

CS215: part-2:-

• Topics for the Sem:

1) Transformation of random variables

- univariate

- multivariate

2) Multivariate Statistics

- Multivariate Gaussian

- functional form, properties

- Relationship with Linear algebra

- singular value decomposition (SVD), of matrices.

- Eigen decomposition of matrices.

3) Principle Component Analysis (PCA)

- Relies on covariance of multiple random variables.

- need not be multivariate gaussian.

- Special properties when data is multivariate Gaussian.

4) Bayesian Statistics

- improving over MLE estimation.

- prior models

- Some concepts from information theory

5) Measuring quality of estimators
efficiency.

- Fischer information

- Cramer-Rao lower bound

- Bayesian Cramer-Rao lower bound.

2 assignments

1 Quiz

1 endsem.

D Transformation of random variables:-

$X \rightarrow$ random variable; $p(x)$ is the PDF of X .

$g(\cdot)$ is the transformation function!

$$Y = g(X)$$

we need $q(y)$; the PDF of Y .

- take $g(x)$ to strictly increasing; for now.
(we will generalize later).

Principle of probability mass conservation:-

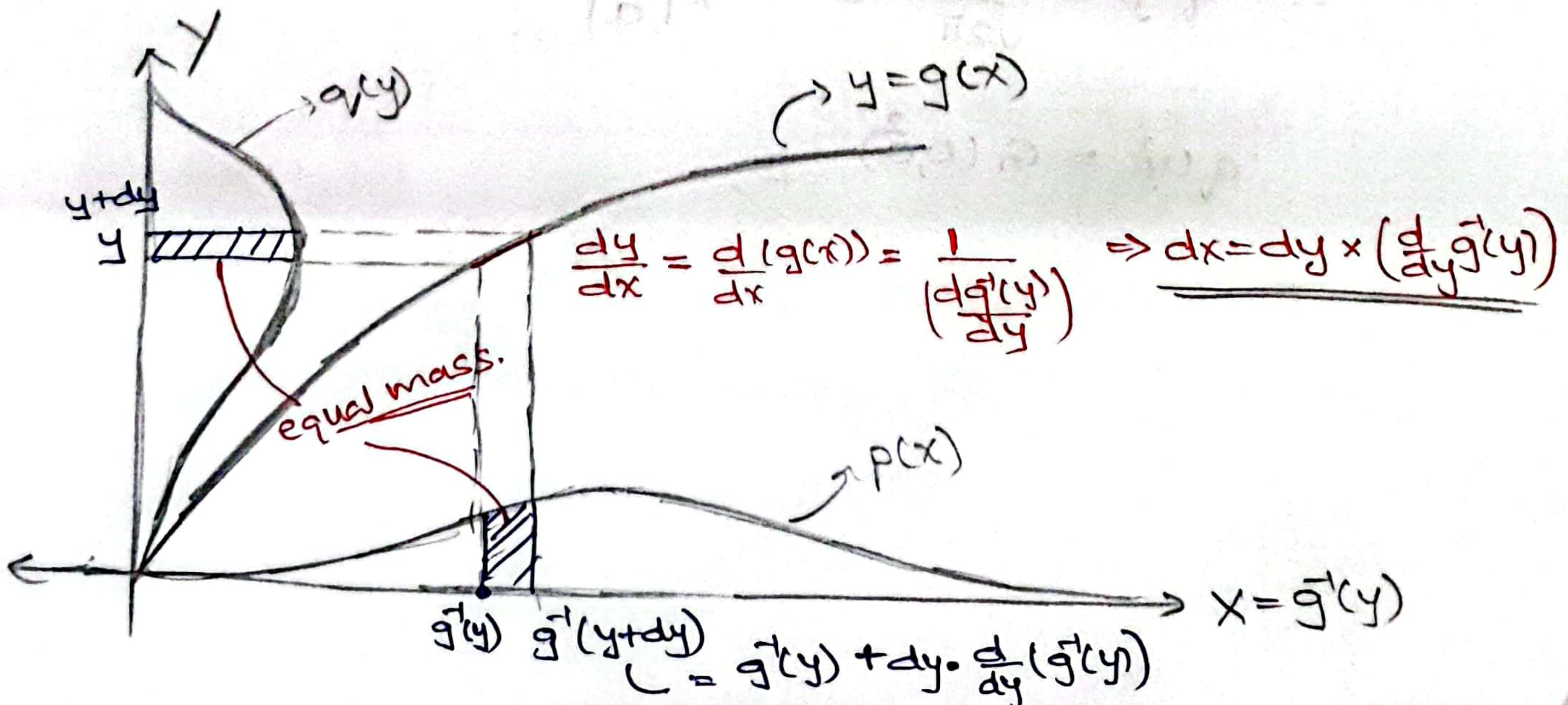
$$P(a < X < b) = P(g(a) < Y < g(b))$$

$$\int_a^b p(x) dx = \int_{g(a)}^{g(b)} q(y) dy$$

$$\int_{g(a)}^{g(b)} p(\bar{g}(y)) \cdot \frac{d}{dy}(\bar{g}(y)) \cdot dy = \int_{g(a)}^{g(b)} q(y) dy$$

true for all $a, b \Rightarrow$ identity.

$$\therefore q(y) = p(\bar{g}(y)) \cdot \frac{d}{dy}(\bar{g}(y))$$



If $g(x)$ is strictly decreasing? ...

$$q(y) = p(g^{-1}(y)) \cdot \left| \frac{d}{dy} g^{-1}(y) \right|$$

works for increasing
 $g(y)$ too!

Examples :-

1) $x \sim U(0,1)$

$$y = -\frac{1}{\lambda} \cdot \log(x)$$

$$g^{-1}(y) = ?$$

$$g^{-1}(y) = e^{-\lambda y}$$

$$\therefore q(y) = p(e^{-\lambda y}) \cdot \left| -\lambda \cdot e^{-\lambda y} \right|$$

$$= \frac{1}{\sqrt{2\pi}} \cdot \lambda \cdot e^{-\lambda y}$$

$\therefore y \rightarrow \text{exponential pdf}, \mu_y = \lambda$

2) $x \sim G(0,1)$ gaussian

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

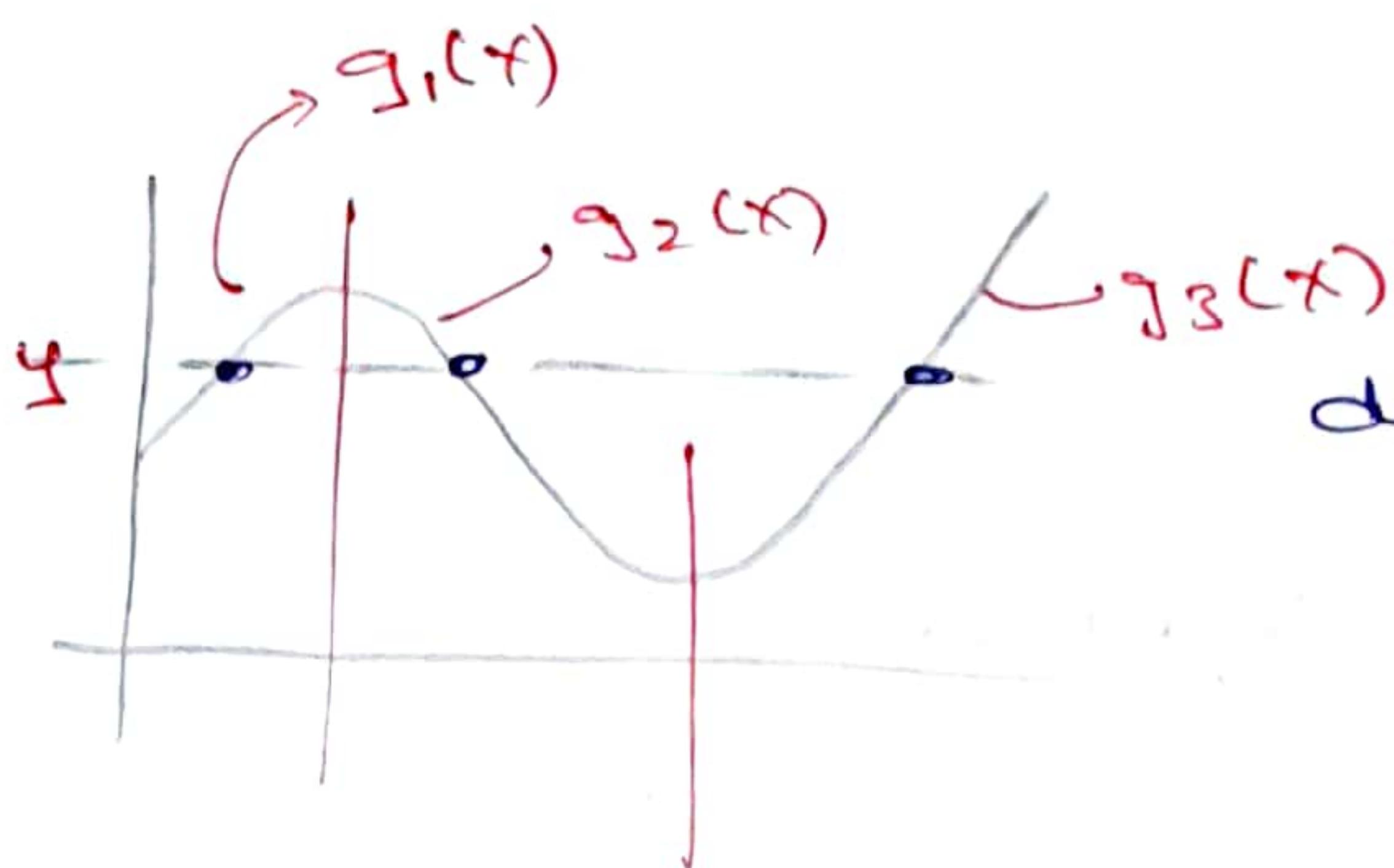
$$y = ax$$

$$\therefore g'(y) = \frac{1}{a}$$

$$\therefore q(y) = \frac{1}{\sqrt{2\pi}} \cdot e^{-x^2/2a^2} \times \left| \frac{1}{a} \right|$$

$$q(y) = G(0, a^2)$$

* General monotonicity functions:-

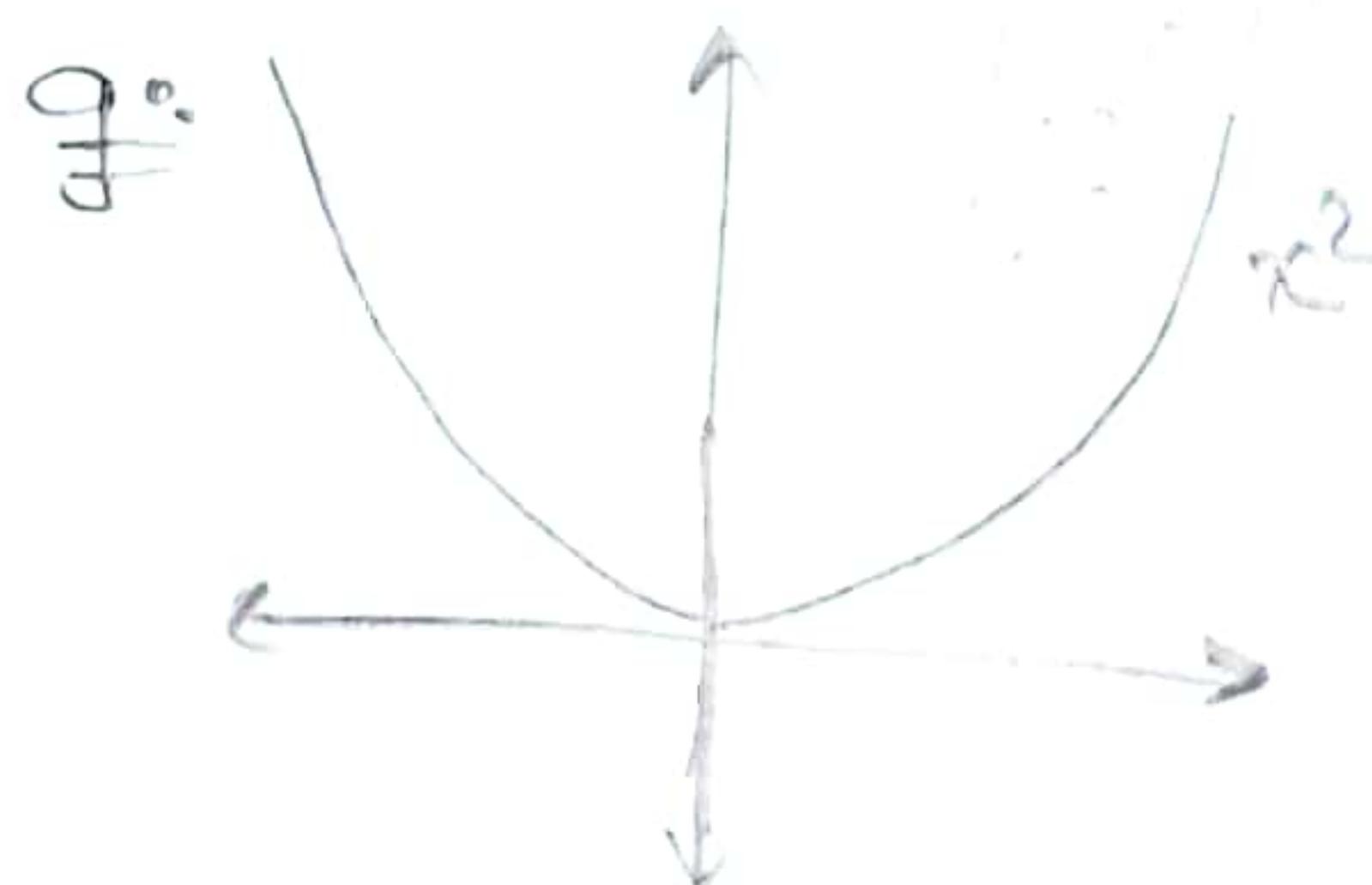


divide $g(x)$ into monotonous functions.
& add their respective $q_i(y)$.

