{P(+|\theta) P(\theta)}{P(+|\theta) P(\theta) + P(+|\theta^c) P(\theta^c)}$$

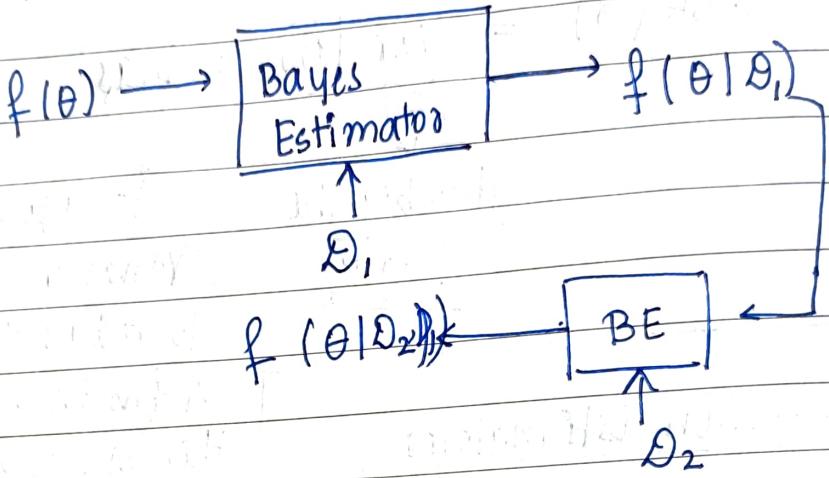
$$= \frac{(0.01)(0.7)}{(0.7)(0.01) + (0.9)(0.01)}$$

$$= \frac{7}{79} = \underline{\underline{0.875}}$$



in MLE - we assumed θ was a fixed unknown.

in Bayes Estimator θ (Random Distribution)
 $f(\theta)$ → prior for θ .
 you have a belief, but it improves more and more as we see more values



$$f(\theta|D_1) = \frac{f(\theta) f(\theta|D_1)}{\int f(\tilde{\theta}) f(\theta|\tilde{\theta}) d\tilde{\theta}}$$

$$= \frac{f(\theta) f(D_1|\theta)}{f(\theta)}$$

[Here MLE cannot be used, because we're getting data piece by piece].

- when we have some prior knowledge, its better to use Bayesian estimator.

conjugate pairs

$$f(\theta | \theta) = f(\theta) f(\theta | \theta)$$

form

coin toss prior:

$$\begin{aligned} f(p) &= \text{Beta}(\alpha, \beta) \quad (\alpha > 0 \\ &\quad \beta > 0) \\ &= \frac{p^{\alpha-1} (1-p)^{\beta-1}}{\int_0^1 p^{\alpha-1} (1-p)^{\beta-1} dp} \\ &= p^{\alpha-1} (1-p)^{\beta-1} (K) \end{aligned}$$

$$f(\theta | p) = p^{\sum x_i} (1-p)^{n - \sum x_i}$$

$$f(p | \theta) = \text{Beta}(\alpha + \sum x_i, \beta + n - \sum x_i)$$

Beta(1,1) \Rightarrow uniform(0,1)

as H, 5T Beta(26,6).

31/1/23

Bayes Estimator

$\theta \sim$ Random Variable

$f(\theta) \sim$ Prior (comes from knowledge of the problem) (we choose it)

$f(\theta | \theta) \sim$ likelihood

$[L(\theta | \theta)]$

→ $f(\theta)$ and $f(D|\theta)$ will form a conjugate pair, if $f(D|\theta)$ has the same form as $f(\theta)$

$$f(\theta|D) = \frac{f(\theta) f(D|\theta)}{\int f(\theta') f(D|\theta') d\theta'}$$

↳ posterior distribution
(goal is to estimate θ)

→ Example:

Beta distribution

$$\text{Beta}(p|\alpha, \beta)$$

α, β are fixed parameters that we don't know.

$$\text{Beta}(p|\alpha, \beta) = \frac{p^{\alpha-1} (1-p)^{\beta-1}}{\int_0^1 p^{\alpha-1} (1-p)^{\beta-1} dp}$$

$$\int_0^1 \tilde{p}^{\alpha-1} (1-\tilde{p})^{\beta-1} d\tilde{p}$$

↳ density function for p .

$$f(p) \sim \text{Beta}(p|\alpha, \beta)$$

$f(\theta|p)$, via θ
sequence of coin
tosses

$$f(\theta | p) = p^{\sum x_i} (1-p)^{n-\sum x_i}$$

number of heads

$$f(p|\theta) = \frac{f(\theta|p)}{f(p)}$$

$$\int_0^1 f(\tilde{p}) f(\theta|\tilde{p}) d\tilde{p}$$

↳ this is a

$$f(p|\theta) = c \times p^{\sum x_i} (1-p)^{n-\sum x_i} p^{(\alpha-1)} (1-p)^{\beta-1}$$

$$= c \times p^{\sum x_i + \alpha - 1} (1-p)^{n - \sum x_i + \beta - 1}$$

$$= \text{Beta}(\alpha + \sum x_i, n - \sum x_i + \beta)$$

Beta is

$$\int_0^1 \tilde{p}^{\alpha-1} (1-\tilde{p})^{\beta-1} d\tilde{p} = \frac{\Gamma(\alpha) \Gamma(\beta)}{\Gamma(\alpha+\beta)}$$

$$\Gamma(a) = \int_0^\infty e^{-x} x^{a-1} dx \rightarrow \text{solve the integration by parts}$$

$$\Gamma(a+1) = a \Gamma(a) - \text{recurrence relation}$$

$$\Gamma(n) = (n-1)!$$

↳ positive integer

" Γ will come up randomly in various problems, its weird way in which nature arranges itself."

