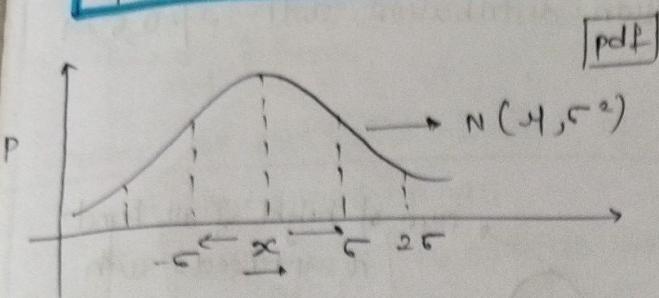


Gaussian / Normal distribution



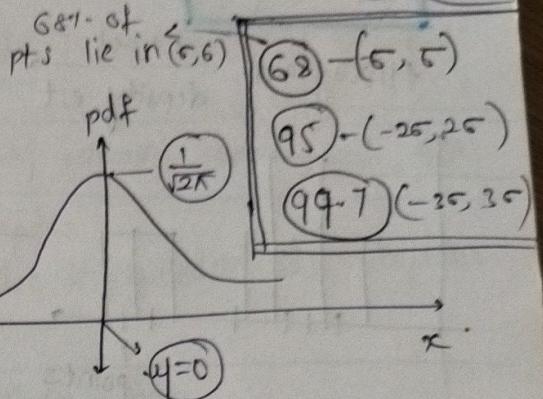
$\text{let } \mu = 0, \sigma = 1$

$p(x) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2} \cdot x^2\right\}$

$p(x) \propto \exp^{-x^2} \propto \frac{1}{e^{x^2}}$

$$p(x=k) = \frac{1}{\sqrt{2\pi} \cdot \sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

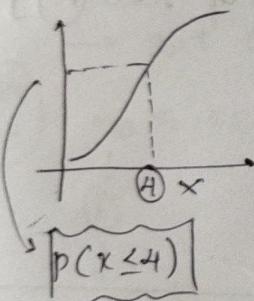
①



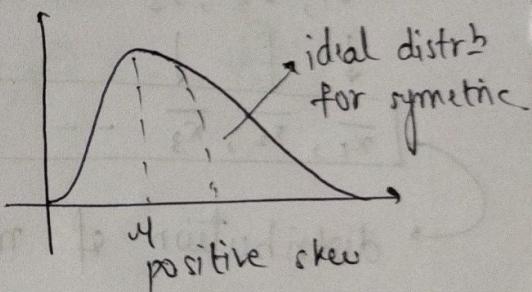
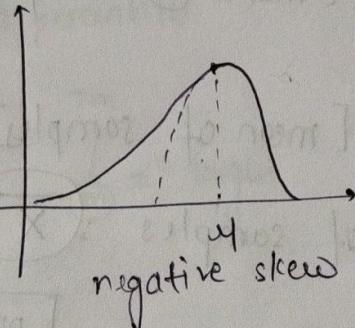
[props] ① symmetric

② as x moves away, pdf falls $\rightarrow e^{-x^2}$ → square exponential

CDF



Skewness

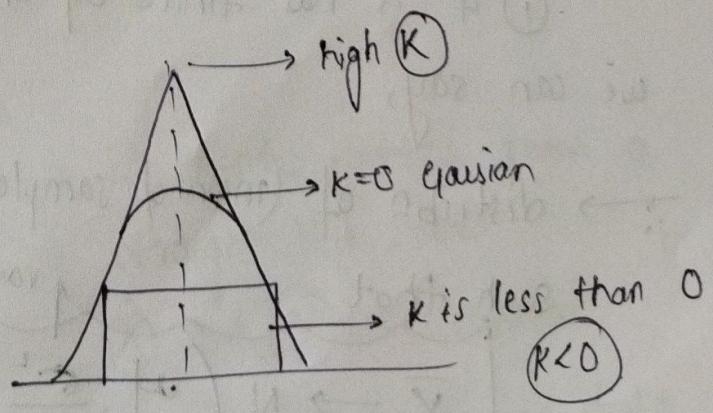


Kurtosis

→ How sharp / peaked your distribution is

$Kurt = 3$

for gaussian distrib?

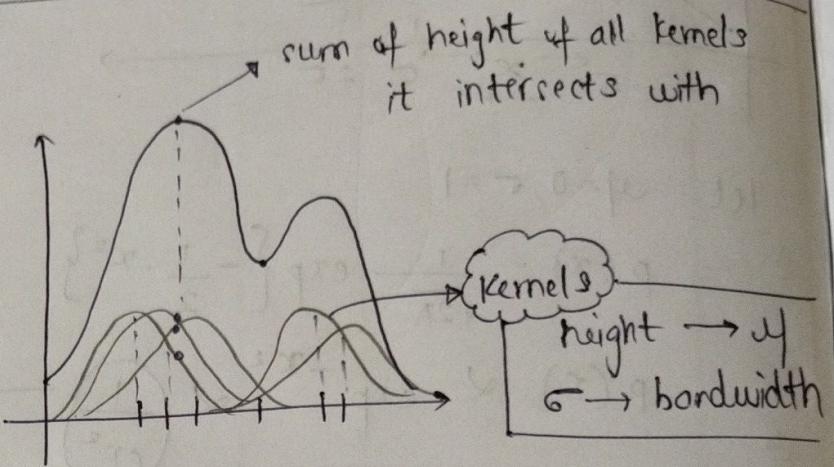
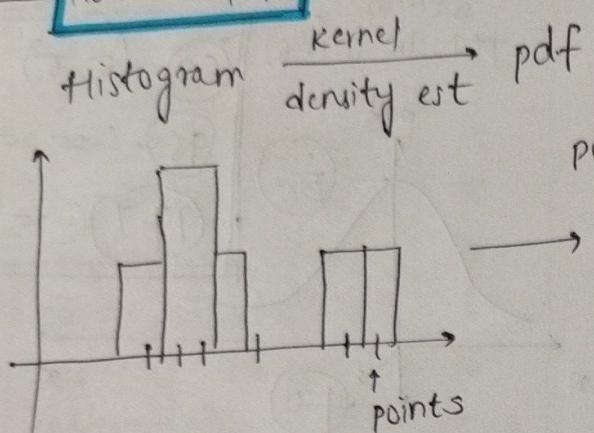


② standard normal variate

$$Z = \frac{x_i - \mu}{\sigma}$$

convert to gaussian distribution with $\mu=0, \sigma=1$

How pdf?



CENTRAL LIMIT THEOREM

* population distrn : $(X) : (\mu, \sigma^2)$ ← can be any type of distribution

* $s_1, s_2, s_3, \dots, s_m$ [make m samples each of size n]

$\downarrow \downarrow \downarrow$

$\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots, \bar{x}_m$ [mean of samples]

→ distribution of means of samples : (\bar{X})

∴ using central limit theorem,

i) if X has finite μ and σ^2

we can say,

$$\text{pop} \rightarrow X(\mu, \sigma^2)$$

after CLT

$$\text{sample} \rightarrow \bar{X}(\mu, \frac{\sigma^2}{n})$$

→ distribn of (mean of sample) follows gaussian distribution

such that

$$\bar{X} \rightarrow N\left(\mu, \frac{\sigma^2}{n}\right)$$

$n \geq 30$ then it follows gaussian

case case 1

(X) population size \rightarrow 3 million
 $\rightarrow \sigma^2$ large

may or may not be a gaussian distribution (3)

assume we plot distribution of heights

let $m = 10000$, $n = 40$

\rightarrow with just $10,000 \times 40 \rightarrow 400$ samples we can find original (4)

$\rightarrow (\bar{x}, s)$ population $\rightarrow (\bar{y}, \frac{s^2}{2})$ (4) remains same (3 mill, 40k)
distn of samples also in gaussian distribution

QQ plot Quartile Quartile plot}

\rightarrow to check if data follows gaussian distribution

steps

① X : data which is to be checked (2) Y : data which follows gaussian distn

(X) x_1, x_2, \dots, x_n

\rightarrow sort and find percentiles

$y: y_1, y_2, \dots, y_n$

\rightarrow sort and find Percentiles

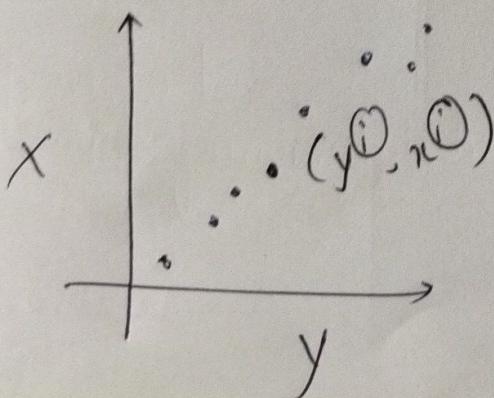
$y^{(1)}, y^{(2)}, \dots, y^{(n)}$

$x_{(1)} x_{(2)}, \dots, x_{(m)}$

\downarrow
 $x_{(1)}, x_{(2)}, \dots, x_{(m)}$

first percentile

③ plot graph of $y^{(1)}, x^{(1)}$



\rightarrow if points plotted form a straight line, then

(X) and (Y) are similar

\therefore (X) follows gaussian dist.