Eg: $P(x)$ is:
0 for $x \leq -1$
0.5 for $x \in [-1, 0]$
 $0.5 e^x$ for $x > 0$.

$$\text{Ex } \underline{Y = x^2}$$

$$q(y) = ?$$



case-1: $x > 0; y \in (0, \infty)$

$$x = \sqrt{y}$$

$$\therefore q(y) = P(\sqrt{y}) \cdot \left| \frac{1}{2\sqrt{y}} \right|$$

$$q(y) = \underbrace{\frac{0.5 \cdot e^{-\sqrt{y}}}{2\sqrt{y}}}_{\text{case-1 formula}}$$

case-2: $x \in (-1, 0) \quad y \in (0, 1)$

$$x = -\sqrt{y} \quad y \in (0, 1)$$

$$\therefore q(y) = P(-\sqrt{y}) \cdot \left| \frac{1}{2\sqrt{y}} \right|$$

$$q(y) = \underbrace{\frac{0.5}{2\sqrt{y}}}_{\text{case-2 formula}}$$

\therefore i) for $y \in (0, 1)$

$$q(y) = \underbrace{\frac{1}{4\sqrt{y}}(1 + e^{-\sqrt{y}})}_{\text{case-1 formula}}$$

(jump discontinuity).

ii) for $y \geq 1$

$$q(y) = \underbrace{\frac{1}{4\sqrt{y}} \cdot e^{-\sqrt{y}}}_{\text{case-2 formula}}$$

Eq 2:

$$x \sim \mathcal{N}(0,1)$$

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

$$\underline{y = x^2}$$

what is $p(y)$ (chi-square.)

1) $x > 0; y > 0$

$$x = \sqrt{y}.$$

$$\therefore q(y) = p(\sqrt{y}) \cdot \frac{1}{2\sqrt{y}}$$

$$= \frac{1}{\sqrt{2\pi}} \cdot e^{-y/2} \cdot \frac{1}{2\sqrt{y}}$$

2) $x < 0; y > 0$

$$x = -\sqrt{y}$$

$$q(y) = p(-\sqrt{y}) \cdot \frac{1}{2\sqrt{y}}$$

$$= \frac{1}{\sqrt{2\pi}} \cdot e^{-y/2} \cdot \frac{1}{2\sqrt{y}}$$

$$\therefore \text{for } y > 0; \underline{q_{\text{net}}(y) = \frac{1}{\sqrt{y}\sqrt{2\pi}} \cdot e^{-y/2}}$$

→ Gamma pdf :- (not poission pdf)
discrete pdf.

$$p(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} \cdot x^{\alpha-1} \cdot e^{-\beta x} = \text{Gamma}(x | \alpha, \beta)$$

where Γ is 1.

- α (shape parameter) $> 0, x > 0$
- $\Gamma(\alpha) = \int_0^\infty e^{-x} \cdot x^{\alpha-1} dx; \alpha$ is any complex number.
- $\text{gamma}(n) = (n-1)! \text{ for int } n.$

say $Y = Y_X$.

$Y \rightarrow$ inverse gamma pdf.

$$\text{then } \text{pdf}(Y) = \frac{\beta^\alpha}{\Gamma(\alpha)} \cdot y^{\alpha-1} e^{-\beta/y}$$

you got to do

$$\frac{d}{dy} g'(y) = -\frac{1}{y^2}$$

not
 $\frac{d}{dy} p(y) = (\text{big})$

take care;
be present!

1. Transformation of variables:-

$$Y = f(X); \quad \text{PDF}_Y = \text{PDF}_X(f^{-1}(Y)) \cdot \frac{d}{dy} f^{-1}(y)$$

2. Multivariate Statistics:-

Multivariate Gaussian:

terrific!

conservation
of probability
mass.

* vector random variable X :

$$X = [x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_D]$$

joint pdf:

- the rand. var. X has a multivariate gaussian PDF if \exists a finite set of

iid univariate standard-normal RVS

w_1, w_2, \dots, w_n such that $N(0, 1)$
(with $n \geq D$)

each x_d can be expressed as

If $n > D$ then we
can always reduce
to D .

$$x_d = \mu_d + \sum_{k=1}^n A_{dk} w_k$$

non-singular!

If $n < D$, easy to
fill null
rows with
to

$$\Rightarrow X = \mu + A \cdot W \quad \text{where } W \text{ is } [w_1 \\ w_2 \\ \vdots \\ w_n]$$

rank = D if $D \neq n$.

only for m.v. gaussian;
not any multivariate

$|A| = 0$
then!

* if $\mu = 0$ & $A = I_{D \times D}$; then $x_d = w_d$

$$\therefore P(X=w) = P_W(w=w) = \prod_{d=1}^D P_{W_d}(w_d=w_d) \quad (\text{Joint probability})$$

0 mean,

isometric/Spherical
gaussian.

$$= \frac{1}{(2\pi)^{D/2}} e^{-0.5 \cdot w \cdot w}$$

→ level sets:-

locus of points in D-dimension; which have the same $p_x(\mathbf{x})$
 $\hookrightarrow \mathbf{x}$ is a datum.

So; a column vector! $\mathbf{x} = (x_1, x_2, \dots, x_D)$

$$L_c(f) = \{(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) \mid f(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) = c\}$$

+ Hence, $\mu = 0$; $A = I_{D \times D}$;

then

level sets are hyper spheres.

∴ isometric/spherical gaussian.

+ what is pdf; for $\mu \neq 0 \in A \equiv$ non-singular; ~~diagonal~~ matrix.

$$\text{say } \mathbf{x} = A \cdot \mathbf{w} + \mu$$

then we need $\text{pdf}(\mathbf{x} = \mathbf{I})$.

Step 1) If $\mathbf{x} = \mathbf{z} \in X = A \cdot \mathbf{w} + \mu$, then $\mathbf{w} = ?$

$$\mathbf{w} = \bar{A}^{-1}(\mathbf{I} - \mu)$$

Step 2) Apply law of conservation of probability mass

* but what will be 'volume' factor? jacobian!
scaling of \bar{A}^{-1}

$$\text{volume factor scaling} = |\bar{A}^{-1}|$$

$$\therefore \text{pdf}_{\mathbf{x}}(\mathbf{x}) = \frac{\text{pdf}_{\mathbf{w}}(\bar{A}^{-1}(\mathbf{x} - \mu))}{|\bar{A}^{-1}|}$$

$$\therefore \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \bar{A}^{-1} \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

if $\mu = 0$:-

$$\therefore q(\mathbf{x}) = p(\bar{A}^{-1}\mathbf{x}) \cdot \frac{1}{|\bar{A}^{-1}|}$$

don't forget.

$$= \frac{1}{(2\pi)^{D/2} |\bar{A}|} \cdot \exp(-0.5 (\bar{A}^{-1}\mathbf{x})^T (\bar{A}^{-1}\mathbf{x}))$$

$$= \frac{1}{(2\pi)^{D/2} |\bar{A}|} \cdot \exp(-0.5 \mathbf{x}^T (\bar{A} \cdot \bar{A}^T)^{-1} \cdot \mathbf{x})$$

let ($C = \bar{A} \cdot \bar{A}^T$) We will see; C is covariance matrix of \mathbf{x} .

$$\therefore q(\mathbf{x}) = \frac{1}{(2\pi)^{D/2} |C|^{1/2}} \cdot \exp(-0.5 \mathbf{x}^T C^{-1} \mathbf{x})$$

if $\mu \neq 0$:-

$$q(\mathbf{x}) = \frac{1}{(2\pi)^{D/2} |C|^{1/2}} \cdot \exp(-0.5 (\mathbf{x} - \mu)^T C^{-1} (\mathbf{x} - \mu))$$

we are in mean centered world...

* Facts:- A is a square matrix.

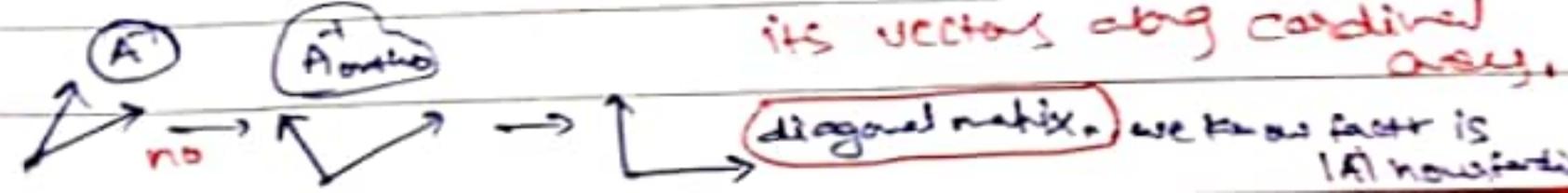
- is A transformation! Yes...

also A is non-singular \Rightarrow Gram-Schmidt orthogonalisation is possible.

∴ make $A^T \rightarrow \bar{A}^T$ orthogonal column matrix.

now; rotate \bar{A}^T orthogonal; to have

its vectors along cardinal axes.



→ Covariance matrix; for a general

- * if \mathbf{Y} is multivariate gaussian;

$\mathbf{Z} = \mathbf{B} \cdot \mathbf{Y} + \mathbf{C}$ is also multivariate gaussian.

$$(coz; \mathbf{Y} = \mathbf{A}\mathbf{W} + \boldsymbol{\mu}; \mathbf{Z} = (\mathbf{B} \cdot \mathbf{A})\mathbf{W} + (\mathbf{C} + \mathbf{B}\boldsymbol{\mu}))$$

fits into definition of

multivariate
gaussian

$$* \mathbf{X} = \mathbf{A}\mathbf{W} + \boldsymbol{\mu}.$$

$$\text{then mean}(\mathbf{X}) = \begin{bmatrix} E(X_1) \\ E(X_2) \\ \vdots \\ E(X_D) \end{bmatrix} = \boldsymbol{\mu}$$

• for general multivariate r.v \mathbf{X} ;

$$E(\mathbf{X}) = \begin{bmatrix} E(X_1) \\ P_x \\ E(X_2) \\ P_x \\ \vdots \\ E(X_n) \end{bmatrix} \equiv \begin{bmatrix} EP_{X_1}(X_1) \\ EP_{X_2}(X_2) \\ \vdots \\ EP_{X_n}(X_n) \end{bmatrix}$$

Important to note,

which pdf, we

$$EP_{X_n}(X_n)$$

are in.

$$C = \underline{E}_{P(Y)} \left[\begin{array}{c} (Y - E(Y))^* \\ \hline D \times 1 \end{array} \right] \left[\begin{array}{c} (Y - E(Y)) \\ \hline 1 \times D \end{array} \right]^T$$

Square matrix
 $D \times D$.

$$\therefore C_{ij} = \underline{E}_{P(Y)} [(Y_i - E(Y_i))(Y_j - E(Y_j))]$$

$$= \underline{E}_{P(Y_i, Y_j)} (Y_i - E(Y_i))(Y_j - E(Y_j))$$

$$= \underline{\text{cov}}(Y_i, Y_j)$$

$$* \text{ for } \mathbf{X} = \mathbf{A}\mathbf{W} + \boldsymbol{\mu}; \quad C(\mathbf{X}) = \mathbf{A} \cdot \mathbf{A}^T$$

- For a general vector r.v \mathbf{X} ; need not be m.gaussian.

let C be covariance matrix of \mathbf{X}

$$1) \quad C = \underline{E}(\mathbf{X} \cdot \mathbf{X}^T) - \underline{E}[\mathbf{X}] \cdot (\underline{E}[\mathbf{X}])^T \sim C_{ij} = \begin{array}{l} \text{sq. mtrx} \\ \hline D \times D \end{array} \quad \begin{array}{l} \text{Dx1} \\ \hline 1 \times D \end{array} \quad \begin{array}{l} \text{DxD} \\ \hline \end{array} \quad - E(X_i) \cdot E(X_j) \\ = \text{cov}(X_i, X_j) \quad \forall$$

2) C is symmetric. (as $C_{ij} = \text{Var}(X_i, X_j) = \text{Var}(X_j, X_i) = C_{ji}$)

3) C is positive semi definite.