$$E[P] = \int_0^1 p f(p) dp$$

$$= \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \int_0^1 p^\alpha (p)^{\beta-1} (1-p)^{\alpha-1} dp$$

$$= \cancel{\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}} \int_0^1 p^\alpha - p^{\alpha+\beta-1} dp$$

$$= \cancel{\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}} \left[\frac{p^\alpha}{\alpha+1} - \frac{p^{\alpha+\beta}}{\alpha+\beta+1} \right]_0^1$$

$$= \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \int_0^1 (p)^\alpha (1-p)^{\beta-1} dp$$

$$= \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \int_0^1 p^{(\alpha+1)-1} (1-p)^{\beta-1} dp$$

$$= \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \cdot \frac{\Gamma(\alpha+1) \Gamma(\beta)}{\Gamma(\alpha+\beta+1)}$$

$$= \frac{\alpha}{\alpha+\beta}$$

$$\text{prior} = \text{Beta}(\alpha, \beta)$$

↓ $\sum x_i, n - \sum x_i$

$$f(p|\theta) = \text{Beta}(\alpha + \sum x_i, \beta + n - \sum x_i)$$

What do we predict in the end?

$$E(f(p|\theta)) = \frac{\alpha + \sum x_i}{\alpha + \beta + n}$$

↳ mean of posterior

option

② MAP

↳ maximum a priori

$$\max_p f(p|\theta)$$

p

$$f'(p|\theta) = 0$$

$$p^{\alpha + \sum x_i - 1} (1-p)^{\beta + n - \sum x_i - 1}$$

⇒

$$f = k p^{\alpha-1} (1-p)^{\beta-1}$$

= k

$$\ln(f(p)) = \ln k + (\alpha-1) \ln p + (\beta-1) \ln(1-p)$$

$$\frac{d \ln(f(p))}{dp} = \frac{\alpha-1}{p} + \frac{\beta-1}{1-p} = 0$$

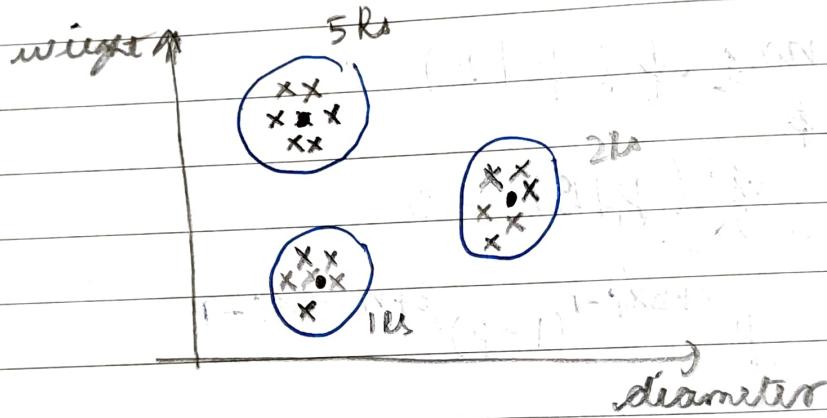
$$p(\alpha-1) + p(\beta-1) = 0$$

$$p\alpha - p + p\beta - p = 0.$$

$$P_{\max} = \frac{d-1}{d+\beta-2}$$

Unsupervised Learning (no labels are given)
K Means Algorithm

problem: you go to an ATM and it has to give out coins



Given this data can we identify which is 1Rs, 2Rs, 5Rs coin.

given features $\{x_i\}_{i=1}^n \rightarrow k$ (fixed)

Labels are not given we have to cluster the features and give them labels.

Required output:

Goal: For all i $Z_{ij} = 1$ if point i belongs to the j^{th} cluster.

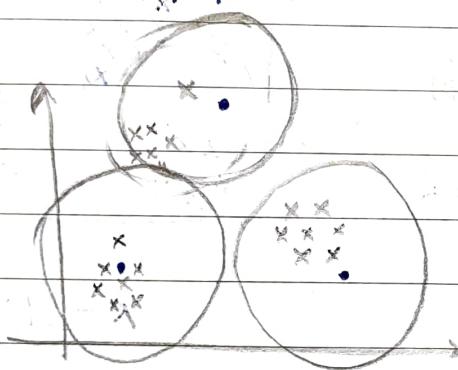
what do we wanna estimate optimise?