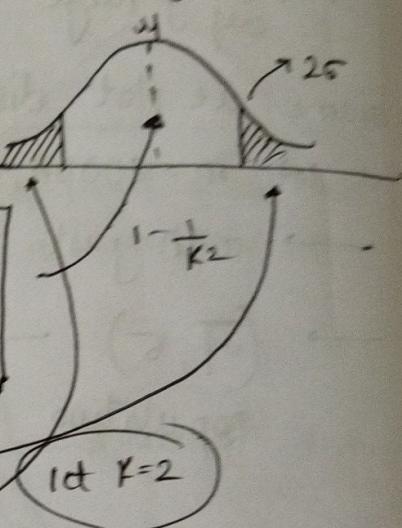
Chebyshov's Inequality

$X \rightarrow$ don't follow gaussian distribution with $\mu = 5$

$$P(|X - \mu| \geq K\sigma) \leq \frac{1}{K^2}$$

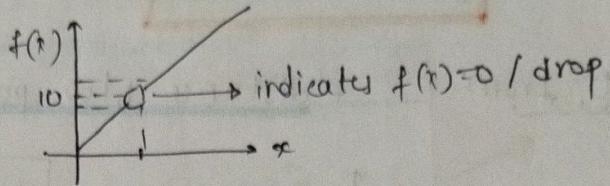
$$P(\mu + K\sigma < X < \mu + K\sigma) > \frac{1}{1 + \frac{1}{K^2}}$$

$$\begin{aligned} P(|X - \mu| \geq K\sigma) &\rightarrow X \geq \mu + K\sigma \\ &\quad X \leq \mu - K\sigma \end{aligned}$$



Limits

$$\text{let } f(x) = \begin{cases} 0, & x=1 \\ x+10, & x \neq 1 \end{cases}$$



$$\therefore \lim_{x \rightarrow 1} f(x) = 10$$

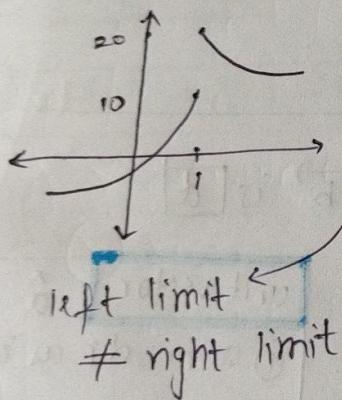
doesn't care abt $f(x)$
instead what values as you approach $x \rightarrow 1$

(defn)

$$\lim_{x \rightarrow a} f(x) = L$$

; as $-f(x)$ gets closer to a , value of $f(x)$ gets closer to L

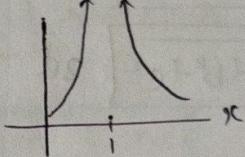
left/right limits



$\lim_{x \rightarrow 1} f(x)$ is not defined

$$\lim_{x \rightarrow 1^-} f(x) = 10$$

$$\lim_{x \rightarrow 1^+} f(x) = 20$$



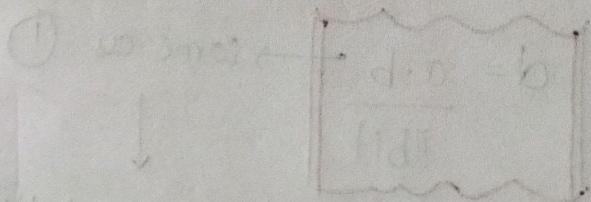
$\lim_{x \rightarrow 1} = \infty / \text{does not exist}$

$x \rightarrow 1$ from left to right

$x \rightarrow 1$ from right to left

0.2011011011 b]

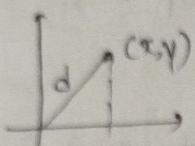
0.2011011011 = 0.2011011011, 0.2011011011



LINEAR ALGEBRA

⑥ *

* dist of pt from origin



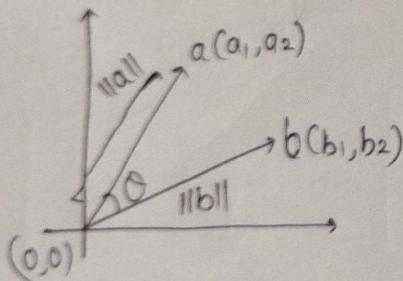
$$d = \sqrt{x^2 + y^2}$$

$$d = \sqrt{x^2 + y^2 + z^2}$$

* vectors are column vectors if not specified

dot product

$$\rightarrow \vec{a} \cdot \vec{b} = a_1 b_1 + a_2 b_2$$



$$\vec{a} \cdot \vec{b} = a_1 b_1 + a_2 b_2$$

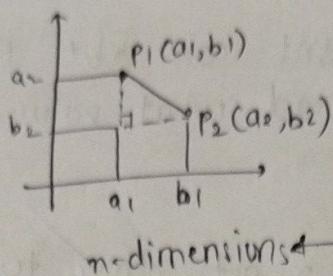
$$= \|a\| \|b\| \cos \theta$$

$$\vec{a} \cdot \vec{b} = \|a\| \|b\| \cos \theta$$

$$\theta = \cos^{-1} \left(\frac{\vec{a} \cdot \vec{b}}{\|a\| \|b\|} \right)$$

$$\text{if } a, b = 0 \Rightarrow \vec{a} \perp \vec{b}$$

* distance b/w points



$$d = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2}$$

n-dimensions

$$d = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}$$

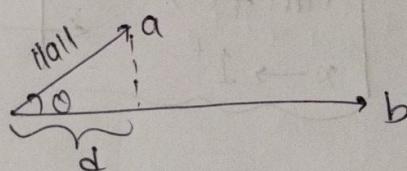
$$\rightarrow \vec{a} \cdot \vec{b} = \sum_{i=1}^n a_i b_i = \|a\| \|b\| \cos \theta$$

n-dimensional vector

$$\rightarrow \vec{a} \cdot \vec{a} = \|a\|^2$$

projection

→ projection of \vec{a} on \vec{b} is d



unit vector

- ① same dir as \vec{a}
- ② $\|a\| = 1$

$$\rightarrow \text{proj of } a \text{ on } b \rightarrow d = \frac{\|a\| \cos \theta}{\|b\|}$$

$$d = \|a\| \cos \theta \rightarrow ①$$

$$\rightarrow \text{also, } a \cdot b = \sum a_i b_i = \|a\| \|b\| \cos \theta$$

$$d = \frac{a \cdot b}{\|b\|} \rightarrow \text{same as } ①$$

$$\|\vec{x}\|_2 = \sqrt{\sum_{i \in X} x_i^2}$$

$$\frac{a \cdot b}{\|b\|} \rightarrow \frac{\|a\| \|b\| \cos \theta}{\|b\|} \downarrow \|b\| \\ \|a\| \cos \theta$$

Equation of line, plane, hyperplane

① Line - 2d

$$y = mx + c$$

$$ax + by + c = 0$$

$$y = \frac{c}{b} + \frac{a}{b}x$$

$$\rightarrow 2d : ax_1 + bx_2 + c = 0$$

$$\Rightarrow w_1x_1 + w_2x_2 + w_0 = 0$$

② plane - 3d

$$w_1x_1 + w_2x_2 + w_3x_3 + w_0 = 0$$

③ hyperplane - nd

$$\rightarrow w_0 + \sum_{i=1}^n w_i x_i \xrightarrow{\text{dot product}} w \cdot x$$

$$[w_1, w_2, \dots, w_n] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

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

$$w_0 + w^T x = 0 \rightarrow \text{Equation of plane } \pi$$

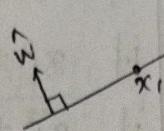
$$w \rightarrow 1 \times n$$

$$x \rightarrow n \times 1$$

if plane passes origin, $w_0 = 0$

$$w^T x = 0$$

what is \hat{w}



$$w = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

* w is unit vector \perp to π

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

$$\text{as } \theta = 90^\circ \quad \cos 90^\circ = 0$$

$$\therefore w^T x_i = 0 \quad \forall x_i \in \pi \text{ if } w \perp \pi$$

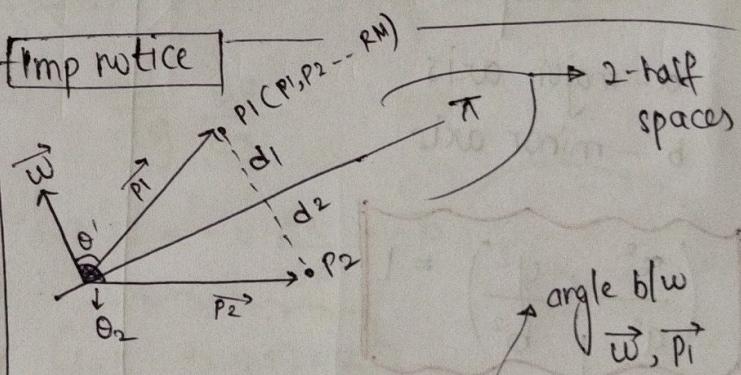
dist of pt from plane

point P (P_1, P_2, \dots, P_n)

dot product

$$d = \frac{w^T x}{\|w\|} = \frac{w^T p}{\|w\|} = \frac{w \cdot p}{\|w\|}$$

Imp notice



$$d = \frac{w \cdot p}{\|w\|} \rightarrow \frac{\|w\| \|\vec{p}\| \cos \theta}{\|w\|} = \|\vec{p}\| \cos \theta$$

$$\cos(0 < \theta) \rightarrow +ve$$

$$\cos 90^\circ = 0$$

$$\cos(\theta > 90^\circ) \rightarrow -ve$$

$$\theta_1 = < 90^\circ [90^\circ - \theta_{P_1}] \rightarrow \cos \theta \geq 0 \quad (+ve)$$

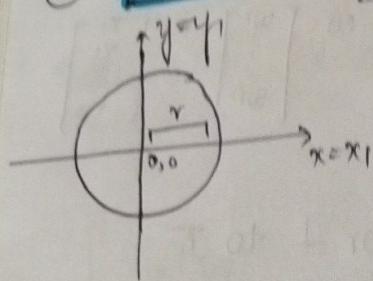
$$\theta_2 = > 90^\circ [90^\circ + \theta_{P_2}] \rightarrow \cos \theta \leq 0 \quad (-ve)$$

$$\therefore d_1 = +ve \quad d_2 = -ve$$

pts in same dir as \vec{w}
then $d = +ve$

else $d = -ve$

⑧ circle



$$x^2 + y^2 = r^2$$

origin (0,0)

if origin is (h,k)

$$(x-h)^2 + (y-k)^2 = r^2$$

→ consider point $P(x_1, x_2)$

if $x_1^2 + x_2^2 < r^2 \rightarrow$ lies inside circle

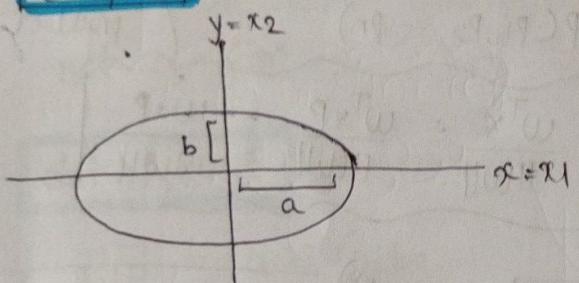
if $x_1^2 + x_2^2 > r^2 \rightarrow$ outside circle

if $x_1^2 + x_2^2 = r^2 \rightarrow$ on the circle

sphere - 3d - $x_1^2 + x_2^2 + x_3^2 = r^2$

hyper sphere - n-d - $\sum_{i=1}^n x_i^2 = r^2$

ellipse



a - major axis

b - minor axis

$$\left(\frac{x^2}{a^2} + \frac{y^2}{b^2} \right) = 1$$

$< 1 \rightarrow$ inside

$> 1 \rightarrow$ outside

$= 1 \rightarrow$ on ellipse

$$\textcircled{7} \quad \frac{d}{dx} (\sin x) = \cos x$$

$$\textcircled{8} \quad \frac{d}{dx} \cos x = -\sin x$$

$$\textcircled{9} \quad \frac{d}{dx} \tan x = \sec^2 x$$

$$\textcircled{10} \quad \sec x = \sec x \cdot \tan x$$

$$\textcircled{11} \quad \csc x = -\csc x \cdot \cot x$$

$$\textcircled{12} \quad \cot x = -\cot^2 x$$

ellipsoid n dimension

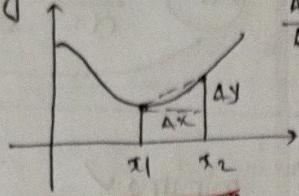
$$\sum_{i=1}^n \left(\frac{x_i^2}{a_i^2} \right) = 1$$