⇒ for any nonzero $a \in \mathbb{R}^{D \times 1}$; $a^T \cdot C \cdot a \geq 0$

$$\text{since; } a^T \cdot C \cdot a = a^T \cdot E((\mathbf{X} - E(\mathbf{X})) \cdot (\mathbf{X} - E(\mathbf{X}))^T) \cdot a$$

$$= E(a^T \cdot (\mathbf{X} - E(\mathbf{X})) \cdot (\mathbf{X} - E(\mathbf{X}))^T \cdot a) = E((E(\mathbf{X}))^T \cdot a) = 0$$

→ terms:-

1) diagonal matrix.

2) orthogonal matrix.

$$\text{if } Q \cdot Q^T = Q^T Q = I$$

then $Q \rightarrow$ orthogonal

$$|Q| = \pm 1$$

3) Rotation matrix...

Q is orthogonal; and

$|Q|=1$ then rotational matrix.

Improperly, can also take

$|Q|=-1$ to be rotational;

But its rotation + reflection...

4) reflection:-

A symmetric orthogonal matrix is

$$\begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$$

reflection matrix.

lets discuss level sets :-

• Spectral theorem Says:

a real symmetric matrix (a self-adj matrix).

has a orthonormal basis of eigen vectors.

→ there exist orthogonal matrix Q ;
such that

$$C \cdot Q = Q \cdot D$$

diag. matrix of eigen values.

$$\therefore C = Q \cdot D \cdot Q^T = Q \cdot \underline{D} \cdot Q^T$$

* covariance matrix C

is symmetric & positive definite.

→ all its eigen values

are real & positive.

there exist
orthogonal Q ;

such that

$$Q^T \cdot C \cdot Q = D$$

$$\therefore C = Q \cdot D \cdot Q^T$$

rotation
matrix

Scaling
matrix for level
sets.

↳ D is all
positive.

Now;

$$P(x) = \frac{1}{(2\pi)^{D/2} |C|^{1/2}} \exp(-0.5 (x-\mu)^T \cdot C^{-1} \cdot (x-\mu))$$

$$= \frac{1}{(2\pi)^{D/2} |C|^{1/2}} \exp(-0.5 (x-\mu)^T \cdot Q \cdot D^{-1} \cdot Q^T (x-\mu))$$

$$\text{let } Q^T \cdot (x-\mu) = y$$

if eigenvalue is more;

variance is more;
ellipsoid axis's length more.

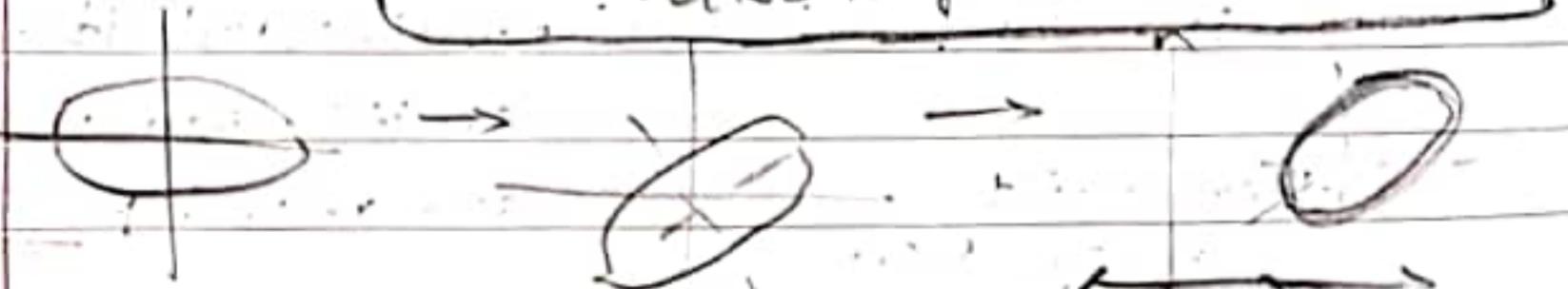
so in y's frame of axis;

level sets are hyperellipsoidal at origin;

with cardinal axes as
axis.

Now; $x = Q \cdot y + \mu$

\rightarrow translation of hyperellipsoid.
 \curvearrowleft rotation of hyperellipsoid.
 \curvearrowright $\because Q$ is orthogonal with x .



level sets in
[x's frame].
our original frame.

* say $n > d$!

$$X = A \cdot W + \mu$$

$D_{n \times n}$ $D_{n \times d}$ $\mu_{d \times 1}$ $D_{d \times 1}$

now also; $C = A \cdot A^T$

$$\text{also } C = Q \cdot D \cdot Q^T$$

$$= Q \cdot (JD)(JD)^T \cdot Q^T$$

\therefore possibly;

$$Q \cdot JD = A$$

no matrix;
only feel...

$A = (e_{11}, e_{12}, \dots)$

eigen value

$A \cdot (I) \cdot Q^T = \text{diagonal} \dots$

like

$$\text{new } W = \begin{bmatrix} 1 & 0 & 0 & 4 & 0 & 0 \\ 0 & 2 & 0 & 0 & 5 & 6 \\ 0 & 0 & 3 & 0 & 0 & 0 \end{bmatrix}$$

NOW;

$(W_1 + W_2)$ is
a new W_{n1}

$(W_2 + W_3 + \dots)$
is new W_{n2}

\downarrow
we can reduce

n to D .

add "of"

$$a_1 \cdot N(0, 1) + b_1 \cdot N(0, 1)$$

independent

$\approx A \cdot N(0, 1)$
ageingaussian!

so $n \rightarrow d$ ✓

since A is non-singular;

gram-Schmidt orthogonalisation

possible.

then rotate A onto;

to get

diagonal matrix D .

then

$$A \cdot (B) = D \cdot w \quad |B| \neq 0$$

$$\text{if } A = D \cdot B^T \quad |B^T| \neq 0$$

$$\therefore x = D \cdot B^T(w) + \mu \quad B^T \stackrel{\text{rank}}{=} D \cdot v$$

Hence, any ways; we
can again have single! instead of linear
comb.
 $x_1 = \text{exit } A_1 \cdot W_1; x_2 = \text{exit } A_2 \cdot W_2 \dots$

→ Marginal PDFs:

meaning: $\text{pdf}_{x_1}, \dots, \text{pdf}_{x_2, x_1, \dots}$

- * the 1D marginal PDF is univariate gaussian, for any multivariate gaussian X .

(Joint).

- * marginal pdf of multivariate N in N dimensions, over any chosen subset of variables of size M is multivariate gaussian pdf.

By definition:

$N \geq D$; not $N = D$.

If $N > D$; then we can "reduce" N to D .

Proof:

let $X = A w + \mu$
 $\overset{N \times P}{\text{invertible}}$

B = projection matrix, such that

$B \in \mathbb{R}^{M \times N}$

$B \cdot X$ is the subset.

$\text{rank}(B) = M$

$\therefore A$ is invertible.

$\Rightarrow \text{rank}(AB) = M$

$Bx = BAw + B\mu$

$y = (BA)w + (\beta A)_{M \times 1}$

$\text{rank} = M \cdot \checkmark$

hence

multivariate gaussian!

$B = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}_{2 \times 4}$

or

$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}_{2 \times 4}$

* But margins being gaussian
won't imply joint is gaussian.

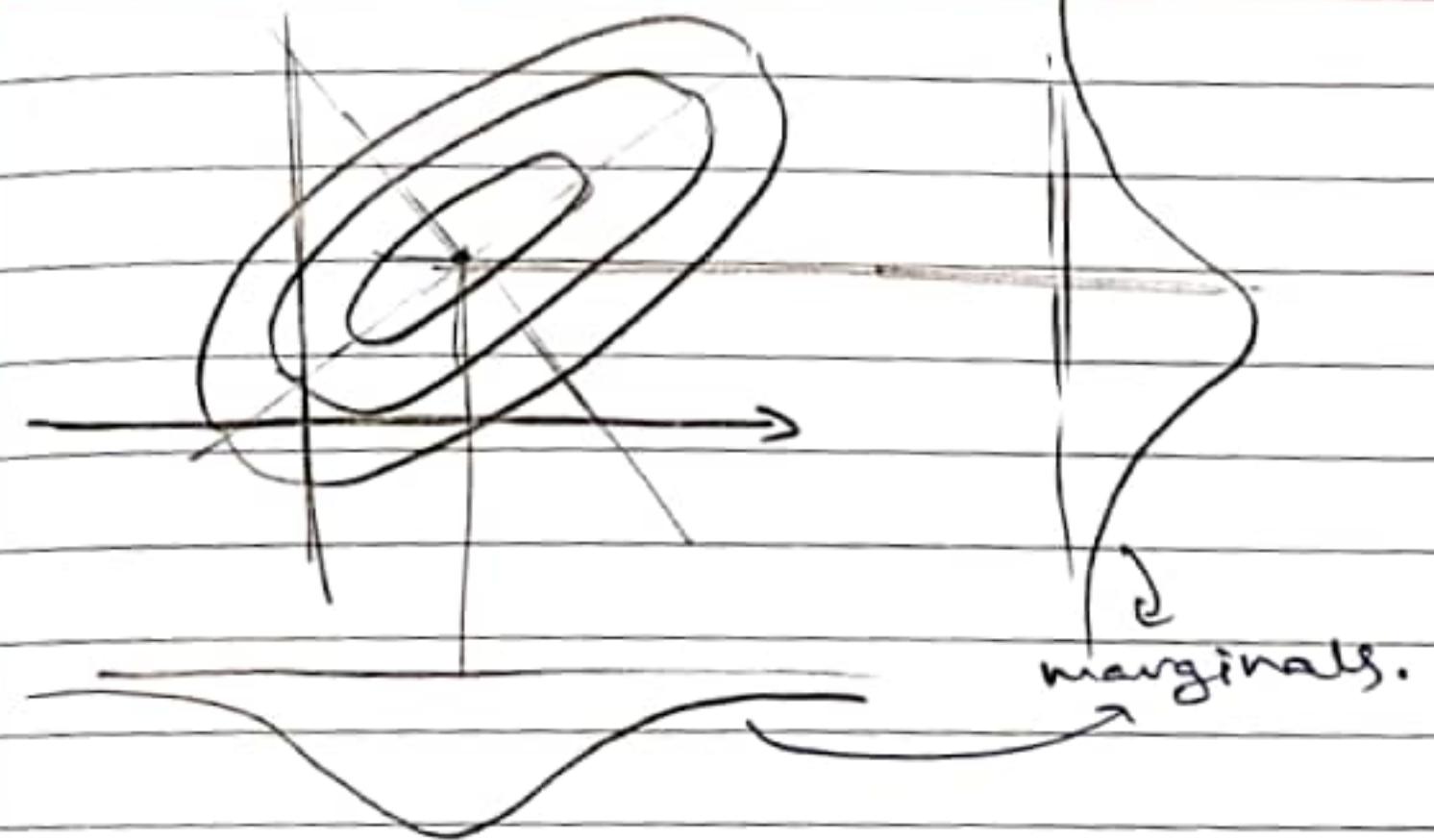
$$\text{if } X_1 = w_1$$

$$X_2 = w_2$$

$$\text{then; } X = \begin{bmatrix} 10 \\ 10 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

vector plot 2

Hence; X is not gaussian.



