

Statistical Foundation for Machine Learning

Evaluation Pattern:

- 1. Assignments / Quiz (10)
- 2. Class Test - 2 (20)
- 3. Mid sem (30)
- 4. End sem (40)

Linear Algebra in Machine learning

x → vector (either $-$ or bold)

x → scalar

* $a^T x = b$ $a^T \rightarrow 1 \times n$ $x \rightarrow n \times 1$ & $b \Rightarrow 1 \times 1$

$T \rightarrow$ Transpose (ie converting rows into columns or vice versa)

Here a & x both are vectors / matrices

b is a scalar

* Array is a vector $1D \rightarrow 1 \times n$ & $2D \rightarrow m \times n$ (list of numbers)

* Linear Algebra is used throughout ML

Scalar

- single number (1×1)

Vector

- Array of numbers arranged in order
- Each identified by an index
- written in lower case
- Vectors as points in space
 - ↳ can be represented
- $1 \times n$ (1 axis)

Matrices

- 2D array
 - Denoted by bold typeface \mathbf{A}
 - If A has shape of height m & width n with real values then
 $A \in \mathbb{R}^{m \times n}$
- $R \rightarrow$ Real number
- $C \rightarrow$ complex

Tensor

- 2+ dimension array (more than 2 axes)
eg: RGB color images has 3 axes
Represented by $\underline{C}_{i,j,k}$
- It is an array of numbers arranged in a regular grid with variable no. of axes

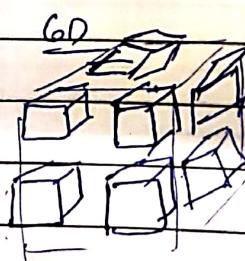
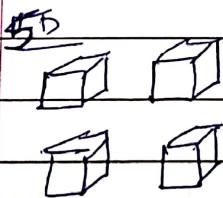
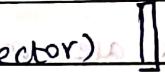
Shape of Tensors

1D \Rightarrow vector (1D-tensor = vector)

2D \Rightarrow Matrix (2D-tensor = matrix)

3D \Rightarrow Tensor (3D-tensor = tensor)

5D/6D / 4D - tensor \Rightarrow tensor



Eg: 3D-tensor = images

4D-tensor = video

* When we say 1D / 2D

D \rightarrow no. of dependent variables

Transpose of Matrix

→ Mirror image across diagonal

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{12} & a_{22} & a_{23} \\ a_{13} & a_{23} & a_{33} \end{bmatrix}$$

$$A^T = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

- vectors are special case of matrix
- scalar is a matrix with 1 element

Matrix Addition

- Add matrices to each other if they have same shape by adding corresponding elements
- If the matrices don't have same shape, we can add 0 to right

Matrix Multiplication

$$A_{m \times n} \times B_{n \times p} = C_{m \times p}$$

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 2 & 3 \\ 4 & 5 \end{bmatrix} = \begin{bmatrix} 10 & 13 \\ 22 & 29 \end{bmatrix}$$

- zeroes can be added to adjust shape of matrix & is called as up sampling

* Up sampling improves resolution by adding no. of samples

not product

- Also called Hadamard product
- Represented by $A \otimes B$
- Application: cosine similarity

Matrix Product Properties

- Distributivity over addition

$$\underline{A} (\underline{B} + \underline{C}) = \underline{AB} + \underline{AC}$$

- Associativity:

$$\underline{A} (\underline{B} \underline{C}) = (\underline{A} \underline{B}) \underline{C}$$

- Not commutative:

$$\underline{A} \underline{B} \neq \underline{B} \underline{A} \Rightarrow \text{Not true always}$$

• True in case of identity matrix

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 4 & 6 \\ 5 & 7 \end{bmatrix} = \begin{bmatrix} 14 & 20 \\ 27 & 46 \end{bmatrix}$$

- Dot product between vectors is

commutative

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

$$\begin{bmatrix} 4 & 6 \\ 5 & 7 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} = \begin{bmatrix} 22 & - \\ - & - \end{bmatrix}$$

$$\text{Eg: } \underline{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

$$\underline{x}^T = \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \quad \underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

$$(x_1 y_1 + x_2 y_2 + x_3 y_3)$$

$$\underline{y}^T = \begin{bmatrix} y_1 & y_2 & y_3 \end{bmatrix} \quad \underline{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \Rightarrow \begin{bmatrix} x_1 y_1 + x_2 y_2 + x_3 y_3 \end{bmatrix}$$



$$-(AB)^T = A^T B^T$$

Applications of Linear Algebra in Machine Learning

1. While designing a classifier (linear) using Neural Network.
eg:

- Consider there is a labeled data set which consists of 3 types of images - Dogs, cats & ships from which we derive weights & bias (2D tensor W_{1D})
- Now image being 2D/matrix is first converted into 1D i.e. into vector (for testing phase) & is applied weights & biases are applied to it to get its relevant class.