Differentiation

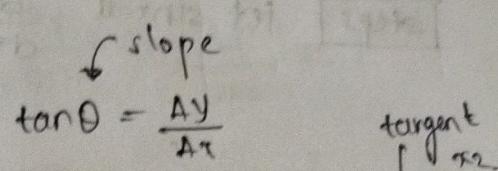
(9)

$\frac{dy}{dx}$ rate of change of y ,
as x changes

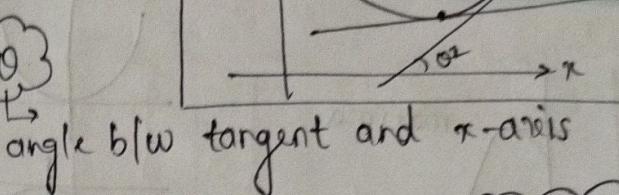
$$\rightarrow \frac{dy}{dx} = \lim_{\Delta x \rightarrow 0} \frac{\Delta y}{\Delta x}$$



$$\frac{\Delta y}{\Delta x} = \frac{y_2 - y_1}{x_2 - x_1}; \tan \theta = \frac{\Delta y}{\Delta x}$$



$$\rightarrow \frac{dy}{dx} = \text{slope of tangent to } f(x)$$

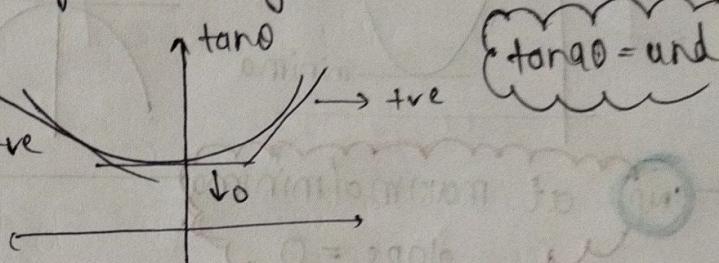


slope = $\tan \theta$

$$\textcircled{1} \quad \theta = 0, \tan \theta = 0$$

$$\textcircled{2} \quad 0 < \theta < 90^\circ \rightarrow \tan \theta > 0$$

$$\textcircled{3} \quad 180^\circ < \theta < 90^\circ \rightarrow \tan \theta < 0$$



Basic Formulae

$$\textcircled{1} \quad \frac{d}{dx}(xn) = nx^{n-1}$$

$$\textcircled{2} \quad \frac{d}{dx}(c) = 0$$

$$\textcircled{3} \quad \frac{d}{dx}(cx^n) = c \cdot nx^{n-1}$$

$$\textcircled{4} \quad \frac{d}{dx}(\log x) = \frac{1}{x}$$

$$\textcircled{5} \quad \frac{d}{dx}(e^x) = e^x$$

$$\textcircled{6} \quad \frac{d}{dx}(f(x) + g(x)) = \frac{d}{dx}f(x) + \frac{d}{dx}g(x)$$

$$\text{Ex: } g(x) = (a-bx); f(x) = x^2; f(g(x)) = (a-bx)^2$$

$$\rightarrow \frac{d}{dx}f(g(x)) = \frac{df}{dg} \cdot \frac{dg}{dx}$$

$$\rightarrow \frac{dg}{dx} = \frac{d(a-bx)}{dx} = -b \cdot \frac{dx}{dx} = -b$$

$$\frac{dg}{dx} = -b$$

$$\rightarrow \frac{df}{dg} = \frac{d}{dz} \quad \left(\begin{array}{l} \text{let } (a-bx) = z \\ \frac{df}{dz} = 2z^2 \end{array} \right) \rightarrow \frac{df}{dg} = \frac{d(z^2)}{dz} = 2z \rightarrow 2(a-bx)$$

$$\frac{df}{dg} = \frac{d(a-bx)^2}{d(a-bx)} = 2(a-bx)$$

$$\therefore \frac{dy}{dx} = (-b)(2(a-bx))$$

$$\frac{dy}{dx} = -2b(a-bx)$$

$$\frac{d}{dx}f(g(x)) = \frac{df}{dg} \cdot \frac{dg}{dx}$$

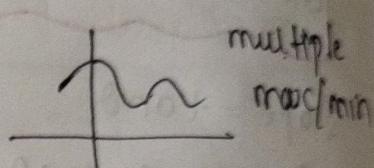
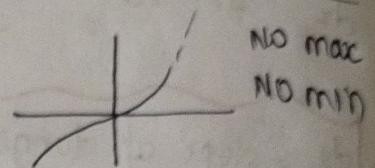
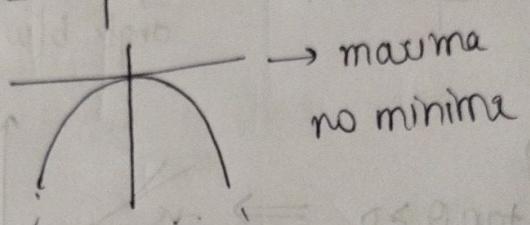
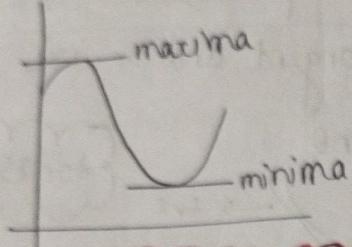
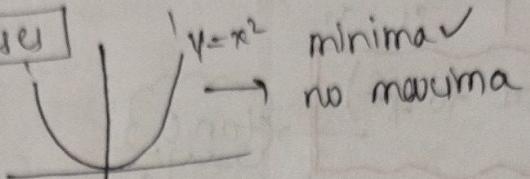
Q1) $\frac{d}{dx} (\sin x)^2 = 2 \cdot \sin x \cdot \cos x$

steps: let $\sin x = u$, $\frac{du}{dx} = \frac{d}{dx} u^2 \times \frac{du}{dx} = \frac{d(u^2)}{du} \cdot \frac{du}{dx} = 2u \cdot \frac{du}{dx}$

$$\frac{du}{dx} = \frac{\cos x}{2u} = \frac{\cos x}{\sin x} = \cot x$$

minima / maxima

case 1



at maxima/minima
slope = 0

$f(x) = x^2 - 3x + 2$

for max/min slope = 0

$$\frac{d}{dx}(x^2 - 3x + 2) = 0$$

$$2x - 3 = 0$$

$$x = \frac{3}{2} = 1.5$$

at $x = 1.5 \rightarrow \text{slope} = 0$

$$\begin{aligned} \min f(x) &= \max -f(x) \\ \max f(x) &= \min -f(x) \end{aligned}$$

maxima

minima

② $f(x) = \log(1 + e^{ax}) = y$

$$\frac{df(x)}{dx} = 0$$

$$\frac{dy}{dx} = \frac{1}{(1 + e^{ax})} \cdot e^{ax} \cdot a = 0$$

$$\frac{a \cdot e^{ax}}{1 + e^{ax}} = 0 \rightarrow a \cdot e^{ax} = 0$$

function is complex and
not as easy as previous.

gradient descend can
do it easily

$$\rightarrow f(1) = -0.25$$

$$f(1) = 0$$

as $f(1) > f(1.5)$ (somept is higher)

$\therefore f(1.5)$ is minima

minima at 1.5

Vector differentiation

(grad / dell)

$$\rightarrow x: \text{vector} \quad y = a^T x; \quad x = \langle x_1, x_2, \dots, x_d \rangle \\ a = \langle a_1, \dots, a_d \rangle \\ f(x) = y = \sum_{i=1}^n a_i x_i$$

$$a^T x = a_1 x_1 + a_2 x_2 + \dots + a_n x_n \\ \frac{\partial f}{\partial x_1} = a_1, \quad \frac{\partial f}{\partial x_2} = a_2$$

$$\frac{df}{dx} = \nabla_x f \rightarrow \text{grad / dell}$$

vector

grad / dell

$$\frac{\partial f}{\partial x_1} \rightarrow \frac{\partial f}{\partial x_1} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = a$$

vector $\in \mathbb{R}^d$

partial diff?

Example 2

$$\text{* logistic loss } L = \sum \log(1 + \exp(-y_i w^T x_i)) + \lambda w^T w$$

$x_i, y_i \rightarrow \text{constants}$

$$\frac{dL}{dw} = \nabla_w L = \frac{\exp(-y_i w^T x_i) \cdot (-y_i \cdot x_i)}{1 + \exp(-y_i w^T x_i)} + 2\lambda w = 0$$

solved by gradient descend

$$w = [w_1 \ w_2 \ w_3]$$

$$x = [x_1 \ x_2 \ \dots \ x_3]$$

$$k_{t+1} = k_{t+1} - \eta (\nabla L)_{w_{t-2}}$$

$$\begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} - \eta \begin{bmatrix} \frac{dw_1}{dw_1} \\ \frac{dw_2}{dw_2} \\ \frac{dw_3}{dw_3} \end{bmatrix} = \begin{bmatrix} w_1 - \eta dw_1 \\ w_2 - \eta dw_2 \\ w_3 - \eta dw_3 \end{bmatrix}$$

$dw_i = \frac{dL}{dw_i}$

→ run forward-prop

→ run back-prop → you get (dw_1, dw_2, dw_3) $w \cdot g = [dw_1, dw_2, dw_3]$

(12)