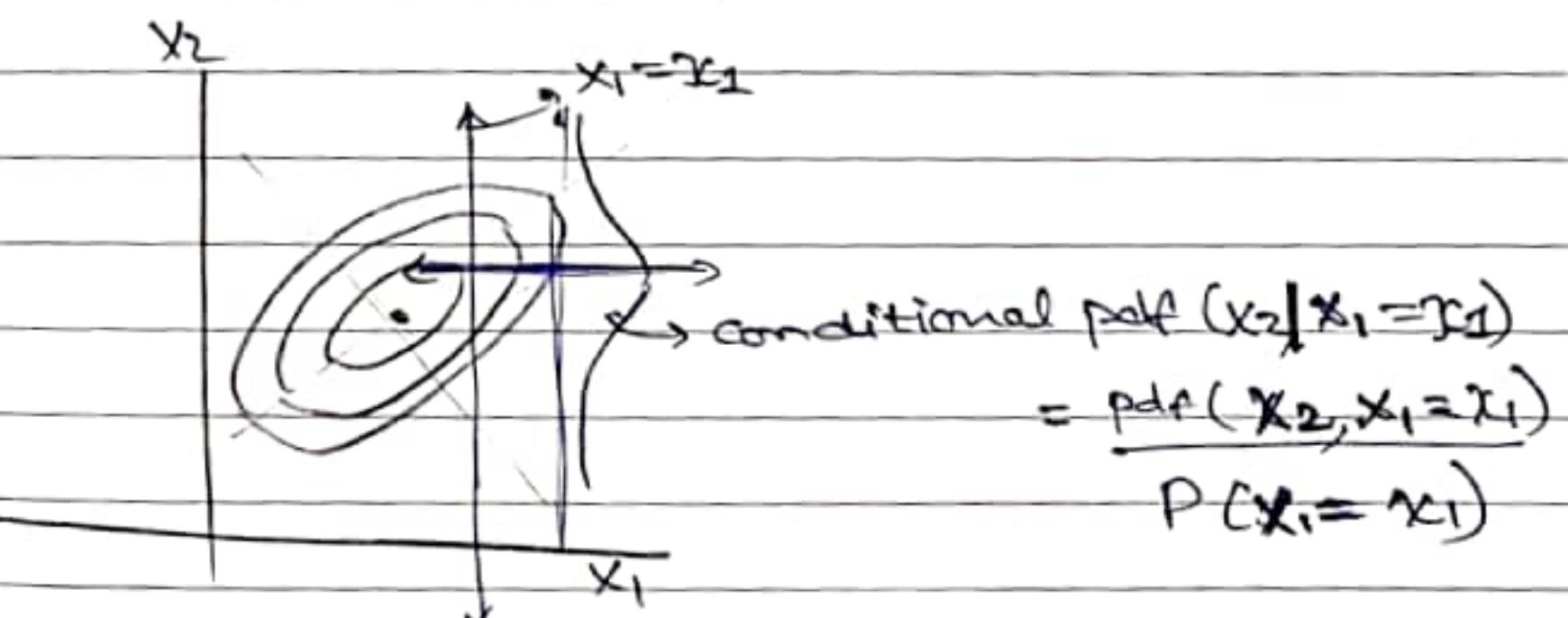
→ conditional pdfs

$$\text{let } X = [x_1 \ x_2]$$

again multivariate

then $P(X_1=x_1 | X_2=x_2)$ is also a m.v.gaussian.
conditional

$$P(X_1 | X_2=x_2) = \frac{P(X_1, X_2=x_2)] \text{ a multivariate gaussian}}{P(X_2=x_2) \text{ normalising factor}}$$



classmate

Date _____

Page _____

classmate

Date _____

Page _____

\tilde{C} ; becoz we do $Pw(\tilde{g}(y))$. classmate
factor



NOT

→ Maximum Likelihood estimation → for general!

for gaussian.

$$p(y) = \frac{1}{(2\pi)^{\frac{N}{2}} |C|^{\frac{1}{2}}} \cdot \exp^{-0.5(y-\mu)^T \cdot \tilde{C}^T (y-\mu)}$$

$\tilde{C} = Q \cdot D^{-1} \cdot Q^T$

data is $\{y_1, y_2, \dots, y_N\}$

these datums are vectors.

1) * MLE for mean vector....

(matrix calculus)

$$L_L = \sum (y_i - \mu)^T \cdot \tilde{C} \cdot (y_i - \mu) + \{\text{const}\} \sqrt{N}$$

$$\frac{\partial L_L}{\partial \mu} = 0.$$

$$\frac{\partial}{\partial \mu} (x - \mu)^T \cdot \tilde{C}^T \cdot (x - \mu) = 2 \tilde{C}^T (x - \mu)$$

$\begin{matrix} 1 \times N & N \times N & N \times 1 \\ N \times 1 & N \times N & N \times 1 \end{matrix}$

$$\rightarrow \sum 2 \tilde{C}^T (y_i - \mu) = 0$$

$\Rightarrow \mu = \text{mean of } y_i$

Note:

scalar fun. $f(a)$ of multiple variables in column vector

Jacobian; $\frac{df}{da}$ will be a row vector of same length a .

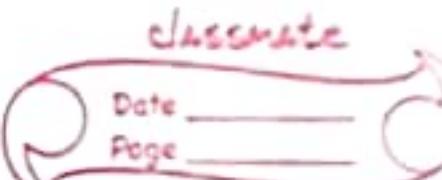
$\therefore (\hat{\mu}) = \text{Sample mean}$

a vector random variable.

$$\therefore df = \frac{df}{da} * da$$

$1 \times N \quad N \times 1 \quad N$

symmetric positive definite



2) MLE for C := Sample covariance.

biased one!

$$LL = \sum \frac{(y_i - \mu_i)^T \cdot \tilde{C} \cdot (y_i - \mu_i) + \log |C|}{2} + \text{const...}$$

$\sqrt{N}/n \text{ not } \sqrt{n-1}$

$$\frac{\partial LL}{\partial C} = 0$$

⇒ partial diff wrt C : for a scalar fun., $f(C)$

1) $C_{N \times N}$ to column vector $a_{N^2 \times 1}$

2) find $\frac{\partial f}{\partial a} \cdot a$ row vector $1 \times N^2$

3) reshape into $N \times N$

$$\therefore df = \frac{\partial f}{\partial C} * dc$$

element wise product.

$$\frac{\partial}{\partial C} (x - \mu)^T \cdot \tilde{C} \cdot (x - \mu) = -\tilde{C}^T \cdot (x - \mu)(x - \mu)^T \cdot \tilde{C}^T$$

$$\frac{\partial}{\partial C} \log |C| = \frac{1}{|C|} \cdot |C| \cdot \tilde{C}^T = \tilde{C}^T$$

∴ we get

$$-\tilde{C}^T (\sum (x_i - \mu)(x_i - \mu)^T) \cdot \tilde{C} + n \cdot \tilde{C}^T = 0$$

taking μ known,
sample covariance

$$\tilde{C} = \frac{\sum (x_i - \mu)(x_i - \mu)^T}{N}$$

x_i is a vector with x_{ij}

→ Mahalanobis distance: - for multivariate gaussian.

the term: $(y-\mu)^T \bar{C}^{-1} (y-\mu)$ is ^{square} mahalanobis distance of y from μ .

$$d(y, \mu; C) = \frac{(y-\mu)^T \bar{C}^{-1} (y-\mu)}{C}$$

* when $C = I$:

mahalanobis dist. = euclidean distance.

or else;

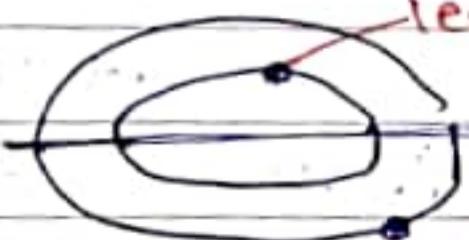
$d(y, \mu; C)$ scales units along each dimension based on standard deviation along that not variance. marginal.

eigen values = variance.
(σ^2).

locus is hyper ellipsoid

= level set for distribution.

less mahalanobis dist. from μ .



more Mahalanobis dist.

* Mahalanobis distance; is a true distance metric function $d(\dots, \dots) \rightarrow \mathbb{R}$ satisfies metric.

(i) identity of indiscernibles:

$$d(x, y) = 0 \Leftrightarrow x = y$$

(2) Symmetric:

$$d(x, y) = d(y, x)$$

(3) triangle inequality.

$$d(x, y) \leq d(x, z) + d(z, y)$$

(These three imply non-negativity!)

for $d(x, y; C)$

- (i) ✓ identity of indiscernibles
- (ii) ✓ symmetric.

(iii) let $C = D \cdot (\text{diag component})$

see $d(x, y) \leq d(x, z) + d(y, z)$

$$\text{LHS} = \sqrt{(u+v)^T \cdot C^{-1} \cdot (u+v)}$$

$$\text{RHS} = \sqrt{u^T \cdot C^{-1} \cdot u} + \sqrt{v^T \cdot C^{-1} \cdot v}$$

\int is needed!!

$$\text{LHS}^2 = u^T C^{-1} u + v^T C^{-1} v + u^T C^{-1} v + v^T C^{-1} u$$

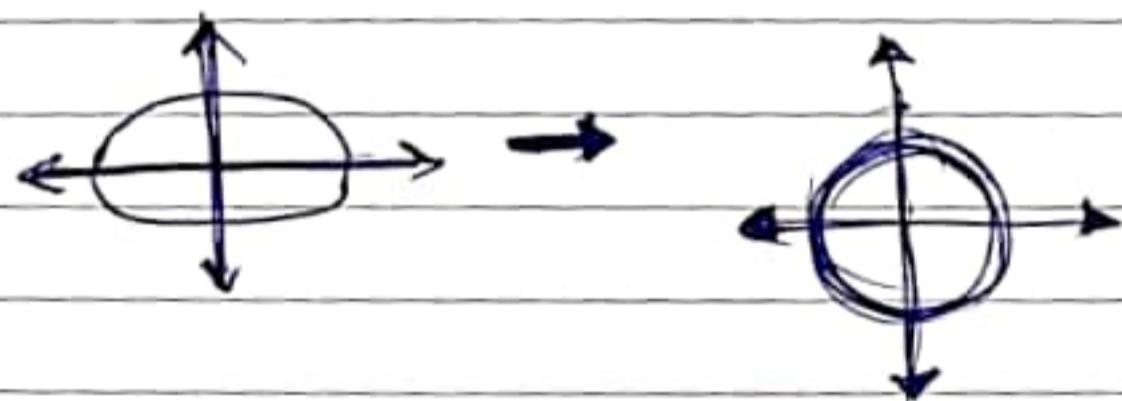
$$\text{RHS}^2 = u^T C^{-1} u + v^T C^{-1} v + 2 \cdot \sqrt{u^T C^{-1} u} \cdot \sqrt{v^T C^{-1} v}$$

now: $C^{-1} = D^{-1}$

$$\text{RHS}^2 - \text{LHS}^2 = 2 \cdot \sqrt{\sum \frac{u_i^2}{\sigma_{u_i}^2}} \cdot \sqrt{\sum \frac{v_i^2}{\sigma_{v_i}^2}} - 2 \cdot \sum \frac{u_i v_i}{\sigma_{u_i} \sigma_{v_i}}$$

Cauchy-Schwarz inequality.

* Mahalanobis scales the coordinate plane



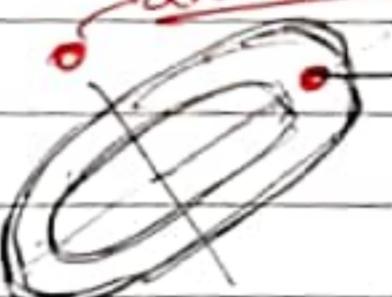
Multivariate Gaussian Applications:-

1) Anomaly detection: whether new datum

fits in (or)
no.

anomaly.

not anomaly.



if Mahalanobis
distance
low \rightarrow fits in,
otherwise

anomaly.

2) maximum likelihood classification.