$$\text{Formula: } \underbrace{w^T x + b}_{\text{linear equation}} = \text{output}$$

linear equation

Linear Transformation

$$- Ax = b$$

$$- A \in \mathbb{R}^{n \times n} \text{ & } b \in \mathbb{R}^n$$

$$- A^{-1} A x = A^{-1} b$$

Identity & Inverse Matrices

$$- AA^{-1} = A^{-1} A = I$$

I dimension is obtained for A in case of $A^{-1} A$

$$- A^{-1} = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \quad A^{-1} = \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

There is an issue with this method of calculating inverse i.e. if determinant $(ad-bc)$ becomes 0.

linear equations: closed form solution
 → since calculation of A^{-1} has shortcoming hence this method is used by Gaussian elimination method followed by back substitution

Disadvantage:

- numerical instability due to division by small no.
- $O(n^3)$ for $n \times n$ matrix

How many solutions exist for $Ax = b$

1. Unique solution

$$2x + 3y = 9$$

$$9x + 2y = 11$$

x & y are unique

different slope & intercept \rightarrow unique
for both eqns

2. Infinite solution

$$2x + 3y = 4$$

$$4x + 6y = 8$$

Both the eqns are same & slope of eqn 1 = slope of eqn 2

Intercept of eqn 1 = Intercept of eqn 2

3. No solution

$$2x + 3y = 4$$

$$2x + 3y = 9$$

No solution since slope of eqn 1 = slope of eqn 2

Intercept of eqn 1 \neq Intercept of eqn 2

Span of set of vectors

- linear combination of set of vectors

$$Ax = \sum x_i A_i$$

- A linear combination of vectors $\{v^0, \dots, v^n\}$ if their coefficients c_i in $\sum c_i v_i$

- linearly independent of v

$$c_1 v_1 + c_2 v_2 + \dots + c_n v_n = 0$$

This is possible only if $c_1 = c_2 = \dots = c_n = 0$

Cosine similarity : Tells how close 2 vectors are

DATE _____

If 2 vectors are close or same, angle is 0

If 2 vectors are different, they're orthogonal ie angle is 90°

- Linearly Dependent

$$c_1 v_1 = -(c_2 v_2 + \dots + c_n v_n)$$

- Application of span of vectors : noise cancellation

Use of vector in Regression and Least square method

$$F_2 = f_2 + \epsilon_{22}$$

$$F_3 = f_3 + \epsilon_{33}$$

$$F_4 = f_4 + \epsilon_{44}$$

$$\begin{array}{|c|c|c|c|} \hline & f_2 & f_3 & f_4 \\ \hline F_1 & 2 & 0 & 0 \\ \hline F_2 & 2 & 0 & 0 \\ \hline F_3 & 2 & 1 & 2 \\ \hline \end{array}$$

$$EJ \leftarrow \frac{1}{2}(F_2 - 1)^2 + \frac{1}{2}(F_3 - 1)^2 + \frac{1}{2}(F_4 - 1)^2$$

$$\begin{array}{|c|c|c|c|} \hline & f_2 & f_3 & f_4 \\ \hline F_1 & 2 & 0 & 0 \\ \hline F_2 & 2 & 0 & 0 \\ \hline F_3 & 2 & 1 & 2 \\ \hline \end{array}$$

$$2F_2 - 2 \cdot 2 - 4 = 0$$

$$\begin{array}{|c|c|c|c|} \hline & f_2 & f_3 & f_4 \\ \hline F_1 & 2 & 0 & 0 \\ \hline F_2 & 2 & 0 & 0 \\ \hline F_3 & 2 & 1 & 2 \\ \hline \end{array}$$

$$\begin{array}{|c|c|c|c|} \hline & f_2 & f_3 & f_4 \\ \hline F_1 & 2 & 0 & 0 \\ \hline F_2 & 2 & 0 & 0 \\ \hline F_3 & 2 & 1 & 2 \\ \hline \end{array}$$