GRADIENT DESCENT

let $r=1$

Step 1 pick any random x_0 as initial point

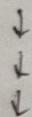
Step 2

$$x_1 = x_0 - r \left[\frac{df(x)}{dx} \right]_{x_0}$$

UPDATE FUNCTION

$$x_{i+1} = x_i - r \left[\frac{df(x)}{dx} \right]_{x_i}$$

r / step size $\frac{df}{dx}$ at x_i



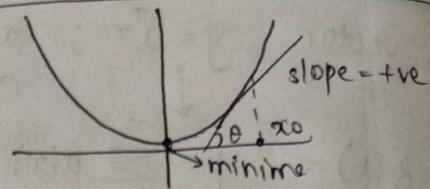
x_0, x_1, \dots, x_k

Step 4

→ if $(x_{k+1} - x_k)$ is very small,

$$\text{declare } x^* = x_k$$

case 1



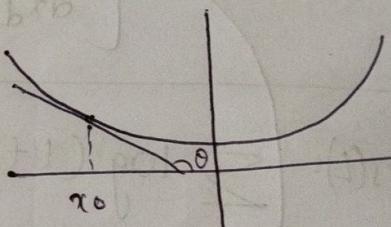
$x_0 \rightarrow \text{slope} = +ve$

$$\rightarrow x_1 = x_0 - r \left(\frac{df}{dx} \right)_{x_0} \rightarrow +ve$$

$\rightarrow x_1 = x_0 - \text{positive value}$

$\therefore x_1 < x_0$ closer to minima

case 2



$$\rightarrow \left[\frac{df}{dx} \right]_{x_0} < 0 \quad [\tan \theta < 0]$$

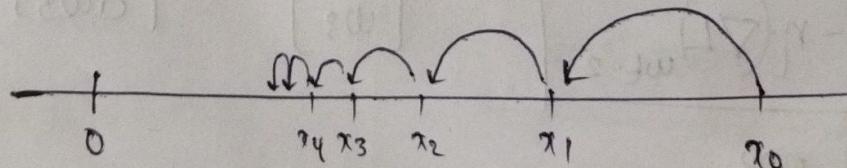
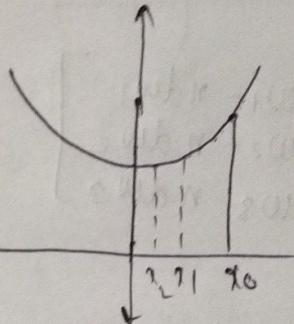
$$x_1 = x_0 - r (-ve) = x_0 + (val)$$

$\therefore x_1 > x_0$ closer to minima

Note

$$\left[\frac{df}{dx} \right]_{x_0} \geq \left[\frac{df}{dx} \right]_{x_1} \geq \left[\frac{df}{dx} \right]_{x_2} \rightarrow \text{because } 0 \rightarrow 0$$

or as $0 \rightarrow 90^\circ, \left[\frac{df}{dx} \right] \uparrow$



→ step distance decreases after each iteration

problem : oscillation \rightarrow updation may lead to same 2 value
 ex $\rightarrow 0.5, -0.5, 0.5, -0.5$ (13)

Soln :- (1) change r with each iteration
 → reduce r with each iteration

stochastic gradient descent for Linear regression

→ In 2D,

$$x_{i+1} = x_i - r \left[\frac{df}{dx} \right]_{x_i}$$

$$\text{In case of Linear regression} \rightarrow \frac{df}{dx} = \frac{df}{dw} = \sum_{i=1}^n -2x_i (y_i - w^T x_i)$$

→ you can calculate $i=1 \rightarrow n$ for all the points
 it is computationally expensive

(2) pick any K points instead of n where $K \ll n$

$$x_{i+1} = x_i - r \sum_{i=1}^K (\dots) \Rightarrow x_{i+1} = x_i - r \sum_{i=1}^K (\dots)$$

if ($K=1$) \rightarrow stochastic gradient descent

if ($K \gg 1$ and $K \ll n$)

Batch-stochastic gradient descent (BGD)

Constrained Optimization

Ex: max $\frac{1}{m} \sum_{i=1}^n (u^T x_i)^2$ given
 $u^T u = 1$

maximize $f(x)$ \rightarrow objective function
 such that

$g(x) = c$ \rightarrow constraint function
 $h(x) \leq 0$

Lagrangian Multipliers

(14)

[pb1m]

$$x^* = \max f(x) \quad \text{given} \quad g(x) = c \\ h(x) = d$$

$$L(x, \lambda, \mu) = f(x) - \lambda \{g(x) - c\} - \mu \{h(x) - d\}$$

→ we can prove that,

diff L wrt x, λ, μ

$$\frac{\partial L}{\partial x} = 0; \quad \frac{\partial L}{\partial \lambda} = 0; \quad \frac{\partial L}{\partial \mu} = 0$$

we get ③ equations → solving → $\tilde{x}, \tilde{\lambda}, \tilde{\mu}$

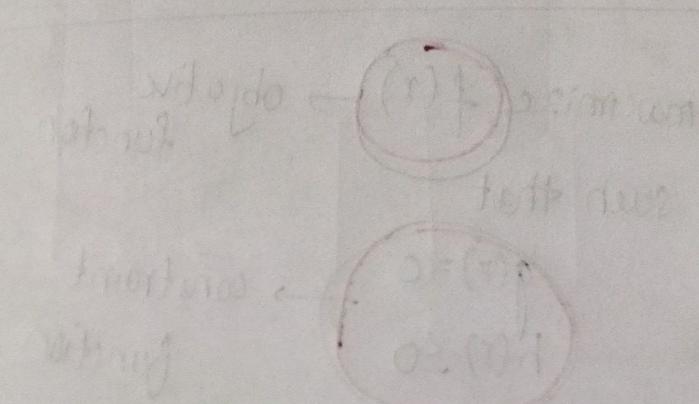
$$\therefore \tilde{x} \equiv x^*$$

Norms

$$L-p \text{ norm} = \left(\sum |x_i|^p \right)^{1/p}$$

$$L_1 \rightarrow \sum |x_i|$$

$$L_2 \rightarrow \sqrt{\sum x_i^2}$$



$$\text{auto-similarity based} \\ \text{using } f(v) \sum_{i=1}^n \frac{1}{w_i}$$

Uniform distribution

- discrete ud
- continuous ud

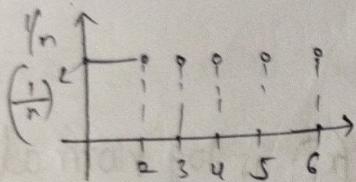
① Discrete uniform distⁿ

Ex:- coin-toss / dice roll

$$U(2,6) \rightarrow 2, 3, 4, 5, 6$$

$$\text{mean} = \frac{a+b}{2} - \text{median}$$

prob-mass function = pdf



Bernoulli's distribⁿ

→ 2 outcomes only.

Ex:- coin-toss

p = prob of outcome 1

q = prob of outcome 2

$$p+q=1$$

Binomial distribution

X = Bernoulli ($p=0.5$)

Y = Binomial (n, p)
no of trials

$$\text{pmf} = \binom{n}{k} p^k (1-p)^{n-k}$$

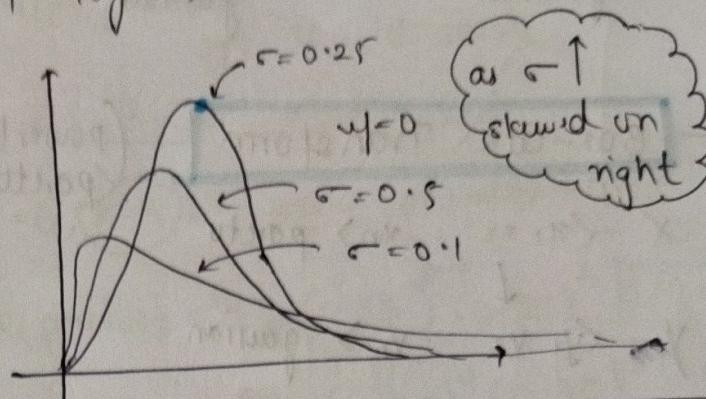
$$k = \text{successes} / (N=k)$$

Equiprobable

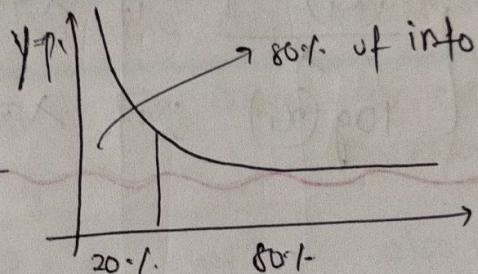
log-Normal

$$X = \log \text{Normal}(\mu, \sigma^2)$$

if $\log(X) \in \text{Normal distribution}$



Power law distribⁿ

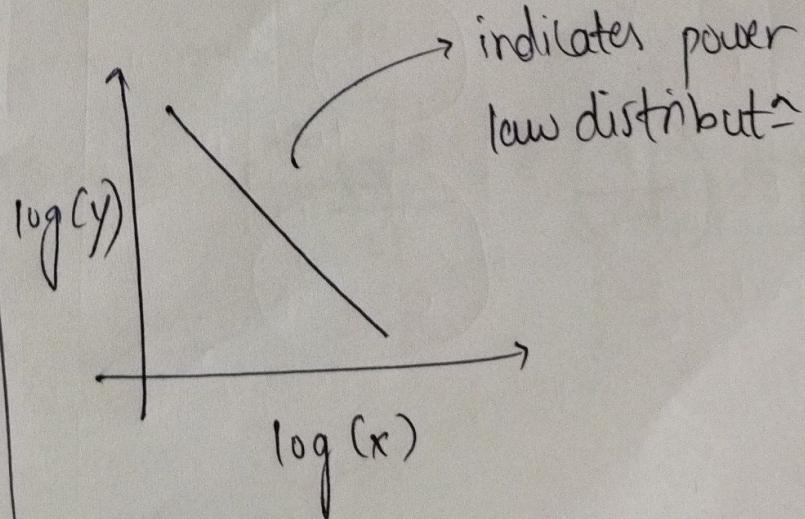


check power law distribution

$$x \rightarrow x$$

$$y \rightarrow f(x)/p(x)$$

log-log graph :-



(16) $\text{lognormal} \rightarrow \text{Gaussian}$

$$X \xrightarrow{\log(X)} Y$$

lognormal gaussian

BOX-COX Transform

$$X = \langle x_1, x_2, \dots, x_n \rangle \text{ pareto}$$



$$Y = \langle y_1, y_2, \dots, y_n \rangle \text{ gaussian}$$

(powerlaw
(pareto)

→ gaussian

step 1 $\text{BOX-COX}(X) \rightarrow \lambda$

step 2

$$y_i = \begin{cases} \frac{(x_i)^\lambda - 1}{\lambda} & ; \text{ if } \lambda \neq 0 \\ \log(x_i) & ; \text{ if } \lambda = 0 \end{cases}$$

→ means x_i is log normal

Probability & distributions

(17)

Binomial distribution

2 outcomes Yes / No

$$P(K|N) = \binom{N}{K} p^k q^{N-K}$$

N → no of trials

p → prob of outcome (K)

$$q = 1 - p$$

k → no of success/fail

Ex: in 100 ppl → 75 die

what is prob that out of 6
ppl 4 will survive

$$\rightarrow P(\text{die}) = 75/100 \quad P(\text{surv}) = 25/100$$

$$\rightarrow N=6, K=4 \leftarrow \text{survive} \quad P(\text{success}) = 0.25$$

$$\begin{aligned} \rightarrow P(\text{survive}) &= {}^6C_4 (0.75)^{6-4} (0.25)^4 \\ &= {}^6C_4 (0.25)^4 (0.75)^2 \\ &= 0.0329 \end{aligned}$$

Ex 2 die rolled 3 times

$$(i) P(5 doesn't appear)$$

$$(ii) P(\text{at least } 1, 5)$$

$$\rightarrow N=3$$

$$(i) P_S = \frac{1}{6} \text{ (success)}$$

$K=0$; no 5

$$P(\text{no 5}) = {}^3C_0 \left(\frac{1}{6}\right)^0 \left(\frac{5}{6}\right)^3$$

$$(ii) P(n=1) = {}^3C_1 \left(\frac{1}{6}\right)^1 \left(\frac{5}{6}\right)^2$$

$$\begin{aligned} P(n \geq 1) &= P(n=1) + P(n=2) \\ &\quad + P(n=3) \\ \text{at least } 1 \text{ time} \end{aligned}$$

Poisson distribution

used in time based or event period
(INTERVALS)

λ → pop size
also $\lambda = np$

$$P(x) = e^{-\lambda} \cdot \frac{\lambda^x}{x!}$$

Ex: person sells 3 cars/week.

what is probability that he sells
2 cars this week?