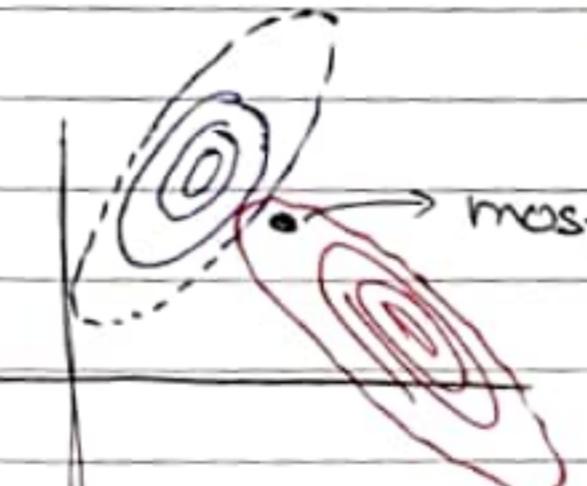
A datum given. belongs to which group?

compare the likelihood of π_i ,

for each group individually,
choose max.

• Mahalanobis dist. isn't enough.

also see $\frac{1}{\kappa_1^{1/2}}$ term.



decision boundary: $(x - m_1)^T C_1^{-1} (x - m_1) = (x - m_2)^T C_2^{-1} (x - m_2)$

decision surface: $+ \log |C_1| + \log |C_2|$

hyperquadric

if $|C_1| = |C_2| \wedge C_1 = C_2 = C$, only mean shift is different; then hyperplane dev.

classmate

Date _____

Page _____

classmate

Date _____

Page _____

if $C_1 = C_2 = C \ \& \ \mu_1 \neq \mu_2$

decision surface
is hyper
plane

Principle Component Analysis:-

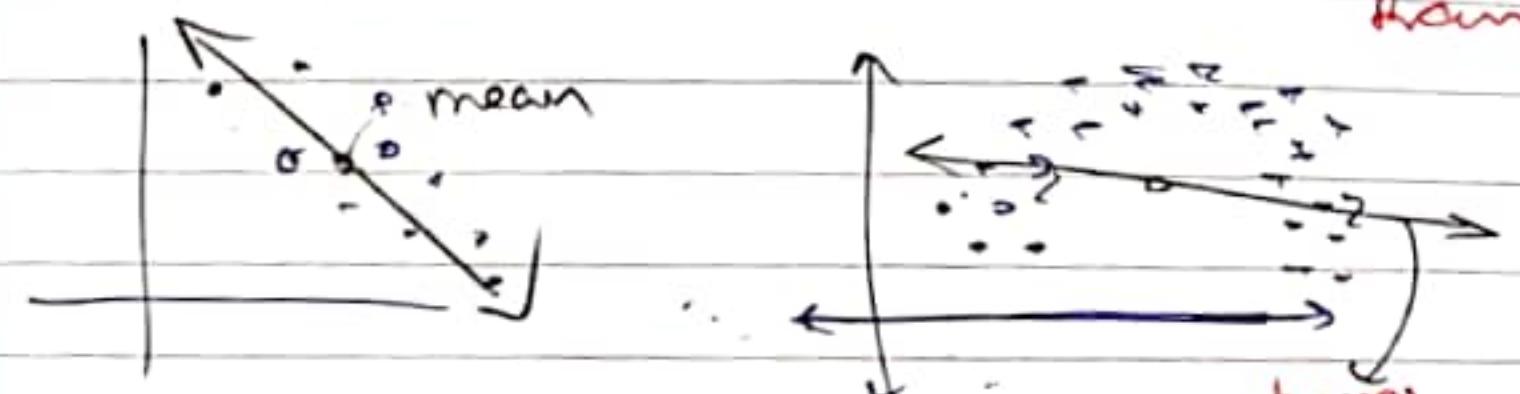
for general vector R.V.
not only gaussian.

modes of variation:-

Set of vectors (so, both mag. & direct)
that are used to depict the variation in
population, around the mean.

∴ these vectors are drawn from
mean.

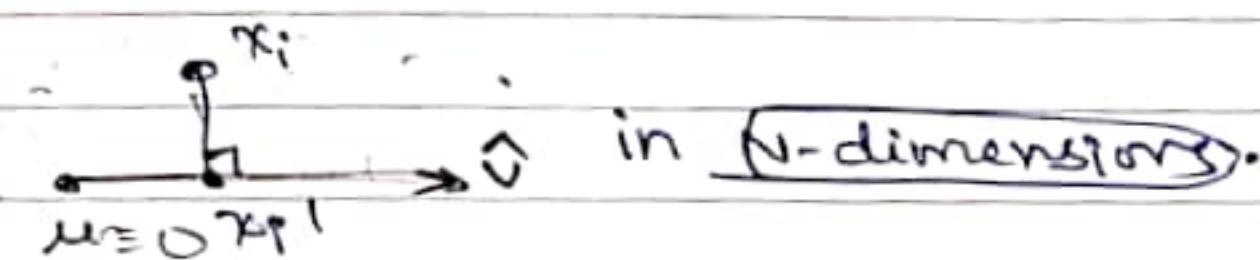
"in mean-centered
form."



1 pca is
disaster.
2 would be
good.

we want the 'direction' which maximizes
variance of points incident on it; wrt mean.

"it captures max. dispersion of dataset!"



now; $x_i^j = \langle x_i, v \rangle v$ j.v is unit vector,

∴ var. of all projections

$$= \frac{\sum_{i=1}^N \langle x_i, v \rangle^2}{N}$$

$$= \frac{\langle x_i, v \rangle \langle x_i, v \rangle^T}{N}$$

$$= V^T \cdot C \cdot V \quad (\text{since } C \text{ is diagonal})$$

now; let $v = (v_1, v_2, \dots, v_n)$

$$\sum v_i^2 = 1 \quad (\text{unit vector})$$

$$\sum v_i^2 \cdot C_{ii} \rightarrow \max$$

∴ chose $v_i = 1$ for argmax Cii

& $v_i = 0$ for others.

* let C = diagonal.

let data $\{x_i\}_{i=1, n}$.

∴ mean = μ .

let vector = v_i .

no need of
gaussian

→ we choose the cardinal axis with max. eigen value.

(for diagonal C_1 ;

coordinate axes are eigen vec;

I is eigen matrix,
the orthogonal Q).

∴ 1 pca direction = eigen vector of map eig. value.

* 2nd pca definition:-

direction, \perp to 1st pca;
& capturing max. dispersion.

then,

as predicted,

= eigen vec; with 2nd map eig. val.

* even if $C = Q \cdot D \cdot Q^T$;

$$U^T \cdot C \cdot U$$

$$U^T \cdot Q \cdot D \cdot Q^T \cdot U$$

$$(Q^T U)^T \cdot D^T \cdot (Q^T U) \text{ is max.}$$

so

now;

for $C_{n \times n}$:

, we have, n pcas.

but, if we have

15 observations

& 200 dimension;

only (14) pca can

represent whole data set.

as

$C_{n \times n}$ is
sym;

n eigen
vectors w/

variance of projected
points = eigen
value

(--).

* now; I want the D-dimensioned space;
which captures max. dispersion.

≡ linear span of first D-pcas.

E, variance is sum of eigen values.
lets do!

$$v_1, v_2, \dots, v_d.$$

$$x_1 = \langle x_1, v_1 \rangle v_1 + \langle x_2, v_2 \rangle v_2 + \dots$$

$$\begin{matrix} x_2 \\ x_3 \\ \vdots \\ x_n \end{matrix}$$

dist. from mean

$$= \sqrt{\sum_{i=1}^n v_i^2} = \sqrt{\sum_{i=1}^n v_i^2} = \dots$$

$$\sqrt{v_1^2 + v_2^2 + \dots}$$

For gaussian multivariate:-

$$X = Aw + b.$$

$$\text{then } C = A \cdot A^{-1}$$

E_p

principle modes; are axis of hyperellipsoid.

closed
set).

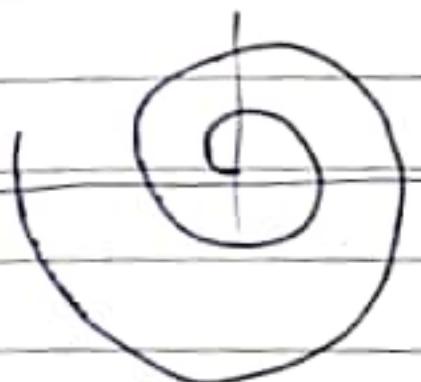
*

4) Dimensionality reduction:-

intrinsic dimension:-

minimum no. of variables
(degrees
of freedom)

required to represent signal.



*

even though every point
is $f(x,y)$; it can be written
some function. as some

$$= (x,y)$$

$$f(O)$$

$$= (x(0), y(0))$$

\therefore degree of freedom = 1.

Intrinsic dim. = 1.

representation
dimension = 2.

* let $X = (x_1, x_2, \dots, x_N)$

& consider function $g(\cdot)$ & $M < N$ variables a_1, a_2, \dots, a_M .
such that

each $x \in X$ is writable as
 $g(a_1, a_2, \dots, a_M)$.

then; instead

of 19200 variables; a fruit just needed
4 variables.
<heat attack>.

- * Hence; by PCA; we find first 4 pca vectors

g,

project the whole data onto this hyperplane.

∴ dimension reduced now!

$$19200 \rightarrow 4$$

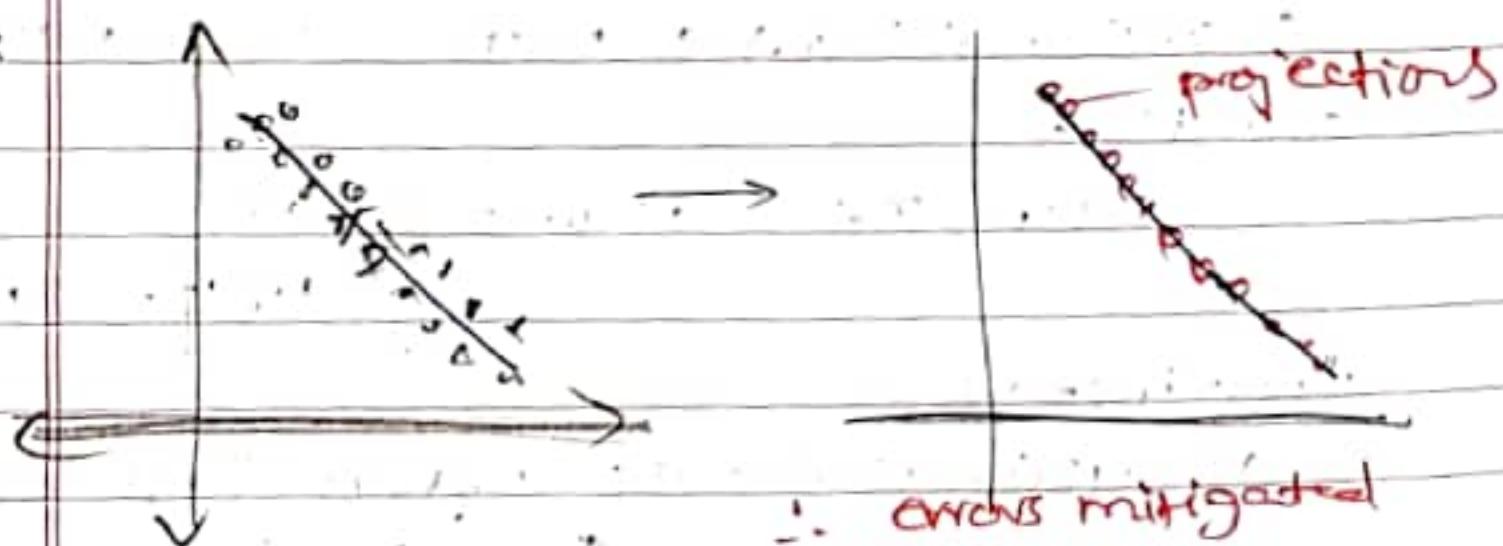
- * i) Acquired data has errors.

- this can make data seem of higher dimns!

∴ PCA performs linear dim. reduction.

(fits hyperplanes only.)