↳ we want close clusters

M_1^* , M_2^* , M_3^* are centroids

$$M_j^* = \frac{\sum z_{ij} X_i}{\sum z_{ij}}$$

goal : $\min \sum_{i=1}^n \sum_{j=1}^k z_{ij} \|x_i - M_j^*\|$



1/2/23

what we have as input? $\{x_i\}_{i=1}^n$, $k > 0$

Output? for every point x_i to some cluster

$z_{ij} = 1$ (if i^{th} point belongs to the j^{th} cluster)

$i \in \{1, \dots, n\}$ $j \in \{1, \dots, k\}$

Objective: $\min \sum_{i=1}^n \sum_{j=1}^k z_{ij} \|x_i - M_j^*\|$

↳ the average distance from centroids is minimum.

k-means algorithm:

1) Random initialisation η

$$M_j^{*(0)} \{ j = 1, \dots, k \}$$

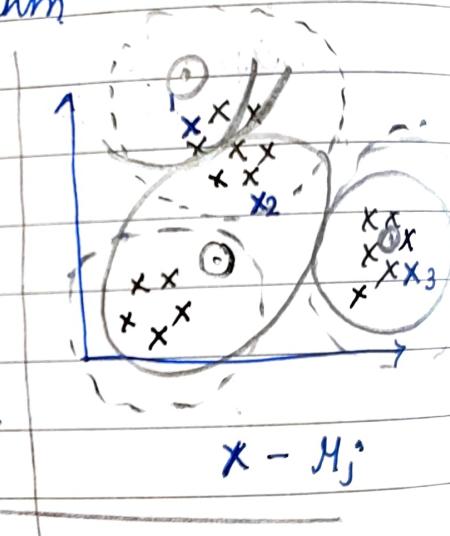
2) Repeat till convergence (t loop variable)

$$\rightarrow z_{ij}^{(t)} = 1 \text{ if } j \in \underset{m}{\operatorname{argmin}} \|x_i - M_m^{(t-1)}\| \text{ (nearest centroid)}$$

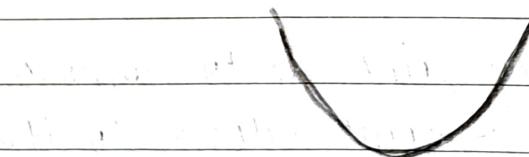
One run of this algorithm

$$\rightarrow M_j^{*(t)} = \frac{\sum_{i=1}^{n(t)} z_{ij}^{(t)} x_i}{\sum_{i=1}^{n(t)} z_{ij}^{(t)}}$$

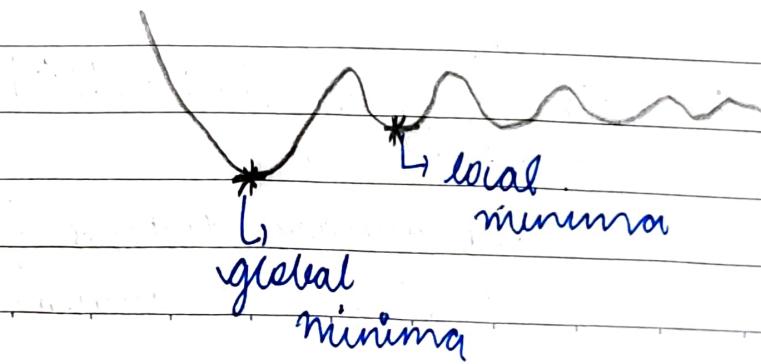
$$\forall j \in 1, 2, \dots, k$$



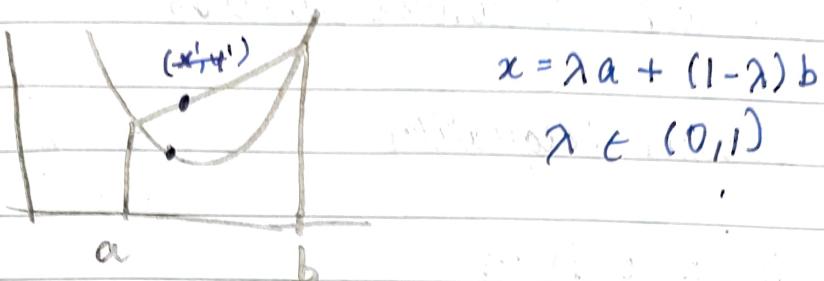
convex functions (shaped like a bowl)



general function



- Solving k-means optimally is an NP-hard.
- K-means gives local minima.



$$x' = \lambda f(a) + (1-\lambda) f(b)$$

$$f(\lambda a + (1-\lambda) b) \leq \lambda f(a) + (1-\lambda) f(b)$$

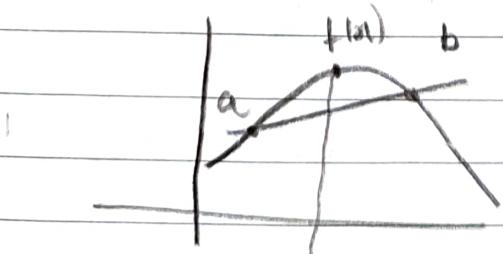
↳ one possible condition for
convexity

easier condition:

↳ $f'(x)$ is increasing

↳ $f''(x) > 0$

Concave function



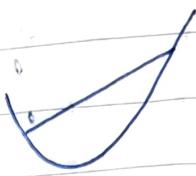
$$f(\lambda a + (1-\lambda) b) \geq \lambda f(a) + (1-\lambda) f(b)$$

$f'(x)$ is decreasing

$f''(x) < 0$.