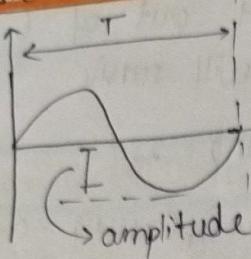
$$\lambda = 3, x = 2$$

$$P(x=2) = e^{-3} \cdot \frac{(3)^2}{2!} \cdot e^{-3} \left(\frac{9}{2}\right)$$

Fourier Transforms

(18)

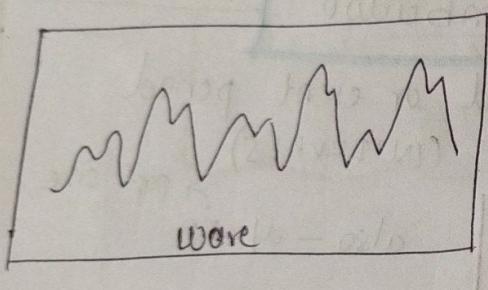
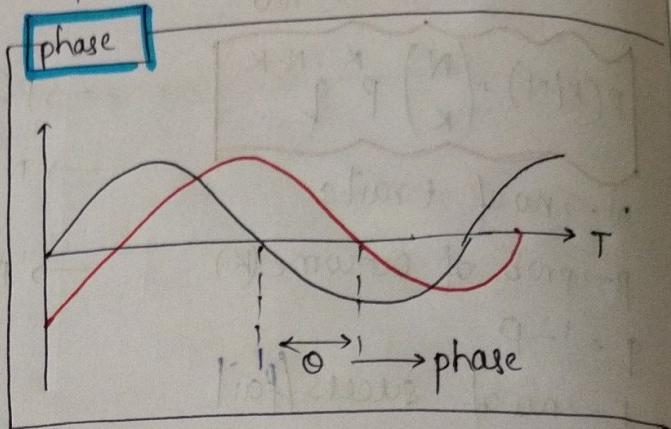
Basic terms



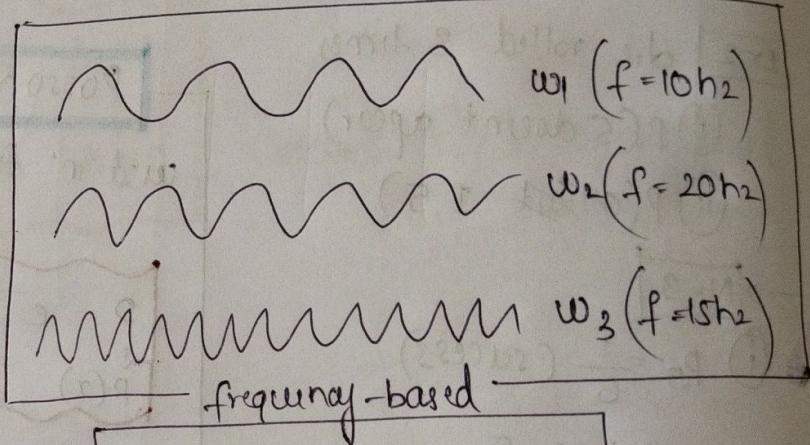
- * amplitude: height of wave
- * T : time to complete one oscillation
- * $f = \frac{1}{T}$

logic

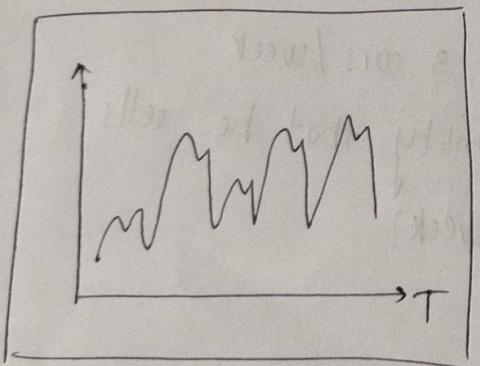
* if pattern is repeating
→ wave can be decomposed into multiple waves



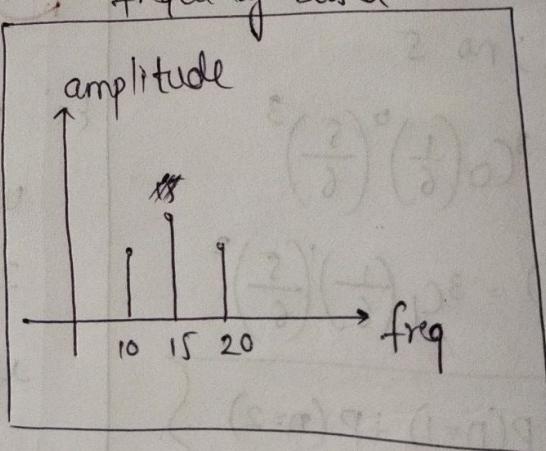
FT



time-based



Fourier transform



steps

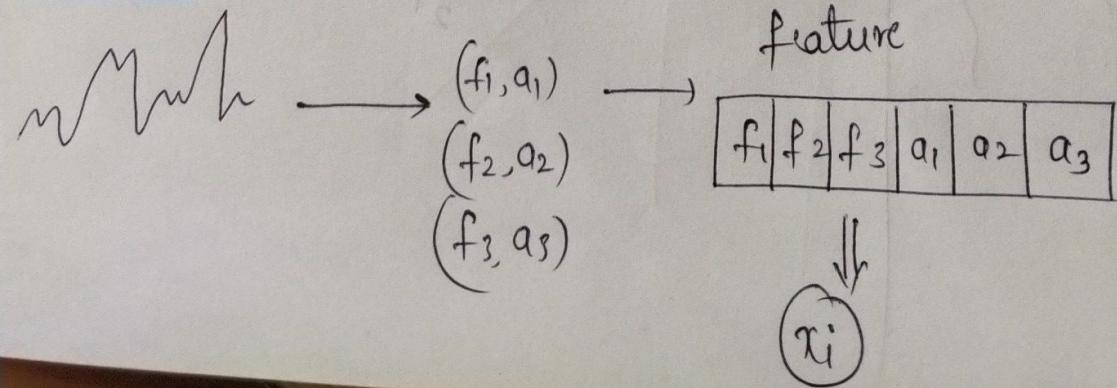
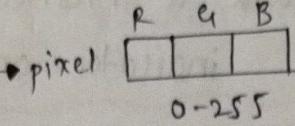
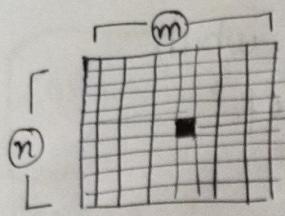


Image Data

(19)



Highly redundant info

Color (R,G,B) p

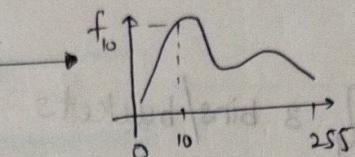
Intensity

Shape

Color Histograms

get Red values for each pixel ($n \times m$)

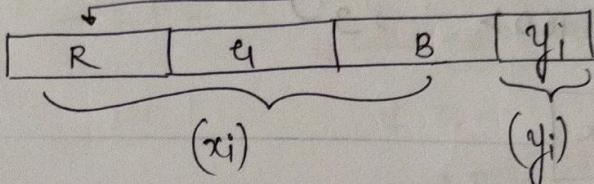
plot hist



convert to vector

0	1	10	10	255
				R

image

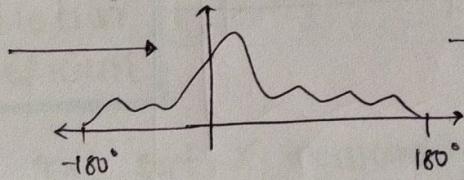


Edge Histograms

divide image to regions for each region

calculate edge-value or edge-angle

Edge \rightarrow drastic change in color \rightarrow value = location
 angle of edge \rightarrow $l = 45^\circ$ $l = 0^\circ$



convert to vector of features

used in face detection \rightarrow Haar features

object detection

SIFT (scale invariant feature transforms)

popular to detect objects in image

Key idea : stores important features and then compares
 SIFT features

main advantage :- scale independent \rightarrow to some extent rotation independent

② Feature indicator

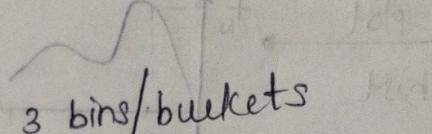
Ex: feature f_1 or height
if ($h < 150$) ret 0
else

→ add a threshold which
indicates on/off.

feature bin

→ multiple indicators

→ Ex: $h < 150 \rightarrow 1$
 $150 - 160 \rightarrow 2$
 $160 > \rightarrow 3$



Conf-matrix

actual	pred
pred	TN FN
actual	FP TP

sklearn

pred	
actual	TN FP
pred	FN TP

Covariance

$$\text{cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{cov}(X, X) = \text{var}(X) = \frac{1}{n} \sum (x_i - \bar{x})^2$$

u can't say how much related (2)

cov = +ve ↗

cov = -ve ↘

drawbacks: changes if metrics for measurement are changed

Pearson's Correlation Coef

$$P_{xy} = \frac{\text{cov}(X, Y)}{\sigma_x \cdot \sigma_y}$$

$$-1 \leq P \leq +1$$

if $P=0$,
→ no correlation

if $P=-1, +1$,
→ perfect linear relat.
→ slope of line doesn't matter

prob.m: biased towards perfect linear relationship

→ want capture non-linear rel= (sin, cos...)