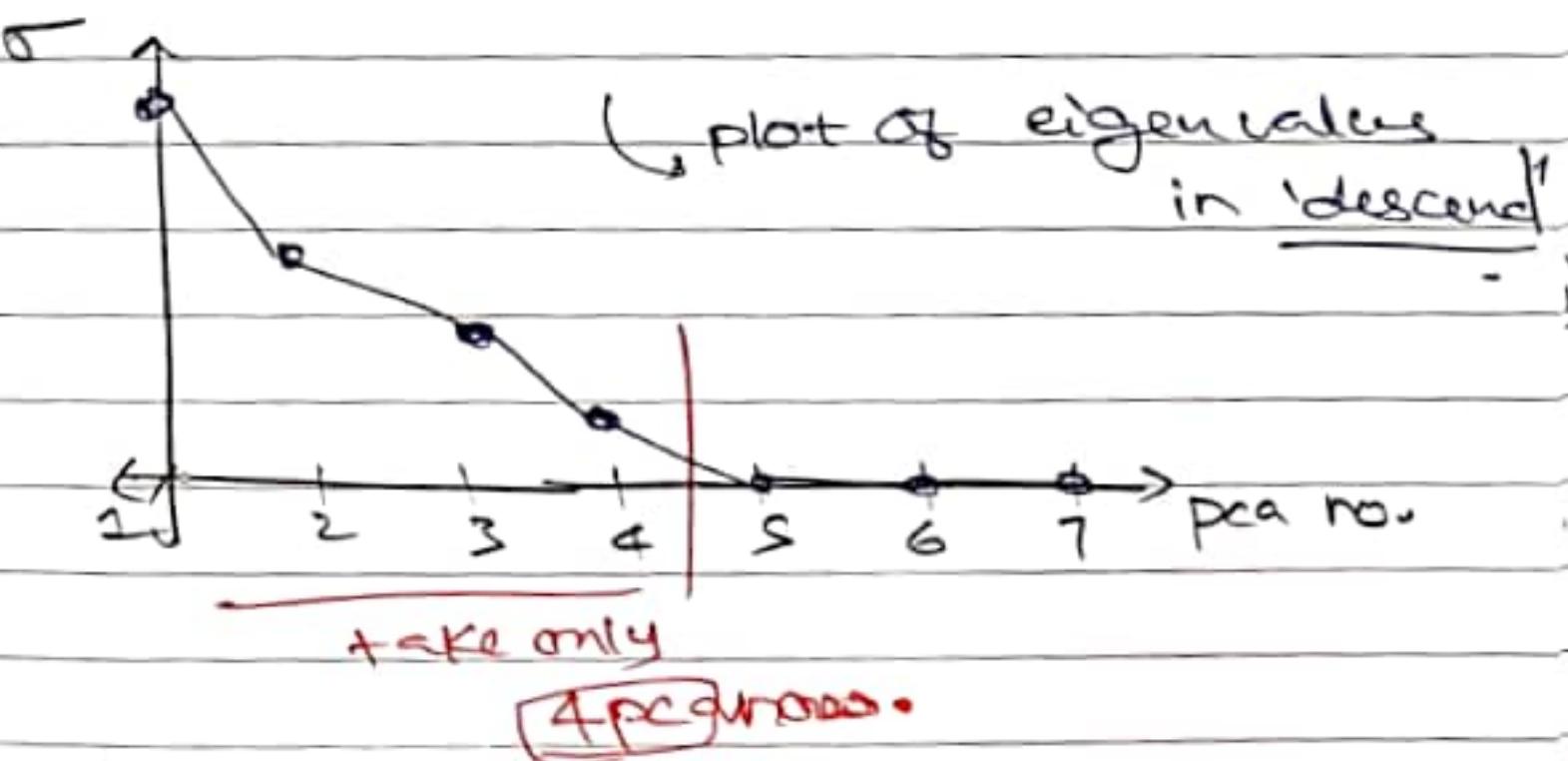
So; if hyperquadric needs 3 dim hyperplane make it 4.



- * How many dimensions should I reduce to?

Scree plot!

- 1) variance along a mode = eigenvalue.
- 2) we need to capture maximum variance of datums.
- 3) Hence; chose the dimensions; which have considerable eigen values.



- 4 pca; can capture most of data's variance.

$A \rightarrow$ linear operator. $U, S, V^T \rightarrow$ steps of the operator.Singular Value Decomposition

(SVD)

* matrix factorisation.

- let A be $M \times N$. our interest is $M \leq N$.

But discussed generally.

when A is real valued, SVD of $A = USV^T$
 $M \times M \quad M \times M \quad N \times N$ $V \rightarrow$ orthogonal $N \times N$ $U \rightarrow$ orthogonal $M \times M$ $S \rightarrow$ rectangular diag matrix. $M \times N$ values of diag \equiv Singular values.* sing. values are
non-negative, real.column vectors of $U, V \rightarrow$ singular vectors.of $U \rightarrow$ left singular vectorsof $V \rightarrow$ right singular vectors.When $M \leq N$; then

$$A = \sum_{m=1}^M S_m \cdot U_m \cdot V_m^T$$

scalar

$$S = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 4 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad 3 \times 5.$$

recto diagonal

* $A = USV^T$ • then $A \cdot z = U \cdot S \cdot (V^T \cdot z)$

rotation
scaling
rotation.

→ Matrix Norm:-

for $x \in \mathbb{R}^n$; 2-norm is $\|x\|_2$.

for matrix A; of $M \times N$

definition:-

$$\|A\|_2 = \max_{x \neq 0} \left\{ \frac{\|Ax\|_2}{\|x\|_2} \right\} \geq 0$$

how to do
manually?

x is $N \times 1$
 Ax is $M \times 1$

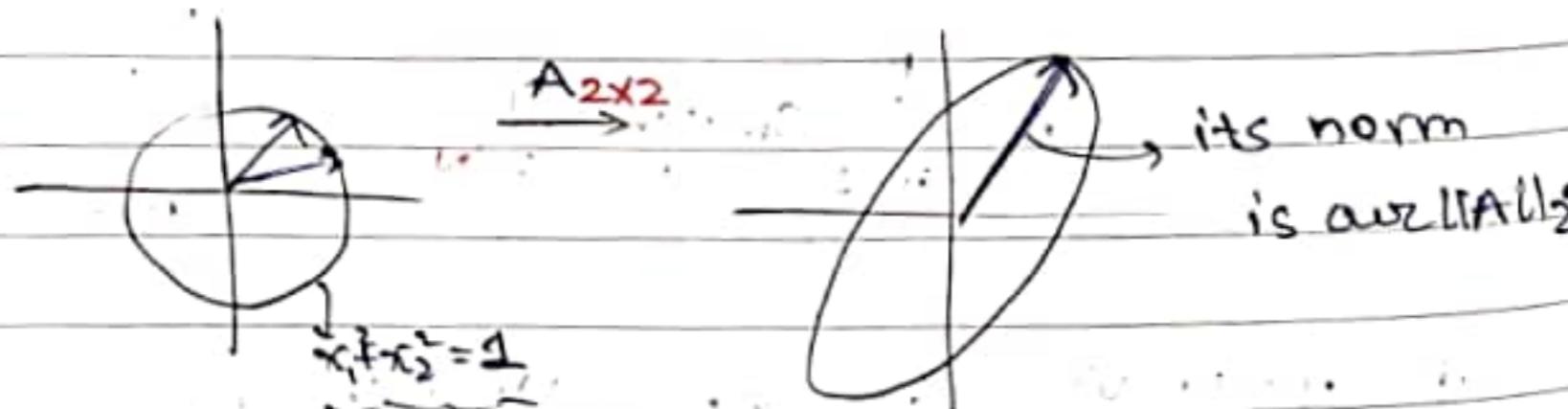
* Geometric interpretation is;

apply linear operator A to all unit vectors x

$N \times 1$

let $y = Ax$; for all such x .

then, pick the longest "y", & report its 2-norm.



→ Existence of Singular value decompos-

exists, for any real matrix $A \in \mathbb{R}^{M \times N}$.

• let $\sigma_i = \|A\|_2$:

$A \in \mathbb{R}^{M \times N}$

E. $\exists v \in \mathbb{R}^n$ such that $\|v\|=1$

$$v, \sigma_i u_i = A \cdot v,$$

$$\|u_i\| = \sigma_i / \sigma_i$$

• consider orthogonal matrix U, a basis of \mathbb{R}^M with first column u_1 .

$\|u_i\|$ not simply u_i .

also orthogonal matrix V; basis for \mathbb{R}^N with first column v_1 .

$$\underline{\underline{\|v_1\|=1}}.$$

$$\text{then, } U^T \cdot A \cdot V = \begin{bmatrix} \sigma_1 u_1 & W^T \\ 0 & B \end{bmatrix}_{M \times N \times M \times N - 1} = S$$

• we will show that $W = [0]_{N-1 \times 1}$.

NOTE: $v_{N \times 1}$ is a column vector; & A is orthogonal $M \times N$

$$\|V\|_2 = \|A \cdot V\|_2$$

for any V . \therefore orthogonal matrix \equiv rotation.

$$* S = U^T A V$$

$$\text{since: } \max_x \|A \cdot x\| = \sigma_1, \|x\|=1$$

$$\Rightarrow \max_x \|A \cdot V \cdot x\| = \sigma_1$$

$$\|V \cdot x\| = \|x\| = 1$$

$$\Rightarrow \max_x \|U^T \cdot A \cdot V \cdot x\| = \|A \cdot V \cdot x\| = \sigma_1$$

$$\Rightarrow \max_x \|S \cdot x\| = \sigma_1$$

$$\therefore \|S\|_2 = \|A\|_2 = \sigma_1$$

$$(\text{if say } \|S\|_2 = 2\sigma_1)$$

$$\Rightarrow \exists x \text{ s.t. } \|x\|=1 \& \|S \cdot x\|=2\sigma_1$$

$$\Rightarrow \|U^T A V \cdot x\|=2\sigma_1$$

$$\Rightarrow \|A \cdot (Vx)\|=2\sigma_1$$

but $\max_y \|A \cdot y\| = \sigma_1$
contradict!

$$\therefore \|S\|_2 = \|A\|_2$$

$$\text{now: let } x = \begin{bmatrix} \sigma_1 \\ w \end{bmatrix}_{N \times 1}$$

$$\therefore \|Sx\| = \sqrt{\sigma_1^2 + w^T w}$$

$$\therefore \|Sx\| = \frac{1}{\sqrt{\sigma_1^2 + w^T w}} \cdot \left\| \begin{bmatrix} \sigma_1 & w^T \\ 0 & B \end{bmatrix} \cdot \begin{bmatrix} \sigma_1 \\ w \end{bmatrix} \right\|$$

$$= \frac{1}{\sqrt{\sigma_1^2 + w^T w}} \cdot \left\| \begin{bmatrix} \sigma_1^2 + w^T w \\ Bw \end{bmatrix} \right\|_{M \times 1}$$

$$\geq \frac{1}{\sqrt{\sigma_1^2 + w^T w}} \cdot (\sigma_1^2 + w^T w)$$

$$\therefore \|Sx\|_{\min} = \sqrt{\sigma_1^2 + w^T w}$$

$$\max_x = \sigma_1$$

$$\Rightarrow w^T w = 0$$

$$\therefore w = \underline{\underline{[0]}}_{M-1 \times 1}$$

$$\therefore \text{now, } UTAV = \begin{bmatrix} \sigma_1 & [0 \quad \dots] \\ [0] & B \end{bmatrix}_{M \times N}$$

claim: g can make S as diagonal, by doing same kind of decompositions.
proof: via induction.

$$MXN \rightarrow M-1 \times N-1 \rightarrow \dots \rightarrow 1 \times K \text{ or } K \times 1$$

$$\text{let } B = [a_1 \ a_2 \ \dots \ a_k]$$

$$\& \|B\| = \sigma_1$$

$$\text{then } U^T B V = [\sigma_1 \ 0 \ 0 \ 0 \ 0]$$

$$\therefore B = U \cdot [\sigma_1 \ 0 \ 0 \ \dots \ 0] \cdot V^T$$

Base case proved

Same for $K \times 1$.

let it hold for $M \times N \times 1$

orthogonal

$$\Rightarrow B = U_1 \cdot D \cdot V_1^T \quad D \rightarrow \text{diagonal}$$

$M \times 1 \quad M \times 1 \quad M \times 1 \quad N \times 1$
 $X \times N \quad X \times N \quad X \times N \quad X \times N$

$$\therefore U^T \cdot A \cdot V = \begin{bmatrix} \sigma_1 & 0 & \dots \\ 0 & \dots & \dots \\ \vdots & \dots & \dots \\ 0 & \dots & \dots \end{bmatrix}_{M \times N} \cdot \begin{bmatrix} U_1 \cdot D \cdot V_1^T \end{bmatrix}_{N \times 1}$$

$$= \begin{bmatrix} 1 & 0 & \dots \\ 0 & \dots & \dots \\ \vdots & \dots & \dots \\ 0 & \dots & \dots \end{bmatrix}_{M \times M} \cdot \begin{bmatrix} \sigma_1 & 0 & \dots \\ 0 & \dots & \dots \\ \vdots & \dots & \dots \\ 0 & \dots & \dots \end{bmatrix}_{M \times N} \cdot \begin{bmatrix} 1 & 0 \\ 0 & V_1^T \end{bmatrix}_{N \times N}$$

diagonal!

$$\therefore U^T A V = U_1 \cdot D \cdot V_1^T$$

$$A = \underbrace{U \cdot U_1 \cdot D \cdot V_1^T \cdot V^T}_{\text{product of orthogonal matrices is orthogonal matrix.}}$$

of orthogonal
matrices is orthogonal matrix.

$$\therefore A = U^T \cdot D \cdot (V_1^T)^T$$

the singular value decomp.