L^p Norm

$$\|\underline{x}\|_p = \left(\sum_i |x_i|^p \right)^{\frac{1}{p}}$$

ℓ^2 Norm (Euclidean norm):

- distance between origin & point \underline{x}
- same as euclidean distance

Use of norm in regression

1. Linear Regression

\underline{x} : a vector, \underline{w} : weight vector

$$y(\underline{x}, \underline{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\underline{x})$$

$$y(\underline{x}, \underline{w}) = w_0 + w_1 x_1 + \dots + w_d x_d \quad \underline{\phi} = \underline{w}^T \underline{x}$$

ϕ_j = Non linear basis fn

2. Loss function

$$\hat{E}(\underline{w}) = \frac{1}{2} \sum_{n=1}^N \{ y_n(\underline{x}_n, \underline{w}) - t_n \}^2 + c_2 \|\underline{w}\|^2$$

second term is a weighted norm, called as regularizer to prevent overfitting

Angle between Vectors

- dot product of 2 vectors can be written in terms of ℓ^2 norms & angle θ between them

$$\underline{x}^T \underline{y} \rightarrow \|\underline{x}\|_2 \|\underline{y}\|_2 \cos\theta$$

- cosine between 2 vectors is measure of similarity between them

efficiency of diagonal Matrix

- $\text{diag}(V) \underline{x} = V \odot \underline{x}$

- Easily calculable inverse matrix

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 9 & 0 \\ 0 & 0 & 3 \end{bmatrix} \quad A^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \frac{1}{9} & 0 \\ 0 & 0 & \frac{1}{3} \end{bmatrix}$$

Symmetric Matrix

- Transpose of symmetric matrix $=$ matrix of T same angle

$$A = AT$$

Eg: $A = \begin{bmatrix} 0 & 1 & 2 \\ 1 & 0 & 3 \\ 2 & 3 & 0 \end{bmatrix}$

Unit Vector

- A vector with unit norm

$$\|x\|_2 = 1$$

$$0 = V(IK - A)$$

Orthogonal Vector

$$x^T y = 0$$

If vectors have nonzero norm, vectors are 90° to each other

Orthonormal vectors

$$AT = A = AA^T = I$$

$$A^T = A^{-1}$$

- orthogonal & have unit norm

Eigen value decomposition / Matrix decomposition

- Decomposition of integer into prime factors
Eg: $12 \rightarrow 2 \times 2 \times 3$

- In case of square matrix, we follow Eigen value decomposition

$$A = V \text{diag}(\lambda) V^{-1}$$

$V \rightarrow$ eigen vectors (orthogonal matrix) i.e. $V^{-1} = V^T$ or $VV^{-1} = I$

$\lambda \rightarrow$ eigen values

scales vector A

eigen vector : It is a non zero vector v such that multiplication by A (original vector) only changes scale of v

$$Av = \lambda v$$

$\lambda \rightarrow$ scalar value / eigen value

If $A_{2 \times 2} \rightarrow$ we get 2 λ 's & λ can be +ve, -ve & 0

λ as 0 \rightarrow no vector is present for decomposition

$$Av = \lambda v \text{ can be written as } Av - \lambda v = 0$$

$$(A - \lambda I)v = 0$$

λ being scalar value hence by we multiply by I

$Av = \lambda v$ has a non zero solution if $|A - \lambda I| = 0$

$$|A - \lambda I| = (\lambda_1 - \lambda)(\lambda_2 - \lambda) \dots (\lambda_n - \lambda)$$

$n \rightarrow$ dimension of A

$$\sqrt{\left(\frac{1}{\sqrt{14}}\right)^2 + \left(\frac{2}{\sqrt{14}}\right)^2 + \left(\frac{3}{\sqrt{14}}\right)^2} = 1$$