convex functions have a global minima (concave - maxima)

lets take a convex function



Jensen's Inequality

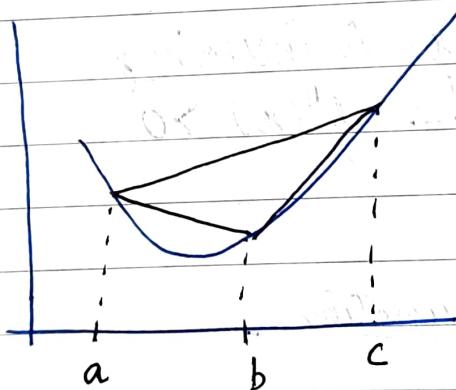
$E[f(x)] \leq$

$$f(\lambda a + (1-\lambda)b) \leq \lambda f(a) + (1-\lambda)f(b)$$

probable proba

$$E[f(x)] \geq f[E[x]]$$

How did we generalize?



$$b = \lambda_1 a + \lambda_2 b + \lambda_3 c \quad (\lambda_1 + \lambda_2 + \lambda_3 = 1)$$

$$a \leq \lambda_1 a + \lambda_2 b + \lambda_3 c \leq c$$

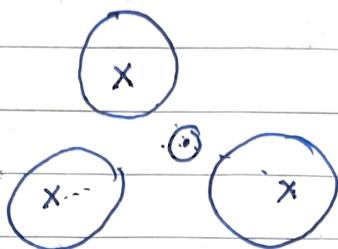
$\downarrow p.$

$$E[x] = p.$$

$$f(E(x)) \leq E[f(x)]$$

K means has a limitations:

1. Assume there are three points



- 1) circular
- 2) same size.
- ↳ ideally
in general
encourages same
size

K means algorithm:

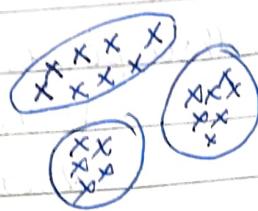
Input: $\{x_i\}_{i=1}^n$, $k > 0$

Output: cluster assignment

$z_{ij} = 1$ if i^{th} pt belongs to j^{th} cluster
 $\hat{\mu}_j^*$ \rightarrow centroid

Random init: choose $\hat{\mu}_j^{(0)}$ (or) $z_{ij}^{(0)}$ at random

Limitations of K-means algorithm

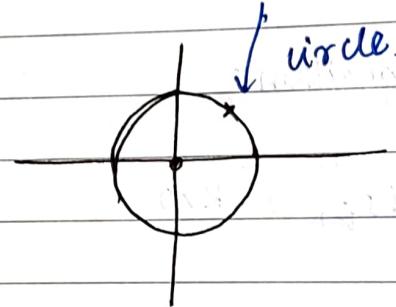


K means don't work

$$\|x_i - \mu_j^*\| = \sqrt{(x - \mu_j^*)^T (x_i - \mu_j^*)}$$

↳ vector norm tells the distance from the origin.

$$(x_i - \mu_j^*)^T I (x_i - \mu_j^*)$$



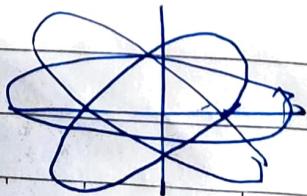
positive (semi) Definite

$$x^T A x \geq 0 \quad (> \text{ for +ve def})$$

We can define new norms

$$(x_i - \mu_j^*)^T A (x_i - \mu_j^*)$$

↳ this can have all kinds of shape.



multivariate gaussian.

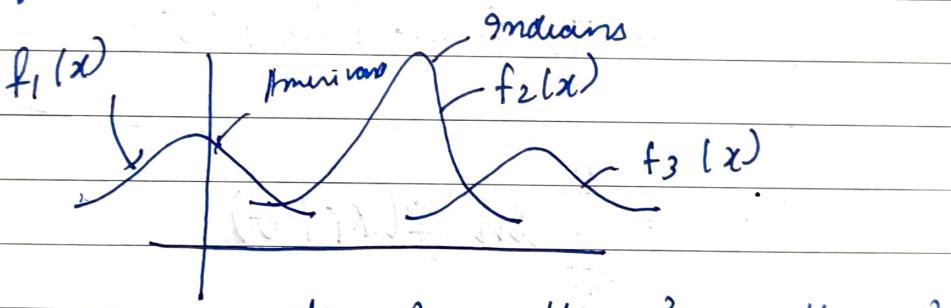
$$\frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{1/2}} \exp(-\frac{1}{2} (x - \mu_j)^T \Sigma^{-1} (x - \mu_j))$$

$$\exp\left(\frac{-1}{2} (x - \mu_j)^T \Sigma^{-1} (x - \mu_j)\right)$$

$\Sigma \rightarrow$ covariance matrix \rightarrow positive definite
 $\Sigma^{-1} \rightarrow$ is also positive definite.