~~$U_1 = A \cdot V_1$~~

$\sigma_i U_i = A V_1$

U_i is unit norm.

* Now; $A = U \cdot S \cdot V^T$ exists.

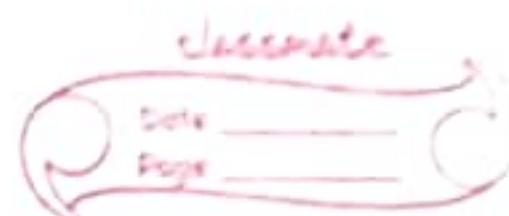
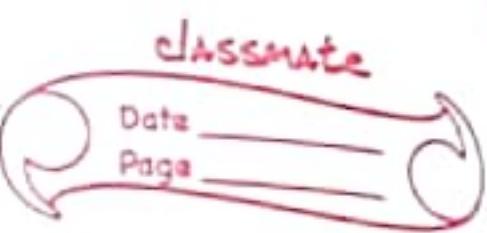
What is $A V_1$? $= \underline{\sigma_1 \cdot U_1}$

this is a column vector.

$\therefore A V_1$ is along U_1 , & orthogonal to all other column vectors of U .

→ orthogonal basis
in \mathbb{R}^n

$\xrightarrow{A_{m \times n}}$ orthogonal basis
in \mathbb{R}^m .
linear operator.



* properties:-

$$\|A\|_2 = \sigma_1 \text{ is unique.}$$

but; the right singular vector v_1

$$\text{S.t. } \|Av_1\| = \sigma_1$$

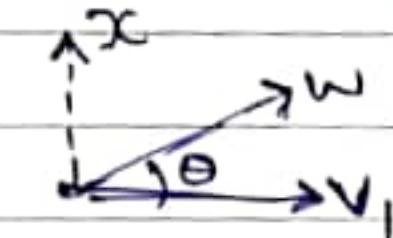
is it unique? (upto sign). $\|v_1\|_2 = 1$

let there be w ; such that $\|Aw\| = \sigma_1$,

$$\|w\| = 1$$

linearly independent with v_1 .

$x \rightarrow$ unit vector in direction
of $w - \langle w, v_1 \rangle v_1$:



$$\therefore \|Ax\| \leq \sigma_1 \text{ (by defn of } \sigma_1)$$

but we can show

$$\|Ax\| = \sigma_1$$

$$\begin{aligned} w &= (\cos\theta)v_1 + \sin\theta(x) \\ &= Cv_1 + Sx \end{aligned} \quad \begin{aligned} \sigma_1^2 &= \|Aw\|^2 = (Aw)^T(Aw) \\ &= (CAv_1 + SAx)^T(CAv_1 + SAx) \\ &= C^2\sigma_1^2 + S^2\|Ax\|^2 \\ &\quad + 2CS \cdot (Ax)^T(Aw) \\ &= C^2\sigma_1^2 + S^2\|Ax\|^2 \\ &\quad + 2CS \cdot (Ax)^T(Aw) \end{aligned}$$

if $x \perp v_1$. does $Ax \perp Av_1$?

$$\begin{aligned} A &= USV^T \cdot (Ax)^T(Av_1) \\ &= (USV^T \cdot x)^T(USV^T \cdot v_1) \\ &= (x^T \cdot V \cdot S^T \cdot S \cdot V^T \cdot v_1) \\ &= \left(\begin{bmatrix} 0 \\ y_1 \\ y_2 \\ \vdots \\ 1 \end{bmatrix}\right)^T S^T \cdot S \cdot \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \end{aligned}$$

$$= [0 \ z_1 \ z_2 \ \dots]^T \cdot \begin{bmatrix} \sigma_1 & & & \\ & 0 & & \\ & & 0 & \\ & & & 0 \end{bmatrix}_{N \times 1} = 0$$

* Hence; $x \perp v$, $\epsilon \|Av\| = \sigma_1$
 $\|Ax\| = \sigma_1$

$\therefore x$ must be right singular vector of A.
 with σ_1 as sing. value.

seen as:-

$$\begin{aligned}\sigma_1 &= \|Ax\| = \|Usv^T x\| \\ &= \|S \cdot v^T x\| \\ &= \|S \begin{bmatrix} 0 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{bmatrix}\|\end{aligned}$$

$$= \left\| \begin{bmatrix} \sigma_1 & 0 \\ 0 & B \end{bmatrix} \cdot \begin{bmatrix} 0 \\ y \end{bmatrix} \right\|$$

$$= \left\| \begin{bmatrix} 0 \\ B \cdot y \end{bmatrix} \right\|$$

$$= \|B \cdot y\|$$

if $\|B\| = \beta$

then $\exists y$ s.t

$$\|B \cdot y\| = \beta$$

construct.

$$y = \begin{bmatrix} 0 \\ y \end{bmatrix}$$

$$\text{then } |S \cdot y'| = \|By\|$$

$$\text{then } |S \cdot v^T (By')| = \sigma_1$$

$$\Rightarrow \|A \cdot (By')\| = \beta$$

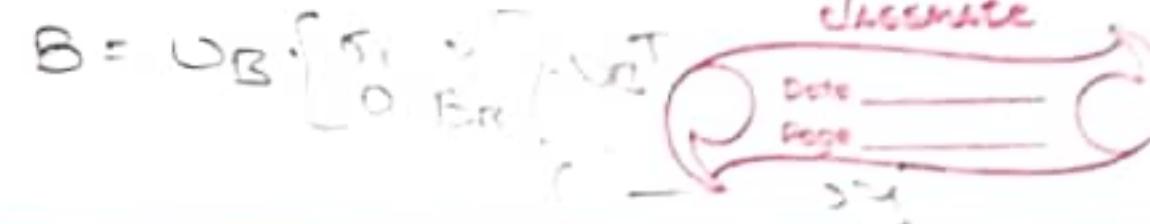
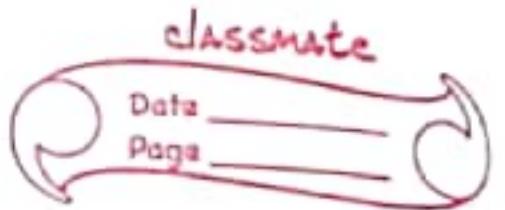
$$\therefore \|A\| \geq \|B\|$$

E: $\|B\|_2 = \sigma_1$ ($\because \|B\|_2 \leq \|A\|_2$)
 first sing. val.

& now we have =).

\therefore if y is the right sing. vector for B ,
 then this would transform to

$$\sqrt{\begin{bmatrix} 0 \\ y \end{bmatrix}} \text{ for } A$$



$$\text{then: } \begin{bmatrix} \sigma_1 & 0 \\ 0 & B \end{bmatrix} \cdot \begin{bmatrix} 0 \\ y \end{bmatrix} = \begin{bmatrix} \sigma_1 & 0 \\ 0 & B \cdot y \end{bmatrix}$$

* hence;

if $\sigma_1 = \|A\|_2$;

$$\begin{bmatrix} \sigma_1 & 0 \\ 0 & y \end{bmatrix}$$

(i) if v_i isn't unique (upto \pm), then
 the singular value σ_1 is repeated in v_i
 diagonal-S notation. not simple!

(ii) if σ_1 is simple (no multiplicity)
 then v_i is unique (upto sign).

* once σ_1, u_1, v_1 are determined;
 the remaining SVD is determined by
 linear operation A on space orthogonal to v_1 .

(again see
 $\|Ax\|_{\max}$
 \vdots)

\therefore we can build $U_1 V_1$.

\therefore if $x \perp v_1$

then $Ax \perp v_1$

$$\Rightarrow U_1 x \perp V_1$$

\therefore orthogonal matrices ...

* SVD in multigaussian:-

$$X = A \cdot W \quad M \leq N \text{ by definition.}$$

$$A = U \cdot S \cdot V^T \quad [\text{possible for any real } A]$$

$$\therefore X = U \cdot S \cdot V^T \cdot W \quad \begin{matrix} M \times N \\ N \times N \end{matrix}, \text{ where comp of } W \text{ are independent.}$$

$$= U \cdot S \cdot W' \quad \begin{matrix} M \times N \\ N \times N \end{matrix}, \text{ where comp. of } W' \text{ are also independent!}$$

$$X \approx (U \cdot S') \cdot W'' \quad \begin{matrix} M \times M \\ M \times N \end{matrix}, \text{ first } M \text{ columns of } S$$

$$X = A' \cdot W''$$

$$N > M \rightarrow N = M$$

$$\text{or: } X = W; A = I$$

covar(X)

$$= I \cdot I$$

\Rightarrow independent

$$X = V^T \cdot W; A = V^T$$

$$\text{covar}(X) = A \cdot A^T$$

$$= V^T \cdot V = I$$

\Rightarrow independent
congrats!

Structure of
an orthogonal
matrix

$$\langle u_i, u_j \rangle = 0 \text{ if } i \neq j$$

$$\text{Var}_t = \sum_{i=1}^P \sigma_i^2 \cdot p_i^2$$

$$= \sqrt{\sum p_i^2}$$

$$= \sqrt{I} \quad \because \text{orthogonal matrix.}$$

(ii) independence;

Structure of
an orthogonal
matrix

classmate
Date _____
Page _____

$$C_{old} = A \cdot A^T = M \cdot M = U S V^T \cdot V S V^T = U S S V^T$$

$$\begin{aligned} C_{new} &= U S V^T \cdot V S V^T \\ &= U S V^T \cdot V V^T \\ &= C_{old} \end{aligned}$$

Here $S \cdot S^T$ is a diagonal matrix.

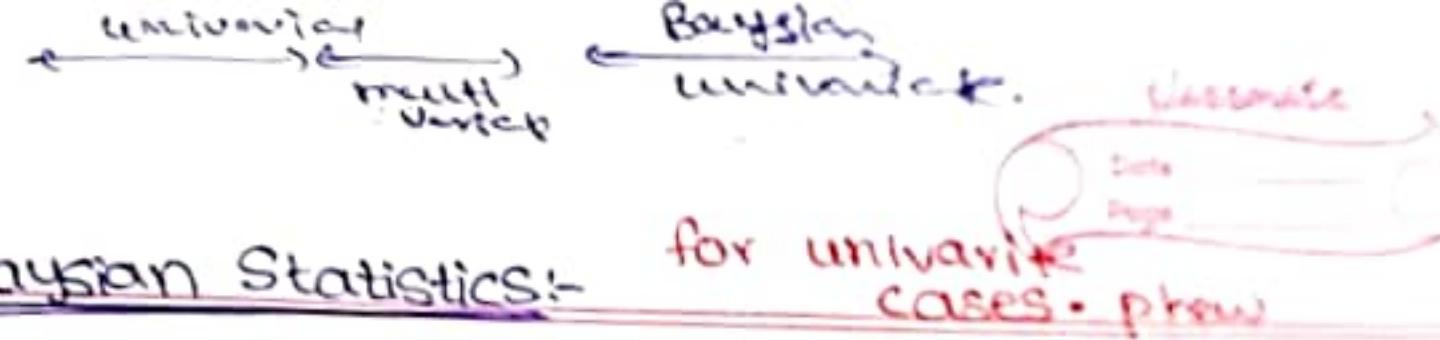
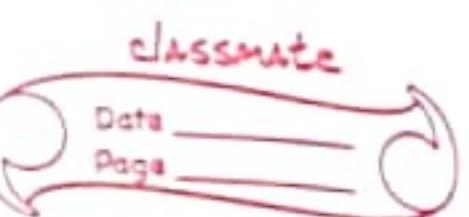
$$\therefore C = U \cdot (S S^T) \cdot V^T$$

\therefore $C = Q \cdot D \cdot Q^T$ form.