Eg:

$$|A - \lambda I| = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}$$

$$= \begin{bmatrix} 2-\lambda & 1 \\ 1 & 2-\lambda \end{bmatrix} \Rightarrow (2-\lambda)^2 - 1$$



$$|A - \lambda I| = 3 - 4\lambda + \lambda^2$$

$$= \lambda^2 - 4\lambda + 3$$

$$= (\lambda - 3)(\lambda - 1) \Rightarrow (\lambda - 1)(\lambda - 3) = 0 \text{ since we assume it has non zero solution}$$

$$\lambda = 1, 3$$

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} v = \lambda v$$

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \lambda \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$2v_1 + v_2 = v_1 \Rightarrow v_1 + v_2 = 0$$

$$v_1 + 2v_2 = v_2$$

$$\therefore v_1 + v_2 = 0 \Rightarrow v_1 = -v_2$$

$$\text{If } v_1 = 1, v_2 = -1$$

when $\lambda = 1$

$$\lambda = 1 \quad \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

~~$$\lambda = 3 \quad \begin{bmatrix} 3 \\ -3 \end{bmatrix}$$~~

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 3v_1 \\ 3v_2 \end{bmatrix}$$

$$2v_1 + v_2 = 3v_1$$

$$v_1 + 2v_2 = 3v_2$$

$$-v_1 + v_2 = 0$$

$$v_1 = v_2$$

when $v_1 = 1$

$$\lambda = 3 \quad \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$



Eigen decomposition of A

$$A = V \text{diag}(\lambda) V^{-1}$$

$$V = \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$$

$$\text{diag}(\lambda) = \begin{bmatrix} 1 & 0 \\ 0 & 3 \end{bmatrix} \Rightarrow \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$$

$$V^{-1} = \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 3 \end{bmatrix} \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 3 \\ -1 & 3 \end{bmatrix} \begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix} = \begin{bmatrix} 4 & 2 \\ 2 & 4 \end{bmatrix} \not\equiv \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

$$\begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

$V \cdot V^T \neq I$
V is not orthogonal

$$V^T = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}$$

$$\therefore V \text{diag}(\lambda) V^T \neq A$$

Decomposition of symmetric matrix

If A is a symmetric matrix:

$$A A^T = I \Rightarrow \text{orthogonal}$$

$$A = V \text{diag}(\lambda) V^T$$

The eigen vectors in this case will be unit vectors

- Eigen decomposition is not unique ie λ_1 & λ_2 are same or all eigen values are same

Application / What does eigen decomposition tell us?

1. If $\lambda = 0 \Rightarrow$ Matrix is singular

Matrix can't be inverted or solved if matrix is singular



2. Useful to optimize quadratic expression of form
 $f(x) = \frac{1}{2} x^T A x \Rightarrow \|x\|_2 = 1$

eigen vector

- 3. A matrix whose eigen values are +ve \Rightarrow +ve definite
- 4. If all eigen values are -ve \Rightarrow -ve definite
- 5. If some eigen values are +ve & some are -ve or 0 \Rightarrow positive semi-definite

Singular value decomposition (SVD)

Shortcomings of EVD

- only defined for square matrix
- Requires calculation of inverse matrix which has $O(n^3)$ time

$$A = UDV^T$$

- U & V are orthogonal

- A is any matrix

$U \rightarrow$ left singular vector

$V \rightarrow$ right singular vector

If A is $m \times n$

$U \rightarrow m \times m$

$D \rightarrow m \times n$

$V \rightarrow n \times n$

Eg:

$$A = \begin{bmatrix} 4 & 0 \\ +3 & -5 \end{bmatrix}$$

$$A^T A = \begin{bmatrix} 25 & -15 \\ -15 & 25 \end{bmatrix}$$

$$|A^T A - \lambda I| = 0$$

$$\lambda = 10 \text{ & } 40$$

$$A^T A \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = 10 \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$\begin{bmatrix} 25 & -15 \\ -15 & 25 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 10v_1 \\ 10v_2 \end{bmatrix}$$

$$25v_1 - 15v_2 = 10v_1 \quad v_1 = v_2$$

$$15v_1 = 15v_2$$



$$V_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \lambda = 10 \quad \rightarrow V_{\lambda=10} \text{ after } L_2 \text{ norm} \Rightarrow \begin{bmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix}$$

When $\lambda = 40$

$$25V_1 - 15V_2 = 40V_1$$

$$-15V_1 = +15V_2$$

$$V_1 = -V_2$$

$$V_1 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

If $V = 1$ or $V = 1$ either
 Note: V_1 or V_2 can be taken & any value
 $\frac{1}{2}$ can be taken)

$$L_2 \text{ norm} = \sqrt{2}$$

$$V_{\lambda=40} = \begin{bmatrix} -1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix}$$

$$\text{Diagonal} = \sum \frac{\sigma^2}{\lambda} = \begin{bmatrix} \sqrt{40} & 0 \\ 0 & \sqrt{10} \end{bmatrix}$$