Spearman rank correlation coefficient

X	Y	r_x	r_y
160	50	3	2
150	70	2	4
140	40	1	1
190	60	4	3

r_x : sort X in ascending order
→ smallest - rank 1

monotonically increasing

$P = 0.78$
spearman = 1

$$\gamma = P_{rx, ry} = \frac{\text{cov}(r_x, r_y)}{\sigma_{rx} \cdot \sigma_{ry}}$$

* it cares only if $x_1 > x_2, y_1 > y_2$ and doesn't care abt linear

* Better when outliers are present

causal models

Note

Correlation \neq causation

X correlated to Y
 \neq X causes Y

(22) confidence Interval

$X \leftarrow$ population (heights of ppl)

$\mu \leftarrow$ population mean

x_1, \dots, x_{10} sample of $n=10$, drawn from X

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (\text{sample mean})$$

* point estimate = $\mu \approx \bar{x}$ (as $n \uparrow$, $(\bar{x} - \mu) \downarrow$)

+ \boxed{CI} $\mu \in [x, y]$ with 95% probability.

↓
confidence

all these on to
find μ

if X is gaussian

$X \sim N(\mu, \sigma^2)$ → wkt 95% points in \mathbb{N} lie b/w $\mu + 2\sigma, \mu - 2\sigma$

let $\mu = 167$ $\therefore \mu = [167 - 5, 167 + 5]$ with 95% prob
 $\sigma = 2.5$

if distribution is unknown

→ $X \sim F$ (any distribution)

u know the pop-std = σ

mean of k samples
each of size n → $[(\bar{x}_1, \bar{x}_2), (\bar{x}_3, \bar{x}_4)]$

then, Central limit theorem gives : $\bar{x} = \text{sample mean} = \sum_{i=1}^n \bar{x}_i$

then, $\boxed{\bar{x} \sim N(\mu, \frac{\sigma^2}{\sqrt{n}})}$

$\mu \rightarrow$ pop-mean

$\sigma \rightarrow$ pop-std (given)

the distribution of means of sample-pop of size n
is also a gaussian distribution

then we can say,

$$M \in \left[\bar{x} - \frac{z}{\sqrt{n}} \cdot 2, \bar{x} + \frac{z}{\sqrt{n}} \cdot 2 \right] \text{ with prob of } 95.1\%$$

case 2 you don't know σ of popⁿ

→ t-distribution

$$\bar{x} \sim t(n-1)$$

degrees of freedom → as df ↑ peakness ↑

How to find confidence intervals
for σ , median, percentile? → } **Bootstrap CI**

① find 95% of CI for median

X → distribution unknown

Sample with size 'n': $\{x_1, \dots, x_n\}$ find CI of median of X ?

Only one S

$S = \{x_1, x_2, \dots, x_n\}$ → sampling with replacement (using uniform r.v.)

→ $S_1 = \{x_1, x_2, \dots, x_m\}$: random sample of size m picked from S where $m < n$

$S_2 = \{ \dots \}$

$S_K = \{x_1^K, x_2^K, \dots, x_m^K\}$ → called Bootstrap sample

Let's take $K=1000$,

→ find medians for each bootstrap sample

CI for median

$m_1, m_2, \dots, m_{1000}$

↓ sort in increasing order

$m'_1, m'_2, \dots, m'_{1000}$

$\begin{matrix} 25^{\text{val}} & | & m'_{25} & \dots & m'_{975} & | & 25^{\text{val}} \\ < m_{25} & & & & & & > m_{975} \end{matrix}$

95% CI of median of

$$X = [m'_{25}, m'_{975}]$$

(24) CI for variance

calculate br variances

$$\sigma_1, \sigma_2, \dots, \sigma_{1000}$$

↓ sort

$$G_1^{\pm}, G_2^{\pm}, \dots, G_{1000}^{\pm}$$

$$95\% = p/c$$

$$\begin{aligned} \text{lower} &= 2.5\% \\ \text{higher} &= 97.5\% \end{aligned}$$

$$\text{Power} = \frac{1-c}{2}$$

$$\text{higher} = \frac{1 + \text{cover}}{2}$$

2.50 ~~1111111111~~ 97.5
CT

qst CI for variance

$$= [\overline{6}_{25}, \overline{6}_{975}]$$

Hypothesis test

given two class student heights. Each class has 50 students

Question: Is there difference in height of C₁ compared to C₂

① choose test measure

lets take mean as measure, $M_1, M_2 \rightarrow$ mean of heights of C and C'

② null hypothesis (H_0)

H_0 : no difference in M_1 and M_2

(alternate) H_1 : there is difference
hyp

③ p-value

[Ex] Given a coin determine if it is biased

(25)

$$\text{biased} = p(\text{H|T}) \neq 0.5$$

→ Expt: flip coin 5 times

Expt obs: Head occurred 5 times ($X = \# \text{ of heads} = 5$)

after performing experiment $x \leftarrow \text{test statistic}$

$$\boxed{\text{obs} = (X = 5)} \quad \text{aftr doing exp}$$

→ H_0 : unbiased, H_1 : biased

$$\begin{aligned} \rightarrow p(\text{obs} | H_0) &= \left(\frac{1}{2}\right)\left(\frac{1}{2}\right)\left(\frac{1}{2}\right)\left(\frac{1}{2}\right)\left(\frac{1}{2}\right) \\ &= \frac{1}{2^5} = \frac{1}{32} \end{aligned}$$

$$\boxed{p(\text{obs} | H_0) = 3.1\%}$$

$$p(\text{obs} | \text{assumption}) = 3.1\%$$

if $p < 5\%$:

→ maybe assumption is wrong
→ reject assumption

as $p < 5\%$, reject H_0 ,
accept H_1

try this for different experiments and take average result

Calculate percentiles

① sort →

12	14	35	43	45	47	48	78	80	98
1	2	3	4	5	6	7	8	9	10

rank

$$\% P = \left(\frac{P}{100} \right) (n+1)$$

25.1.

$$0.25 \times 11$$

2.75

If % is decimal

$$x.y = a[x] + (a[x+1] - a[x]) \cdot (y)$$

2.75

$$\left\{ a[2] + 75\% \text{ of } d(a[2], a[3]) \right\}$$

difference

14 + (0.75) of (14 → 35)

$$14 + (0.75)(21)$$

$$= 39.75$$

$$35 - 14$$

$$\text{ans} = 78 + (0.25 * (80 - 78))$$

$$\text{Ex } 2 \quad 25.1. \rightarrow (0.75)(11) = 8.25$$

$$\text{① MSE} = \frac{1}{n} \sum (y - \hat{y})^2$$

Errors

③

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}} \rightarrow \sum (y_i - \hat{y}_i)^2$$

$$\text{② MAPE (mean abs percentile error)}$$

$$\text{MAPE} = \frac{1}{n} \sum \left| \frac{y - \hat{y}}{y} \right| \quad \text{if } y = 0 \text{ / exception}$$

$$\bar{y} = \frac{\sum y_i}{n}, \quad n\bar{y} = \sum |y_i|$$

modified-MAPE

$$\text{MM} = \frac{1}{n} \frac{|e_1| + |e_2| + \dots + |e_n|}{\bar{y}}$$

$$\text{M-MAPE} = \frac{1}{n} \sum \left(\frac{y - \hat{y}}{\bar{y}} \right)$$

mean(y)

$$\text{MN} = \frac{|e_1| + |e_2| + \dots + |e_n|}{\sum |y_i|}$$

model agnostic

Forward Feature Selection

→ time complexity
is high

- Q8 consider dataset with 10 features $f_1 - f_{10}$
- [it 1] Train with only single feature and save accuracy obtained
 $f_1 \rightarrow a_1, f_2 \rightarrow a_2 \dots$
 → pick feature with highest a
 Ex f_3 which feature along with f_3 would give us most value?
- [it 2] f_3 is fixed, now train each feature in $F - \{f_3\}$ and calculate accuracy
 Now train each feature in $F - \{f_3, f_5\}$ and calculate accuracy
 → $f_1, f_3 \rightarrow a_1$
 $f_2, f_3 \rightarrow a_2$
 $f_4, f_3 \rightarrow a_4$ } best obt fr $\{f_5, f_3\}$
- [it 3] f_5, f_3 is fixed, now pick from $F - \{f_3, f_5\}$ and continue.

pros: can be applied to any algorithm but computationally expensive

Backward Feature Selection

- [step 1] → pick all weights and calculate accuracy

- [it 1] In each it, drop one weight and check which feature dropped cause least decrease in accuracy comp to acc with all feature
- [it n] remove a feature at each iteration which causes lowest decrease in accuracy

Time Complexity

it 1 → train d models

it 2 → train $(d-1)$ models

it 3 → train $(d-2)$ models

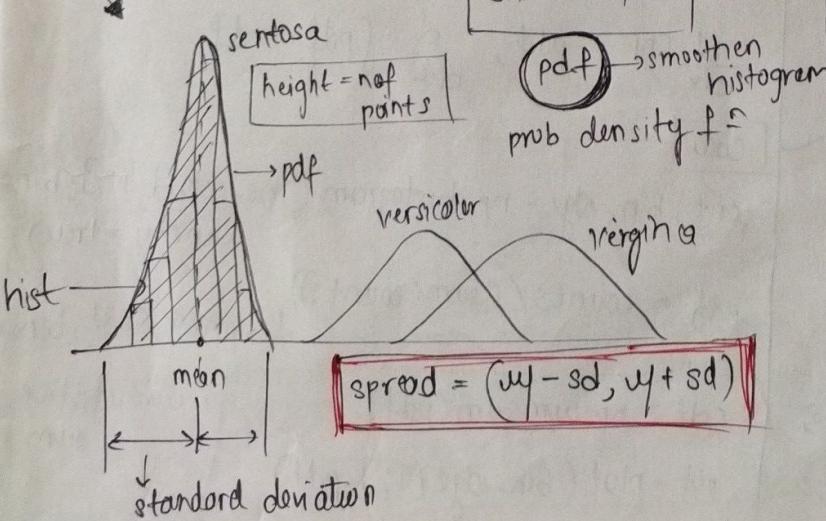
Time complexity

seaborn as sns

[IRIS Dataset]

01

- iris = pd.read_csv()
- iris["species"].value_counts()
- **2d plot**
 - * iris.plot(kind="scatter", x=" ", y=" ", field= " ")
- **3D plot (color)**
 - * sns.set_style("whitegrid")
 - * sns.FacetGrid(iris, hue="species", size=4) \
 .map(plt.scatter, 'sepal-length', 'sepal-width') \
 .add_legend();
 - * plt.show()
- **PAIR PLOT**
 - no. of plots = nC_2
 - useful upto 6 features
 - sns.pairplot(iris, hue='species', size=3)
- **UNIVARIATE ANALYSIS**
 - single variable analysis (pick 1 but feature)
 - * sns.FacetGrid(iris, hue='species') \
 .map(sns.distplot, 'petal-length')
 - ↓ histogram on petal length



EXPLORATORY DATA ANALYSES

$$\text{MEAN} \rightarrow \bar{x} = \frac{\sum x_i}{\text{len}(x_i)} = \frac{\sum x_i}{n}$$

outlier cause deviation

CSV
Header, ---
data1, data2, -

MEDIAN ① sort in ascending order
② pick middle value

some problem

VARIANCE → spread → how far from mean.

$$\text{var} = \frac{1}{n} \sum_{i=0}^n (x_i - \bar{x})^2$$

to remove -ve
-ve \bar{x} +ve
 $\sum x_i / n$, avg dist of pts from mean

$$\text{std} = \sqrt{\text{var}} = \sqrt{\frac{1}{n} \sum_{i=0}^n (x_i - \bar{x})^2}$$

average spread from mean

np.mean()
np.std()
np.median()

MEDIAN BASED

percentiles → where do you lie in sorted list

$$50^{\text{th}} \text{ percentile} = \text{arr.sort()} \quad [\text{len} = 100] \\ \text{arr}[50]$$

meaning → 49 values are smaller than it

50th percentile = median ****

Quantile

25th percentile → 1st Quantile
50th " → MEDIAN → 2nd Quantile
75th " → 3rd Quantile
100th " → 4th Quantile (max value).