$S \cdot S^T$ is the eigen vector

value
of C .

= square of sing.
value of A .



Bayesian Statistics:-

for univariate
cases - p. no.

Bayes theorem: $P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$

↳ very imp. to statistics.

NOW;

$$P(\mu = \mu_0 / \text{data}) = \frac{P(\text{data} / \mu = \mu_0) \cdot P(\mu = \mu_0)}{P(\text{data})}$$

likelihood prior
evidence

$$= \int_{\mu=-\infty}^{\infty} P(\text{data} / \mu) f(\mu)$$

X = random var. (discrete)

modelling some unknown parameter.

Y = cont./disc. r.v. - modelling observed data.

1) Likelihood: $P(Y=y | X=x)$
- conditional

2) Evidence : $P(Y) = \sum_x P(Y=y, X=x)$

- marginal pdf. (All the ways, the data could

3) Prior : $P(X)$

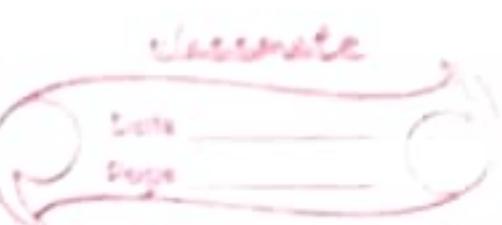
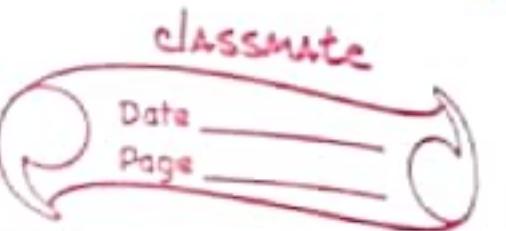
didn't
hardly
before.
was
uniform.

have been
generated
doesn't depend
on parameter
 X .

4) Posterior: $P(X=x | Y=y)$

- After data is observed, our beliefs are updated.

Improvement over MLE



* ML estimation:

- this assumes that prior info is not available, hence; uniform from $-\infty$ to ∞ :

(or) over the domain
 δx .

$$P(x=x_0 | \text{data}) \propto P(\text{data} | x=x_0)$$

likelihood function.

& we take the mode (argmax(LL)) of posterior function

as our estimate of x^*

* if we have prior belief (either a pdf) about x 's value, then

the probability of true value would be even amplified; benefiting us.

∴ "improvement over MLE":

"only when we take a good prior fu"
not a non-sense one."

Eg: data $\{x_i\}$ drawn from Gaussian known var. = σ^2
unknown mean = μ

prior belief: μ is drawn from $N(\mu_0, \sigma_0^2)$

∴ posterior

$$P(\mu | \text{data}) = P(\text{data} | \mu) \cdot P(\mu)$$

(const) normalizing.

$$P(\mu) = c_0 \cdot e^{-\frac{1}{2} \frac{(\mu - \mu_0)^2}{\sigma_0^2}}$$

$$P(\text{data} | \mu) = c_1 \cdot e^{-\frac{1}{2} \frac{\sum (x_i - \mu)^2}{\sigma^2}}$$

joint
pdf..

$$\text{LL function.} = c_0 \cdot e^{-\frac{1}{2} \frac{\sum (x_i - \mu)^2}{\sigma^2}}$$

$$= e^{-\frac{1}{2} \frac{\sum (x_i - \bar{x})^2 + (\bar{x} - \mu)^2}{\sigma^2/N}}$$

+!
constant.

$$\therefore P(\text{data} | \mu) \propto e^{-\frac{1}{2} \frac{(\mu - \bar{x})^2}{\sigma^2/N}}$$

$$\therefore P(\mu | \text{data}) = C \cdot e^{-\frac{1}{2} \frac{(\mu - \bar{x})^2}{\sigma^2}} \cdot e^{-\frac{1}{2} \frac{(\mu - \mu_0)^2}{\sigma_0^2}}$$

$$= e^{-\frac{1}{2} \frac{(\mu - \bar{x})^2}{\sigma^2}}$$

$$\hat{\mu} = \frac{\bar{x} \cdot \sigma_0^2 + \mu_0 \cdot \sigma^2 / N}{\sigma_0^2 + \sigma^2 / N}$$

mean / mode for gaussian.

the MAP
maximum
a posteriori
estimate.

NOW

$N \rightarrow \infty$: law of large nos. ignore prior.

$\sigma_0 \rightarrow \infty$: weak prior. see data

$\sigma_0 \rightarrow 0$: strong prior. see prior

→ Loss & risk function:-

if true value was Θ ; our estimator produces $\hat{\Theta}$;

we incur a loss' based on

loss function $L(\hat{\Theta}|\Theta)$

* if asked MAP; then posterior middle.

what if asked something else...

given $\{x_i\}_{i=1}^n$ from $P(x|\theta^*)$

posterior pdf

$$P(\theta | \{x_i\}) = \frac{P(\{x_i\} | \theta) \cdot P(\theta)}{P(\{x_i\})}$$

∴ Risk function

$R(\hat{\theta})$ = expectation of loss

$$= E_{P(\theta | \{x_i\})} L(\hat{\theta}, \theta)$$

Now, find $\hat{\theta}^*$ to minimize the MSD

$$\cdot E_{P(\theta | \{x_i\})} (\hat{\theta} - \theta)^2$$

could have been any

'risk function'

$$\text{would be } \hat{\theta}^* = E_{P(\theta | \{x_i\})} \theta$$

→ mean of Posterior PDF

= mode;

if Gaussian:

Goal! choose optimal $\hat{\theta}$ ($= \hat{\theta}^*$) to minimize risk function.

$R(\hat{\theta})$

again point estimate.

• we can choose $L(\hat{\theta}|\theta)$ appropriately;

such that

$$\hat{\theta} = \text{mode of } P(\theta | \{x_i\}) \quad L(\hat{\theta}|\theta) = I(\hat{\theta} \neq \theta)$$

(indicator fun.)

$$\hat{\theta} = \text{mean of } P(\theta | \{x_i\}) \quad L(\hat{\theta}|\theta) = (\hat{\theta} - \theta)^2 \frac{\text{MSD}}{P(\theta)}$$

$$\hat{\theta} = \text{median of } P(\theta | \{x_i\}) \quad L(\hat{\theta}|\theta) = |\hat{\theta} - \theta|_{\text{abs}}$$

order:-

- 1) take θ_{true}
- 2) take dataset X
- 3) do slope of LL at $\theta = \theta_{\text{true}}$.
- 4) expectation over all X .

classmate

Date _____

Page _____

→ Fisher Information:-

1) If likelihood function, $P(\text{data}|\theta)$ changes sharply

wrt change in θ , around $\theta = \theta_{\text{true}}$;

then easier to estimate θ_{true} .

for fisher info, of a parameter

classmate

Date _____

Page _____

expectation over all possible

observations.

∴ more fisher info \Rightarrow larger negatives of

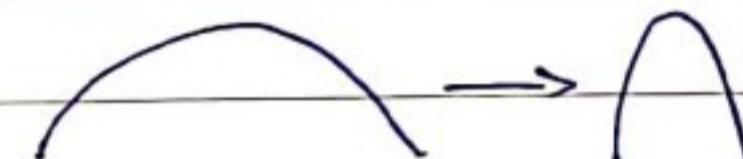
$$\frac{\partial^2}{\partial \theta^2} -$$

\Rightarrow the likelihood is more concave

2) If $P(\text{data}|\theta)$ has large spread wrt $\Delta\theta$, around θ_{true} ;

it will take many N-sized observations to get

a close estimate to θ_{true}



means; more certain about $\theta_{\text{true}}..$

1) $E_{P(\text{data}|\theta_{\text{true}})} \left[\frac{\partial}{\partial \theta} \log(P(\text{data}|\theta)) \Big|_{\theta_{\text{true}}} \right]$

$= 0$ expectation, over all possible observations.

* Gaussian; σ^2 known
 $\mu \rightarrow$ determine!

$$\text{then } P(x|\theta) = k \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

2) $E_{P(\text{data}|\theta_{\text{true}})} \left[\left(\frac{\partial}{\partial \theta} \log P(\text{data}|\theta) \right)^2 \Big|_{\theta_{\text{true}}} \right]$

$\neq 0$

$= I(\theta_{\text{true}})$

$$\therefore \log P(x|\theta) = -\frac{(x-\mu)^2}{2\sigma^2} + C.$$

$$\frac{\partial}{\partial \theta} (\quad) = \frac{(x-\mu)}{\sigma^2}$$

$$\frac{\partial^2}{\partial \theta^2} (\quad) = -\frac{1}{\sigma^2}$$

$$\therefore \boxed{I(\mu) = \frac{1}{\sigma^2}} \rightarrow \text{independent of } \mu!$$

position parameter

$$I(\theta) = E_{P(x|\theta)} \left(\frac{\partial}{\partial \theta} \log(P(x|\theta)) \right)^2 = -E_{P(x|\theta)} \frac{\partial^2}{\partial \theta^2} \log(P(x|\theta))$$

if N observations!

$$\boxed{I(\mu) = \frac{N}{\sigma^2}}$$

ANY unbiased estimator for θ .

Cramer-Rao Lower bound:

How good can an estimator for θ get?

Let's see
see
MSE for
estimator.

unbiased
estimators for now.

* Let $\hat{\theta}(x)$ be an unbiased estimator for θ .

$$\Rightarrow E_{P(x|\theta)} \left[(\hat{\theta}(x) - \theta) \right] = 0.$$

We show that;

$$\text{var}(\hat{\theta}(x)) \geq (I(\theta))^{-1}.$$

Same as MSE,
since unbiased.

* For gaussian mean; MLE is unbiased.

$$\text{var} = \left(\frac{1}{\sigma^2} \right)^{-1}; I(\theta) = \frac{1}{\sigma^2}$$

∴ MLE mean is unbiased, minimum variance estimator.

Bayesian Cramer-Rao Lower bound:-

for biased; unbiased estimators.

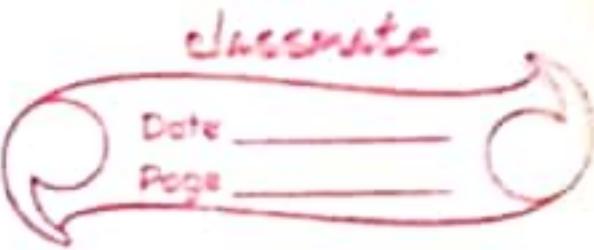
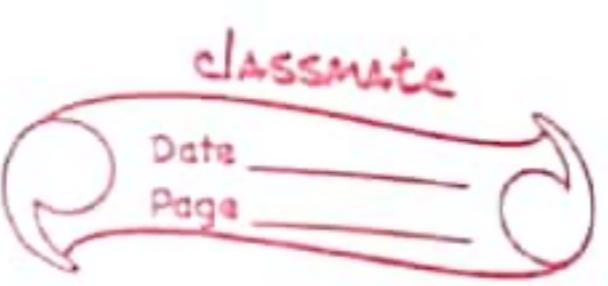
$$E_{Q(\theta|\alpha)} \left[E_{P(x|\theta)} (\hat{\theta}(x) - \theta) \right] \geq \left(E_{Q(\theta|\alpha)} (J_P[\theta]) + J_Q[\theta] \right)$$

MSE.

where $Q(\theta|\alpha)$ is a prior; $\alpha \rightarrow$ hyper parameter.

$$J_Q[\theta] = \left(E_{Q(\theta|\alpha)} \left(\frac{\partial}{\partial \theta} \log Q(\theta|\alpha) \right) \right)$$

not α .



→ Jeffry's prior:-

$$\text{take } Q(\theta) \propto \sqrt{I(\theta)}.$$

then; if your friend takes

$$Q(\beta) \propto \sqrt{I(\beta)}$$

$(\beta = g(\theta))$
inconformable

Both will be in sync!

∴ making the prior "non-informative".

$$Q(\beta) = Q_\theta(g^*(\beta)) \cdot \frac{\partial g^*(\beta)}{\partial \beta}$$

if $Q_\theta(\theta) \propto \sqrt{I(\theta)}$

then; we will get

$$Q(\beta) \propto \sqrt{I(\beta)},$$

same functional form

Harmony everywhere.