We're gonna generalise k-means to Gaussian mixture models:

Single variable.



$$f(x) = \pi_1 f_1(x) + \pi_2 f_2(x) + \pi_3 f_3(x)$$

↓
probability

You are in
the first
population

$$L(\theta | D) = \prod P(X_i | \theta)$$

θ

$$\theta = \begin{bmatrix} \pi_1 & \pi_2 & \pi_3 \\ M_1 & M_2 & M_3 \\ \sigma_1^2 & \sigma_2^2 & \sigma_3^2 \end{bmatrix}$$

If we had multiple variables.

$$\boldsymbol{y}_i = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \rightarrow \text{vector}$$

$$\boldsymbol{\Sigma}_i = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \rightarrow \text{matrix}$$

$$\log(L(\theta|\omega)) = \prod f(x_i|\theta)$$

$z_i \rightarrow \text{latent variable}$

here the category is unknown.

$$f(x_i, z_i=c | \theta) = \pi_c f(x_i | \mu_c, \sigma_c^2)$$

$$l(\theta) = \sum_{i=1}^n \ln f(x_i | \theta)$$

$$= \sum_{i=1}^n \ln \sum_c f(x_i, z_i=c | \theta)$$

$$= \sum_{i=1}^n \ln \sum_c \pi_c f(x_i | \mu_c, \sigma_c^2)$$

$$\hat{\theta}_{MLE} = \underset{\tilde{\theta}}{\operatorname{argmax}} \ell(\tilde{\theta}) \quad \mid \sum \pi_c = 1$$

$$\sum \pi_c = 1$$

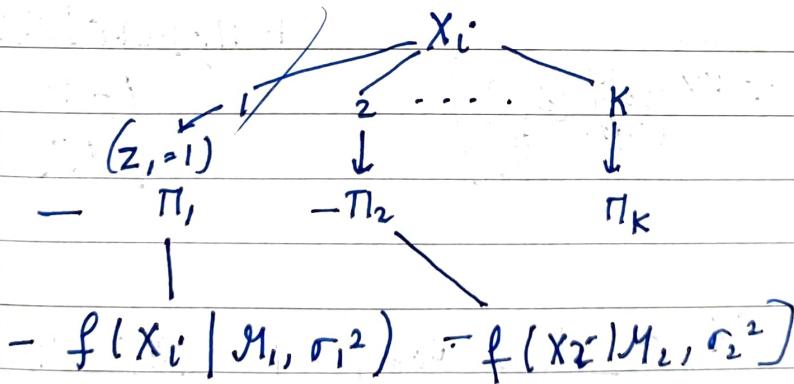
$$\alpha(\theta, \lambda) = \underset{\tilde{\theta}}{\operatorname{argmax}} \ell(\tilde{\theta}) + \lambda (1 - \sum \pi_c)$$

$$\frac{\partial \alpha(\theta, \lambda)}{\partial \mu_c} = \sum_{i=1}^n \frac{1}{\sum \pi_c f(x_i | \mu_c, \sigma_c^2)}$$

$$\propto \frac{1}{2\sigma_c^2} \delta(x_i - \mu_c)$$

06/02/23

every point x_i can come from one of the K clusters.



gaussian mixture model

? compute for MAP

MLE:

$$L(\theta) = \prod_{i=1}^n f(x_i | \theta)$$

$$\begin{aligned}
 l(\theta) &= \sum_i^m \ln f(x_i | \theta) \\
 &= \sum \ln \sum_c f(x_i, z_i=c | \theta) \\
 &= \sum \ln \sum_c \pi_c f(x_i | \mu_c, \sigma_c^2)
 \end{aligned}$$

\Rightarrow lagrangian

$$L(\theta, \lambda) = \sum \ln \sum_c \pi_c f(x_i | \mu_c, \sigma_c^2) + \lambda (1 - \sum \pi_c)$$

$$\frac{\partial L}{\partial \mu_j} = \sum_{i=1}^n \frac{\pi_j f(x_i | \mu_j, \sigma_j^2)}{\sum_c \pi_c f(x_i | \mu_c, \sigma_c^2)} \times \frac{(x_i - \mu_j)}{\sigma_j^2}$$

$$\frac{\partial L}{\partial \sigma_j} = \sum_{i=1}^n \frac{\pi_j f(x_i | \mu_j, \sigma_j^2)}{\sum_c \pi_c f(x_i | \mu_c, \sigma_c^2)} \times \left(\frac{1}{\sigma_j^3} (x_i - \mu_j)^2 \right) + \pi_j f(x_i | \mu_j, \sigma_j^2)$$

$$\boxed{\left(u'v + u''v - \frac{1}{2\sigma_j} \left(\frac{u''}{\sigma_j} \right) \right) \left(\frac{-1}{\sigma_j^3} \right) = 0}$$

$$X_i = \frac{1}{\sqrt{2\pi} \sigma_j} \exp \left(\frac{-1}{2\sigma_j^2} (x_i - \mu_j)^2 \right)$$

$$\frac{\partial L}{\partial \pi_c} = \sum_{i=1}^n \frac{f(x_i | \mu_c, \sigma_c^2)}{\sum_c \pi_c f(x_i | \mu_c, \sigma_c^2)} - \lambda = 0$$

$$\text{Put } \Gamma_{ij} = \frac{\pi_j f(x_i | \mu_j, \sigma_j^2)}{\sum_c \pi_c f(x_i | \mu_c, \sigma_c^2)}$$

$$\sum_i \frac{\Gamma_{ij} (x_i - \mu_j)}{\sigma_j^2} = 0 \quad \text{--- } ①$$

$$\mu_j = \frac{\sum \Gamma_{ij} x_i}{\sum \Gamma_{ij}}$$

$$\sum_i \Gamma_{ij} \left(\frac{1}{\sigma_j^3} (x_i - \mu_j)^2 - \frac{1}{\sigma_j} \right) = 0$$

$$\sigma_j^2 = \frac{\sum \Gamma_{ij} (x_i - \mu_j)^2}{\sum \Gamma_{ij}}$$

$$\sum_j \left(\left(\sum_i \Gamma_{ij} \right) - \lambda \Pi_j \right) = 0$$

$$\lambda = \sum_{j=1}^n \sum_i \Gamma_{ij} = n$$

$$\Pi_j = \frac{1}{n} \sum_{i=1}^n \Gamma_{ij}$$

Gaussian Mixture model

choose $\theta^{(0)}$ at random
to initialize do k-means, $\pi_c = \frac{1}{k}$.
loop $t = 1, \dots$, convergence.