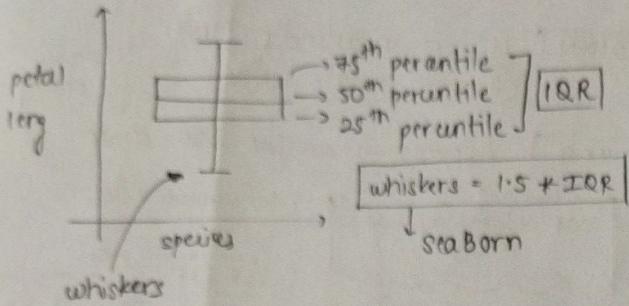
* np.percentile()

→ arr, percentile

BOXPLOTS AND WHISKERS

* `sns.boxplot(x=, y=, data="iris")`
dataset

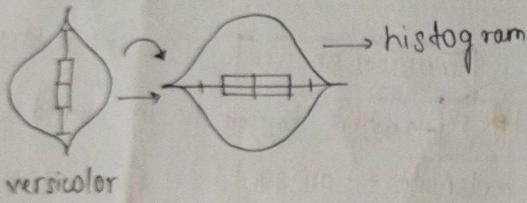
Ex → `x = species, y = petal-length`



VIOLIN PLOT

↳ boxplot + histogram

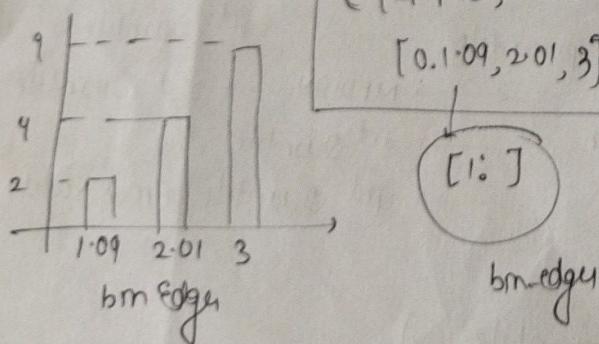
* `sns.violinplot(x=species, y=len, data=iris)`



↳ `np.histogram(arr, bins=10)`

returns

(count, bin-edge)



pdf → count / sum(counts)

contour |

`sns.jointplot(x=" ", y=" ", data=df, kind="kde")`

median
mean absolute deviation

↳ how far away from median
sd → mean

} no outlier problem

→ from statsmodels import robust

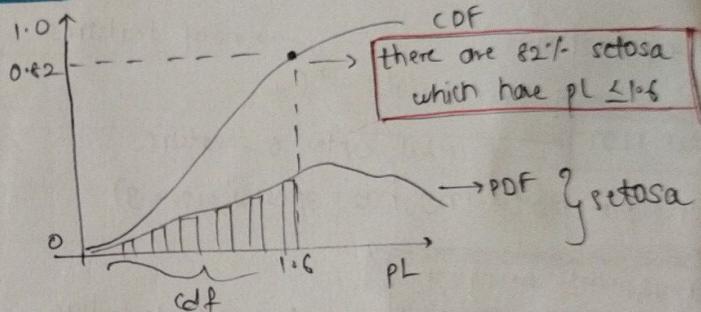
+ robust.mad()

$$MAD = \sum \sqrt{|x_i - \text{median}|} / n$$

inter-Quartile range

75th percentile - 25th percentile = IQR

CDF cumulative distribution function



calculation

→ $x = val$

→ $y = \frac{\text{how many pts have } val \leq x}{\text{total pts}} = \frac{41}{50} = 0.82$

or

count probabilities upto curr point

CDF = area under pdf until that point

$$\frac{d}{dx} \text{CDF} = \text{pdf}, \quad \int \text{pdf} = \text{CDF}$$

Code

cnt, bin-edge = `np.histogram(iris['petal length'], bins=10, density=True)`

pdf = `counts / (sum(counts))`

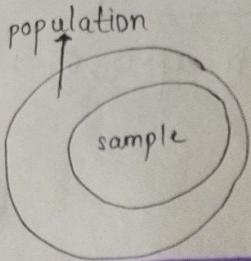
↳ `print(pdf, bin-edge)`

$$\text{cdf} = \text{np.cumsum}(pdf)$$

`plt.plot(bin-edge[1:], pdf, cdf)`

} cumulative sum

Mathematics for ML



Bell Curve

Gaussian/Normal distribution

$\mu \rightarrow$ height, $\sigma^2 \rightarrow$ width
reduced width reduced height.

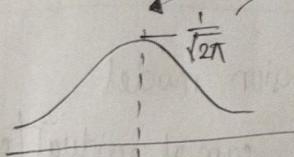
peak at μ

X is Normal distr $\rightarrow X \sim N(\mu, \sigma^2)$

$$p(x) = \frac{1}{\sqrt{2\pi} \cdot \sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

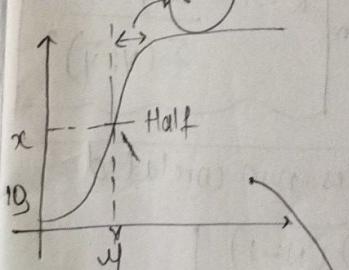
let $\mu=0, \sigma^2=1$

$$p(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} = \left(e^{-x^2} \right) \begin{matrix} (-2)^2 = 4 \\ (-1)^2 = 1 \\ (0)^2 = 0 \\ (1)^2 = 1 \\ (2)^2 = 4 \end{matrix}$$



$$p(x) = e^{-x^2} = \frac{1}{e^{x^2}} = y$$

CDF

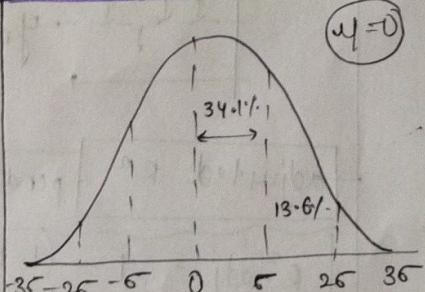


① smaller variance, less apart from mean

probability that val $\leq \mu$ is $\frac{1}{2}$

68-95-99.7 rule

S.D



b/w $-6 - 6 \rightarrow 68\%$ pts
 $-26 - +26 \rightarrow 95.0\%$ pts
 $-36 to 36 \rightarrow 99.7\%$ pts

Covariance

How elements r related?

splits ds

$$\text{hb}-y = \text{hb}[\text{rb}[r]] \text{ [Q-3 yes]}$$

$$\text{cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{cov}(x, x) = \text{var}(x)$$

Chebyshov's Inequality

if $y \not\sim$ Gaussian distribution

$$P(\mu - K\sigma \leq y \leq \mu + K\sigma) \geq 1 - \frac{1}{K^2}$$

Central Limit Theorem

$X \not\sim \text{ID } (\mu, \sigma^2)$

$n \geq 30$

$$\begin{aligned} S_1 &\rightarrow x_1 + x_2 + \dots + x_{30} = \bar{x}_1 \\ S_2 &\rightarrow x_1 + x_2 + \dots + x_{30} = \bar{x}_2 \\ &\vdots \\ S_{100} &\rightarrow x_1 + x_2 + \dots + x_{30} = \bar{x}_{100} \end{aligned} \rightarrow \text{mean} = \bar{X}$$

$$\bar{X} \approx \text{ID } \left(\mu, \frac{\sigma^2}{n} \right)$$

plotting means $\bar{x}_1 - \bar{x}_n$

$$\bar{x} \approx \mu$$

$$\sigma \approx \sigma/\sqrt{n}$$

Random variable

Categorical

Numerical

discrete

Continuous

$$\mu - \sigma \quad \mu + \sigma$$

68% have some quan b/w $-\sigma$ to σ

04 PEARSON CORRELATION COEFFICIENT

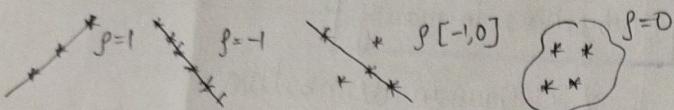
$$\text{cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \rightarrow \text{not quantified}$$

$$PCC = \rho(x, y) = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

how much +ve or -ve

$-1 \leq \rho \leq 1$

\hookrightarrow direction of cov
 \hookrightarrow strength



- if $\rho(x, N) = 1$, u can drop one

SPEARMAN'S RANK CORRELATION

pearson's only focuses on linear aspects

$$r_s = 1 - \frac{6 \cdot \sum d_i^2}{n(n^2 - 1)}$$

no of obs

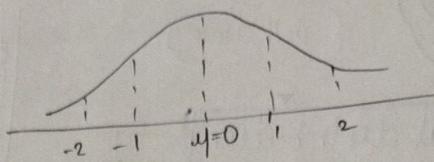
① sort first col, X.