This should be in descending order of eigen values
 whereas we don't care in case of SVD

We used $\sqrt{\lambda}$ values since

$$A^T \cdot A = \|A\|^2 = \sum \sigma^2$$

Hence we take square root to take \sum
 we don't bother doing $\sqrt{\lambda}$ in SVD since eigen values are
 obtained directly from A .

$$V = \begin{bmatrix} -1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}$$

We need to place V first then V
 $\lambda = \text{highest nbr}$ $\lambda = \text{next highest nbr}$

$$A \cdot A^T = \begin{bmatrix} 4 & 0 \\ 3 & -5 \end{bmatrix} \begin{bmatrix} 4 & 3 \\ 0 & -5 \end{bmatrix} = \begin{bmatrix} 16 & 12 \\ 12 & 34 \end{bmatrix}$$

$$|AA^T - \lambda I| > 0$$

$$\begin{aligned} AA^T - \lambda I &= \begin{bmatrix} 16 & 12 \\ 12 & 34 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} \\ &= \begin{bmatrix} 16-\lambda & 12 \\ 12 & 34-\lambda \end{bmatrix} \\ &= (16-\lambda)(34-\lambda) - 144 \end{aligned}$$

$$\lambda_1 = 10, \quad \lambda_2 = 40$$

$$AA^T \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = 10 \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$\begin{bmatrix} 16 & 12 \\ 12 & 34 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 10v_1 \\ 10v_2 \end{bmatrix}$$

$$16v_1 + 12v_2 = 10v_1$$

$$6v_1 = -12v_2$$

$$v_1 = -2v_2$$

$$\text{Let } v_2 = 1 \Rightarrow v_1 = -2, v_2 = 1$$

$$v = \begin{bmatrix} -2\sqrt{5} \\ \sqrt{5} \end{bmatrix}$$

$$\begin{bmatrix} 16 & 12 \\ 12 & 34 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 40v_1 \\ 40v_2 \end{bmatrix}$$

$$16v_1 + 12v_2 = 40v_1$$

24

$$2v_1 = -12v_2$$

$$2v_1 = -v_2$$

$$\text{Let } v_1 = 1; v_2 = -2$$

$$v = \begin{bmatrix} \sqrt{5} \\ -2\sqrt{5} \end{bmatrix}$$

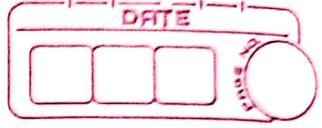
\therefore

$$U\Sigma V^T = \begin{bmatrix} 16 & 12 \\ 12 & 34 \end{bmatrix} \begin{bmatrix} \sqrt{40} & 0 \\ 0 & \sqrt{10} \end{bmatrix} \begin{bmatrix} 1/\sqrt{5} & -2/\sqrt{5} \\ -2/\sqrt{5} & 1/\sqrt{5} \end{bmatrix}$$

$$= \begin{bmatrix} 16\sqrt{40} & 12\sqrt{10} \\ 12\sqrt{40} & 34\sqrt{10} \end{bmatrix} \times V^T$$

$$= \begin{bmatrix} 16\sqrt{40} & 12\sqrt{10} \\ 12\sqrt{40} & 34\sqrt{10} \end{bmatrix} \begin{bmatrix} 1/\sqrt{5} & -2/\sqrt{5} \\ -2/\sqrt{5} & 1/\sqrt{5} \end{bmatrix}$$

$$= \begin{bmatrix} 32\sqrt{2} & -4\sqrt{2} \\ -4\sqrt{2} & 32\sqrt{2} \end{bmatrix}$$



Applications of SVD

① We can calculate inverse of matrix / Moore Penrose Pseudo inverse

$$A = UDV^T$$

$$A^+ = (UDV^T)^+ \rightarrow \text{pseudoinverse}$$

$$= \cancel{U}^{-1} \cancel{D}^+ (V^T)^+ + U^T - (K-N)(K-n) =$$

$$A^+ = V D^+ U^T$$

Since U & V are orthogonal

$$U^T U = I$$

$$V^T V = I$$

Using SVD, we can calculate A^+ very easily

② PCA or dimension reduction

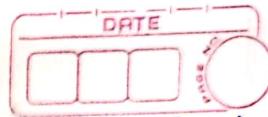
Book References

SVD - Pattern Recognition - Bishop (examples given)