$$\Gamma_{ij}^{(t+1)} = \frac{\pi_j^{(t)} f(x_i | \mu_j, \sigma_j^2)}{\sum \pi_c f(x_i | \mu_c, \sigma_c^2)}$$

→ E step

→ M step.

$$\mu_j^{(t+1)} = \frac{\sum_i \Gamma_{ij}^{(t+1)} X_{pj}}{\sum_i \Gamma_{ij}^{(t+1)}}$$

$$\sigma_j^{(t+1)} = \sqrt{\frac{\sum_i \Gamma_{ij}^{(t+1)} (x_i - \mu_j)^2}{\sum_i \Gamma_{ij}^{(t+1)}}}$$

$$\pi_j = \frac{1}{n} \sum_{i=1}^n \Gamma_{ij}$$

EM algorithm.

Γ_{ij} = probabilistic assignment

We never saw z , so now let's talk about.

Full Info Mode.

$$\sum \sum z_{ij} \ln (\pi_j f(x_i | \mu_j, \sigma_j^2))$$

M	T	W	T	F	S	S
Page No.	YOUVA					
Date:						

Full Information Model.

$$D_i = \{(x_i^j, z_i^j)\}_{j=1}^n$$

$i \dots k$

will replace z_i^j with one hot encoding

$$z_{ij} = 1 \text{ if } z_i = j$$

$$l(\theta) = \prod_{i=1}^n p(x_i, z_i | \theta)$$

$$= \prod_{j=1}^K \prod_{i=1}^n \left\{ \pi_j f(x_i | \mu_j, \sigma_j^2) \right\}^{z_{ij}}$$

$$\mu_j = \frac{\sum z_{ij} x_i}{\sum z_{ij}} \quad \sigma_j^2 = \frac{\sum z_{ij} (x_i - \mu_j)^2}{\sum z_{ij}}$$

$$\pi_j = \frac{\sum z_{ij}}{n}$$

$$\text{Input: } \{(x_i, z_i)\}_{i=1}^n$$

replace z_i with z_{ij} one hot.

$$\mu_j = \frac{\sum z_{ij} x_i}{\sum z_{ij}} \quad \sigma_j^2 = \frac{\sum z_{ij} (x_i - \mu_j)^2}{\sum z_{ij}}$$

$$\pi_j = \frac{\sum z_{ij}}{n}$$

z_{ij} = hand assignment.

$$l(\theta) = \sum \sum z_{ij} \ln (\eta_i f(x_i, y_j, \theta))$$

why does LMM work?

$$l(\theta)$$

$$\sum_i \sum_j \pi_{ij} \ln (P(x_i, z_i = c | \theta))$$

$$= \mathbb{E}_{\pi_{ij}} \ln (P(x_i, z_i = c | \theta))$$

$$\pi_{ij} = P(z_i = j | x_i, \theta^{(n)})$$

$$l(\theta) = \ln \left(\prod_i \prod_j \left(f(x_i, z_i = j | \theta) \right)^{\pi_{ij}} \right) \quad 7/2/23$$

$$l(\theta) = \sum \sum z_{ij} \ln f(x_i, z_i = j | \theta)$$

$$+ \lambda \left(\sum \pi_j - 1 \right)$$

$$\pi_j f(x_i | y_j, \sigma^2)$$

$$\pi_{ij}^{(t+1)} = \frac{\pi_j^{(t)} f(x_i | y_j, \sigma_j^{(t)})}{\sum_c \pi_c^{(t)} f(x_i | y_c, \sigma_c^{(t)})}$$

$$= \frac{f(x_i, z_i = j | \theta^{(t)})}{f(x_i | \theta^{(t)})}$$

$$= P(z_i = j | x_i, \theta^{(t)})$$

$$Q(\theta, \theta^t) = \sum_{i=1}^n \sum_{j=1}^K r_{ij}^{(t+1)} \ln f(x_i, z_i=j | \theta)$$

(E step)

$$= E_{z|x, \theta^t} [\ln f(x_i, z_i | \theta)]$$

$$\theta^{t+1} = \operatorname{argmax} (Q(\theta, \theta^t))$$

→ Exam there'll be a question with Z as a continuous function.