② $r_A, r_B \rightarrow$ rank aw to ascending order

③ $d_i = r_A - r_B$

standard Normal distribution

Gaussian distribution with $\mu=0, \sigma=1$



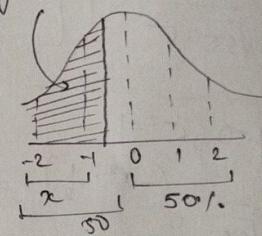
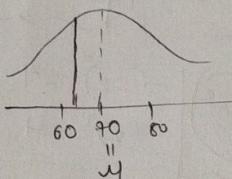
Z-score

Gaussian \rightarrow standard Normal

$$Z\text{-score} = \frac{x_i - \mu}{\sigma}$$

$$P(X > 65) = 50 + (50 - x)$$

get from z-table

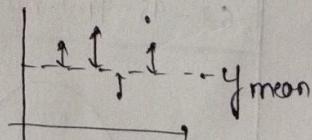


R-SQUARE SCORE

Accuracy of linear regression model

$$R^2 = 1 - \frac{SS_{\text{Res}}}{SS_{\text{tot}}} \quad \begin{array}{l} \text{sum of residual/error} \\ \text{diff of predicted,} \\ \text{original pts} \end{array}$$

$SS_{\text{tot}} \rightarrow$ diff b/w points and y-mean line



$$R^2 = 1 - \frac{\sum (y_i - \bar{y})^2}{\sum (y_i - \bar{y})^2}$$

adjusted R^2

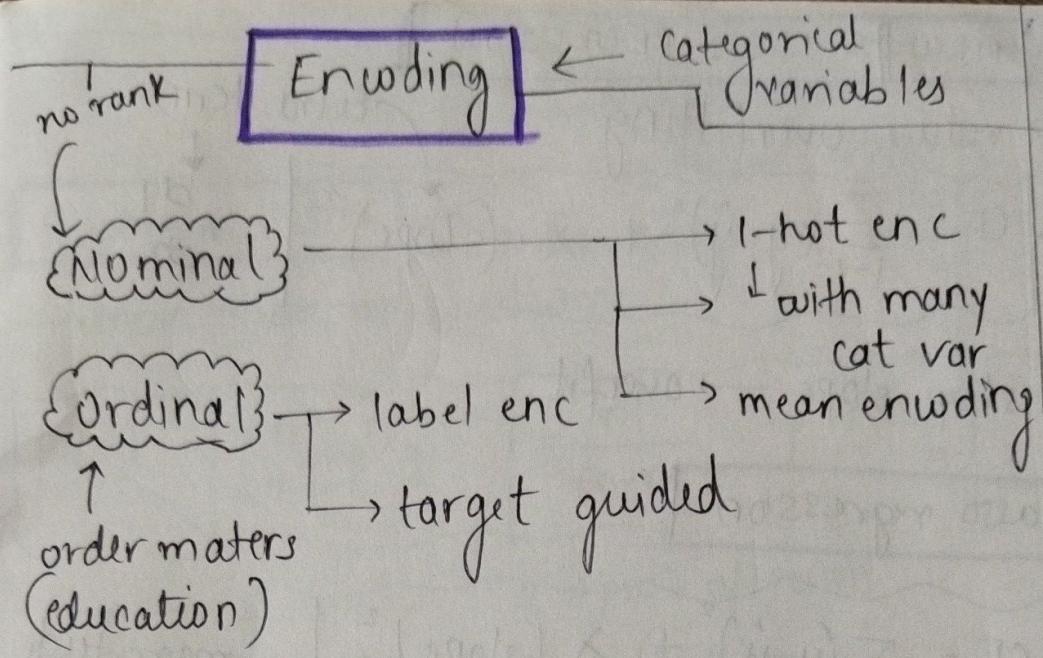
penalizes run correlated

$$R^2_{\text{adj}} = 1 - \frac{(1-R^2)(N-1)}{N-p-1}$$

$p = \text{no of predictor s} - \times$

$N = \text{sample size}$

$P \uparrow \cdot R^2 \downarrow$



BIAIS & VARIANCE

underfitting → error with training data

overfitting → error for test ↑ - less for test

BIAIS → error of training data

VARIANCE → error for ~~training~~^{test} data

generalized model → low variance
low bias

→ Bias - Variance - trade off (manage both)

error

Probability Density functions

Normal	Bernoulli
Students T	uniform
Binomial	Poissons

07

Binomial dist \rightarrow bi → 2 values
 Independent val
 more than one outcome

$$P(x) = \binom{n}{x} p^x q^{n-x} = \frac{n!}{(n-x)!x!} \cdot p^x (1-p)^{n-x}$$

n = no of trials
 x = successes
 no of ↑

p → prob of success

q → (1-p)

Ex: die rolled three, prob of getting 3 twice
 $n=3, p=1/6, q=5/6, x=2$

Poissons distributions

$$p(x) = e^{-\mu} \cdot \frac{\mu^x}{x!}$$

μ → average / mean
 for interval

prob of events in a time period

Normal / Gaussian

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

$\rightarrow \mu = 0, \sigma = 1$ → standard normal distrib

→ Normal $\rightarrow z = \frac{x-\mu}{\sigma}$, standard ND

$$z = \frac{x-\mu}{\sigma}$$

z-statistics → probability in ranges
 $\rightarrow \mu = 10, \sigma = 5$

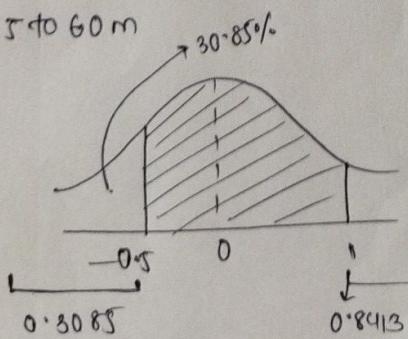
z-table

\rightarrow df → prob of work done in x min $\rightarrow [1 \dots n]$

$\rightarrow p(\text{work done in } n) \text{ from 60m}$

$$\frac{45-50}{10} \leftrightarrow \frac{60-50}{10}$$

$$-0.5 \leq z \leq 1$$



$$0.8413 - 0.3085 \rightarrow 0.5329 \quad (53.29\%)$$

z-score estimation

df = profits for 12 year $\rightarrow n=12, \mu=10, \sigma=40$

profit 13th year = ?

$$\mu \pm z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

$\alpha \rightarrow 1 - (\text{lv of confidence})$
 (lv of significance)

Ex: confidence = 90%, $\alpha = 1 - 0.9 = 0.1/2 = 0.05$

$$\rightarrow z_{0.05} = 1.621 \quad \text{from z-table}$$

$$\rightarrow a \leq x \leq b$$

\therefore profit in b/w (a,b) with 90% confi

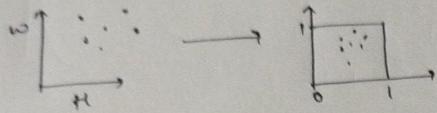
data Normalization - column normalization

feature

$$[col-a] \quad x_i = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}}$$

$$x_i = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}} \quad x_i \rightarrow 0 \leftrightarrow 1$$

gets rid of units
restricts data to grid

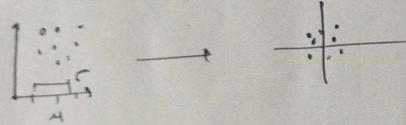


mean vector = $[f_1, f_2, f_3, \dots, f_n]$
central data point

column standardization

$$a_1, a_2, \dots, a_n \rightarrow a_1^*, a_2^*, a_3^*, \dots, a_n^* \quad [\sigma = 1, \mu = 0]$$

$$z\text{-score} = a_i^* = \frac{a_i - \mu}{\sigma}$$



covariance

$$\begin{bmatrix} f_1 & f_2 & \dots & f_n \\ x_1 & x_{12} & \dots & x_{1n} \\ x_2 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & & & & \vdots \\ x_n & x_{n1} & \dots & x_{nn} \end{bmatrix}$$

$$\text{cov}(f_1, f_1) = \text{var}(f_1)$$

$$\text{cov}(f_1, f_2) = \text{cov}(f_2, f_1)$$

$$\text{cov}(f_1, f_2) = \frac{1}{n} \sum (f_{1i} - \mu_1)(f_{2i} - \mu_2) \quad \rightarrow \mathcal{M}(f_1) \quad \rightarrow \mathcal{M}(f_2)$$

if data is col standardized

$$\text{cov}(f_1, f_2) = \frac{1}{n} \sum (x_{1i} \cdot x_{2i}) \quad \downarrow \quad \frac{1}{n} (f_1^T \cdot f_2)$$

$$\therefore \text{CVM} = \frac{X^T \cdot X}{d \times d} \quad \text{if columns are standardized}$$

$$S = X_{n \times d} \cdot X_{d \times n}^T \rightarrow \text{CVM}_{d \times d}$$

$$a \cdot b = a^T b$$

11

PCA

$f_{1,2} \Rightarrow u_1$ $u_1 \rightarrow$ unit vector

$x_i' \rightarrow$ projection of x_i on u_1

$$x_i' = \text{proj}_{u_1} x_i = \frac{u_1 \cdot x_i}{\|u_1\|} = u_1^T x_i$$

\rightarrow find u_1 st $\text{var}(\text{proj}_{u_1} x_i)$ is max

$$\text{var}(x_i') = \frac{1}{n} \sum_{i=1}^n (u_1^T x_i - u_1^T \bar{x}_i)^2 \rightarrow 0 (\text{std})$$

$$\therefore \text{PCA} \rightarrow \text{maximize } \frac{1}{n} \sum u_1^T x_i^2$$

VARIANCE MAXIMIZING APPROACH

Eigen val and vectors

$S_{d \times d} \rightarrow$ Eigen values $\lambda_1, \lambda_2, \dots, \lambda_d$
vectors v_1, v_2, \dots, v_d

$$v_i \perp v_j \quad \downarrow \quad v_i \cdot v_j = 0 \quad \rightarrow \lambda_i v_i = S \cdot v_i, \lambda_2 v_2 = S \cdot v_2, \dots$$

finally $u_1 = v_1 \rightarrow$ largest eigen vector

consider, $d=2$

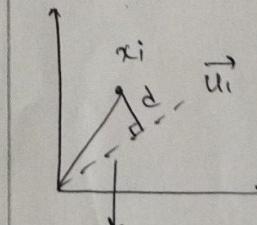
$$\begin{cases} \lambda_1 = 3 \\ \lambda_2 = 0 \end{cases} \rightarrow \lambda_1 \text{ give 100% variance}$$

$$\begin{array}{l} \lambda_1 = 3 \\ \lambda_2 = 0 \\ \therefore \frac{\lambda_1}{\lambda_1 + \lambda_2} = \frac{3}{7} \end{array} \quad \uparrow \text{already proved}$$

PCA - distance minimization

$$\text{minimize } \sum_{i=1}^n d_i^2$$

$$d_i^2 = \|x_i\|^2 - u_1^T x_i^2$$



$$d_i^2 = x_i^T x_i - (u_1^T x_i)^2$$

$$\min \sum_{i=1}^n (x_i^T x_i - u_1^T x_i)^2$$

$$\text{proj}_{u_1} x_i = u_1^T x_i$$

$$\hookrightarrow \text{pythagoras} \rightarrow \|x_i\|^2 = d_i^2 + (u_1^T x_i)^2$$