Linear Algebra - Schaum's

8/9 → Quiz

(Till today's session)



→ Populus Book

Probability Theory in ML

- Quantify the uncertainty from 0-1
- Used in classification within ML
- Bayes theorem is used for evaluating probability in ML

Random Variables

- Enable us to associate numerical quantities with some experiment.
- Random variables are associated always for uncertainty
- Deterministic variable - variable whose outcome is already known.
The probability of such is always 1
- Random variables are independent to each other ie outcomes don't affect each other
- Random Process: associated with Random variable with time $x(n)$

Eg: E - Tossing a coin

$$U = \{H, T\}$$

Marginal Probability

Eg: consider there are 2 R.V. box & fruit

Box - Red, blue

Fruit - apple, orange

Marginal Probability : (Individual probability)

Probability of an apple (only 1 probability)

Eg: $P(F=a)$

conditional Probability:

Given we've orange, what is probability we choose blue box

Eg: $P(B=b | F=O)$

Joint Probability:

Probability of orange & blue box

Eg: $P(B=b, F=O)$

Joint

$$P(x=x_i, y=y_j) = \frac{n_{ij}}{N}$$

$$= \frac{n_{ij}}{c_i} \times \frac{c_i}{N}$$

$$P(y=y_j | x=x_i)$$

Joint conditional

Marginal probability

$$= P(y=y_j | x=x_i) \times P(x_i)$$

$$\therefore P(y=y_i | x=x_j) = \frac{P(x=x_i, y=y_j)}{P(x_i)}$$

↓

Bayes Theorem

$$P(x, y) = \frac{P(x|y) \cdot P(y)}{P(x)} = \frac{P(y|x)}{P(x)}$$

$$P(x|y) = \frac{P(y|x) \cdot P(x)}{P(y)} = \frac{P(x,y)}{P(y)}$$

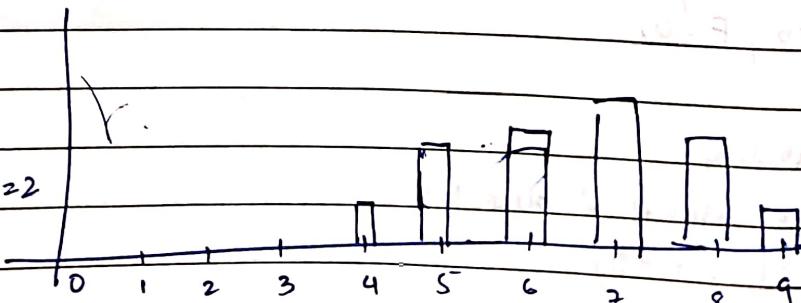
$$P(x) = \sum_y P(x, y)$$

$$= \sum_y \frac{P(x|y) \cdot P(y)}{P(y)} = \sum_x \frac{P(y|x)}{P(x)}$$

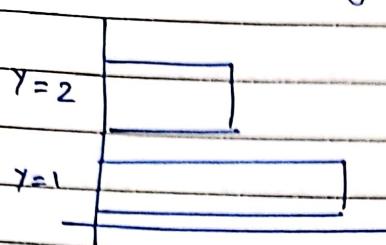
Joint distribution over 2 variables

Histogram of
x given y=2

$$P(x|y=2)$$



Histogram of γ / Marginal Probability of γ



$P(x, \gamma)$

$\gamma=2$
$\gamma=1$
0	•	•	•	•	•	•
1	•	•	•	•	•	•

Example: Find probability of box as red & orange using Bayes theorem
 $P(B=r | F=0)$