→ We have $P(z_i=j | x_i, \theta^{t+1})$

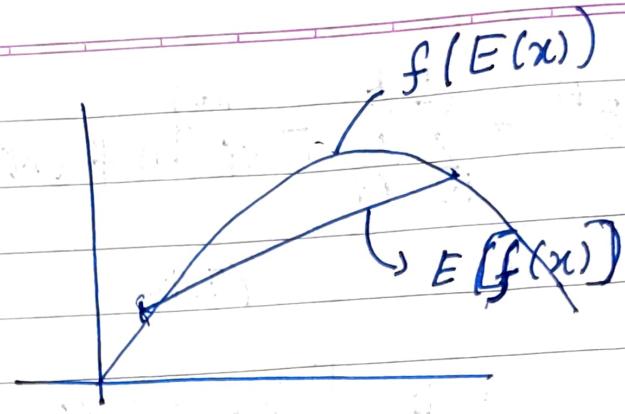
get to $f(x_i, z_i=j | \theta)$
argmax

A PROOF:

claim: EM converges to local maxima

Alternatively: Show that every step

$$l(\theta^{t+1}) \leq l(\theta^{(t+1)})$$



$$\ell(\theta^t) = \sum_{i=1}^n \ln \sum_{j=1}^k \pi_j f(x_i | \mu_j, \sigma_j^2)$$

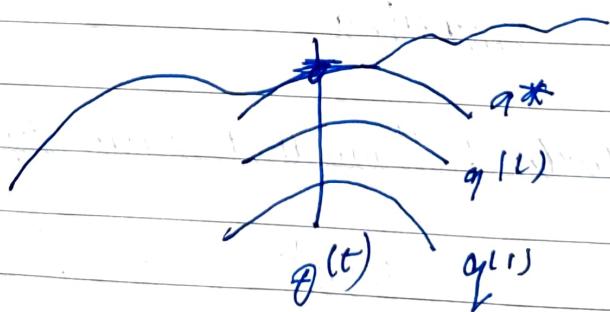
$$= \sum_{i=1}^n \ln \underbrace{\sum q(z_i=j) \pi_j}_{q(z_i=j)} \underbrace{f(x_i | \mu_j, \sigma_j^2)}_{E(x)}$$

q is any given dist.

$$\sum \sum_{i=1}^n \sum_{j=1}^k q(z_i=j) \ln \frac{p(x_i, z_i=j | \theta^{(t)})}{q(z_i=j)}$$

(Jensen's inequality)

$Q(\theta)$



$$\ell(\theta^t) = \sum \sum \frac{q(z_i=j) \ln p(x_i, z_i=j | \theta^{(t)})}{q(z_i=j)}$$

$$= \sum \ln p(x_i | \theta^{(t)}) - \sum \sum q(z_i=j)$$

$$= \sum_{i=1}^n \sum_{j=1}^n q(z_i=j) \ln p(x_i | \theta^{(t)})$$

$$- \sum \sum q(z_i=j) \frac{\ln p(x_i, z_i=j | \theta^{(t)})}{q(z_i=j)}$$

$$= \sum \sum q(z_i=j) \left(\ln p(x_i | \theta^{(t)}) - \frac{\ln p(x_i, z_i=j | \theta)}{q(z_i=j)} \right)$$

$$= \sum \sum q(z_i=j) \ln \frac{p(x_i | \theta^{(t)}) \times q(z_i=j)}{p(x_i, z_i=j | \theta^t)}$$

$$= \sum \sum q(z_i=j) \ln \frac{q(z_i=j)}{p(z_i=j | x_i, \theta_t)}$$

(KL divergence

↳ distance b/w

$q(z_i=j)$ and $p(z_i=j | x_i, \theta_t)$

$$= 0$$

$$q(z_i=g) = P(z_i=g | X_1, \theta^{(t)})$$

$$l(\theta^{(t)}) = \sum_{q_i} \ln \frac{P(x_i, z | \theta^{(t)})}{q_i}$$

$$+ \sum_k KL(q_i || P(z|x_i, \theta^{(t)})$$

given: $\begin{cases} KL(q||P) \geq 0 \\ KL(q||q) = 0 \end{cases}$

$$= Eq_i (\ln (x_i, z | \theta^{(t)})) - Eq_i (\ln q_i) \geq 0$$

$$+ \sum_k KL(q_i || P(z|x_i, \theta^{(t)}))$$

E Step:

$$= Q(\theta^t, \theta^{(t)}) + \underbrace{\text{const}}_{\sum \theta} + 0$$

M Step:

$$\theta^{t+1} = \underset{\theta}{\operatorname{argmax}} \quad Q(\theta, \theta^t)$$

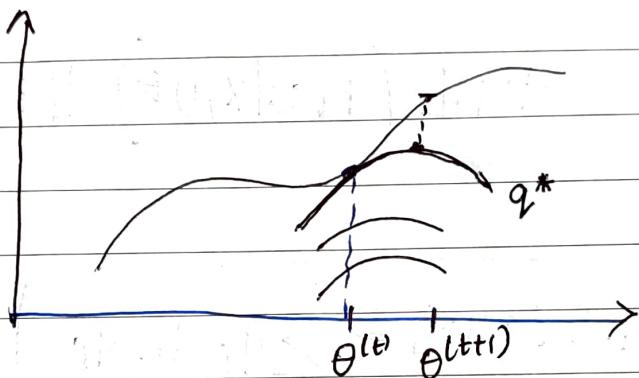
$$Q(\theta^{t+1}, \theta^t) > Q(\theta^{(t)}, \theta^{(t)})$$

$\ell(\theta^{t+1})$

9/2/23

$$\ell(\theta | \mathcal{D}) = E_{q(z)} \frac{\ln p(x, z | \theta)}{q(z)} \rightarrow \begin{array}{l} \text{expected} \\ \text{complete} \\ \text{likelihood} \end{array}$$

$$+ KL(q || p(z | x, \theta)) \rightarrow \text{entropy term.}$$



at $\theta^{(t)}$ we want KL divergence to be 0

$$\ell(\theta^{(t)} | \mathcal{D}) = E_{z|x, \theta^{(t)}} \frac{\ln (p(x, z | \theta^{(t)}))}{p(z|x, \theta^{(t)})}$$

$$Q^*(z) = p(z|x, \theta^{(t)}) + \cancel{KL} \rightarrow \text{The distance term goes to zero.}$$

$$\underbrace{E_{z|x, \theta^{(t)}} \ln p(x, z | \theta)}_{\theta^{(t+1)} \underset{\theta}{\operatorname{argmax}} Q(\theta, \theta^{(t)})} - E_{z|x, \theta^{(t)}} \ln (z|x, \theta^{(t)})$$

↓ intermediate
step for
explanation.

$$\underset{z|x, \theta^{(t)}}{E} \ln P(x, z | \theta^{t+1}) - \underset{z|x, \theta^{(t)}}{E} \ln P(z | x, \theta^{(t)})$$

maximise

constant

$$L(\theta^{t+1}) = \underset{z|x, \theta^{(t)}}{E} \ln P(x, z | \theta^{t+1}) - \underset{z|x, \theta^{(t)}}{E} \ln P(z | x, \theta^{(t)})$$

$$+ KL(P(z|x, \theta^t) || P(z|x, \theta^{t+1}))$$

because
there's a
-ve line
inequality
flips

$$KL(q||p) = \sum q_i \ln \frac{q_i}{p_i}$$

HW: MLE map, qn. Bayesian map ^{midsem}.

- 25% of midsem, \hookrightarrow read this from Bishop
- will be from theory assignments.

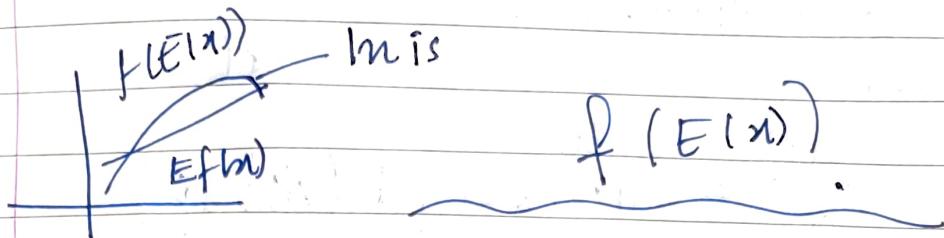
HMM

prop: $KL(q||p) \geq 0 \Rightarrow$ is positive

$$- \sum q_i (z_i=j) \underset{\mathbb{E}}{\underbrace{\ln P(z_i=j | x, \theta^t)}}_{q_i(z_i=j)}$$

Expectation of f(x)

Will prove this using Jensen's Inequality.



$$\geq \frac{\ln \sum q_i(z_i=j) \times P(z_i=j | x, \theta^t)}{q_i(z_i=j)}$$

meaux

There's a
-ve the
tensionality

There's a
-ve one $\frac{>}{=} 0$

> 0

$$\text{flips} \quad \therefore \quad \boxed{KL(q_f || p) \geq 0.} \quad \text{hence proved.}$$

Mixture of Bernoullis.

1) Coin 1

↓

$P_1 = \text{head}$

2) coin 2

1

$$P_2 = \text{head.}$$

$$3) \text{ prob of choosing coin 1} = \pi_1 \quad \pi_1 + \pi_2 = 1$$

" coin 2 = π_2

Toss the coin for d rounds.

$n_i \Rightarrow$ number of heads received.
repeat.

$$\mathcal{D} = \{n_i\}_{i=1}^n$$

$$\theta = \pi_1, \pi_2, p_1, p_2.$$

Enter

$$\Gamma_{ij}^{(t+1)} = \frac{\pi_j^{(t)} \times P(x_i | p_j)}{\sum_j \pi_j P(x_i | p_j)}$$

$$P(x_i | p_j) = \binom{d}{n} (p_j)^n (1-p_j)^{d-n}$$

~~P(X_i = n | p_j)~~ Binomial distribution

$$p_j^{(t+1)} = \sum_i \Gamma_{ij}^{(t+1)} x_i$$

$$\begin{aligned} \Gamma_{ij}^{(t+1)} &= \frac{P(z_i=j, n_i | \theta^t)}{\sum_{k=1}^2 P(z_i=k, n_i | \theta^t)} \\ &= \frac{\pi_j P(n_i | z_i=j, \theta^t)}{\sum_{k=1}^2 \pi_k P(n_i | z_i=k, \theta^t)} \end{aligned}$$

$$= \pi_j \left(\frac{d}{n} \right) p_j^{n_i} (1-p_j)^{d-n_i}$$

$$p_j^{(t+1)} = \frac{\sum_i \Gamma_{ij}^{(t+1)} n_i}{\sum_i (\Gamma_{ij}^{(t+1)}) d}$$

$$\pi_j^{(t+1)} = \frac{\sum_i \Gamma_{pj}^{(t+1)}}{d n}$$