Given:

$$P(F=0 | B=r) = \frac{3}{4} \quad P(B=r) = \frac{1}{10} \quad P(F=0 | B=b) = \frac{1}{4} \quad \frac{12}{48}$$

$$P(F=a | B=r) = \frac{1}{4} \quad P(B=b) = \frac{6}{10} \quad P(F=a | B=b) = \frac{3}{4} \quad 4 =$$

$$P(B=r | F=0) = \frac{P(F=0 | B=r) \times P(B=r)}{P(F=0)} = \frac{\frac{3}{4} \times \frac{1}{10}}{\frac{12}{48}} = \frac{3}{48} \quad \begin{matrix} 3 - \text{apples} \\ 1 - \text{orange} \end{matrix}$$

$$P(F=0) = P(F=0, B=r) + P(F=0, B=b)$$

$$= P(F=0 | B=r) \times P(B=r) + P(F=0 | B=b) \times P(B=b)$$

$$= \frac{3}{4} \times \frac{1}{4} + \frac{1}{4} \times \frac{6}{10} = \frac{3}{16} + \frac{1}{4} = \frac{3}{4} + \frac{1}{4} = \frac{4}{4} = 1$$

$$= \frac{15 \times 6}{48} + \frac{5 \times 4}{48} = \frac{90 + 20}{48} = \frac{110}{48} = \frac{5}{24}$$

$$= P(F=0 | B=r) \times P(B=r) + P(F=0 | B=b) \times P(B=b)$$

$$= \frac{9}{20}$$

$$P(B=b | F=a) = ?$$

$$\frac{P(F=a | B=b) \times P(B=b)}{P(F=a)}$$

$$= \frac{3}{4} \times \frac{6}{10}$$

$$P(F=a)$$

$$P(F=a) = P(F=a | B=r) \times P(B=r) + P(F=a | B=b) \times P(B=b)$$

$$= \frac{1}{4} \times \frac{4}{10} + \frac{3}{4} \times \frac{6}{10}$$

$$= \frac{22}{40} = \frac{11}{20}$$

$$= \frac{3 \times 6}{3} = \frac{9}{12} = \frac{1}{4}$$

$$= \frac{11}{20} \cancel{\neq}$$

Independent Variables

$$P(X \neq Y) = P(X) \cdot P(Y)$$

$X \& Y$ are independent

$$P(Y|X) = \frac{P(X, Y)}{P(X)} = \frac{P(X) \cdot P(Y)}{P(X)} = P(Y)$$

Mutually Exclusive :

$$P(A \cap B) = \emptyset = \{ \} \quad P(A \cap B) = P(A) \cdot P(B)$$

Not independent variables are necessarily mutually exclusive



Sum, Product, Bayes for continuous Random variable

$$p(x) = \int p(x,y) dy \Rightarrow \text{Marginal}$$

$$p(x,y) = p(y|x) p(x) \Rightarrow \text{Joint}$$

$$p(y|x) = \frac{p(x,y)}{p(x)} = \frac{p(x|y) p(y)}{p(x)} \Rightarrow \text{conditional}$$

Expectation: It is average value of function $f(x)$ or x , represents weighted

- Denoted by E
- $E[x] = \sum_x x p(x)$ where weights are probabilities of each value
discrete R.V

Eg: x - 10 values from 0 to 9.

Each value is equiprobable (uniform distribution)

$$P(0) = P(1) = \dots = P(9) = \frac{1}{10}$$

$$\begin{aligned} E[x] &= (0)\left(\frac{1}{10}\right) + (1)\left(\frac{1}{10}\right) + 2\left(\frac{1}{10}\right) + 3\left(\frac{1}{10}\right) + 4\left(\frac{1}{10}\right) + 5\left(\frac{1}{10}\right) + 6\left(\frac{1}{10}\right) \\ &\quad + 7\left(\frac{1}{10}\right) + 8\left(\frac{1}{10}\right) + 9\left(\frac{1}{10}\right) \end{aligned}$$

$$= \frac{45}{10} = \underline{\underline{4.5}}$$

* Expectation of a deterministic signal remains constant

$$- E[f(x)] = \sum_x f(x) p(x)$$

Eg: $f(x) = x^2$

- If $f(x)$ is continuous

$$E[f(x)] = \int f(x) p(x) dx$$

- If x is continuous

$$E[x] = \int x p(x) dx$$



★ conditional expectation?



conditional expectation:

$$E[f] = \sum_x p(x|y) f(x)$$

Variance:

- Define variability in $f(x)$ around its mean value

$$\text{var}[f] = E[(f(x) - E[f(x)])^2]$$

$$= E[(f(x))^2 - 2f(x)E[f(x)] + (E[f(x)])^2]$$

$$= E[f(x)^2] - 2E[f(x)]E[f(x)]$$

$$= E[f(x)^2] - \underbrace{E[2f(x)E[f(x)]]}_{\text{cancel}} + E[E[f(x)]^2]$$

This will be a constant

$$= E[f(x)^2] - 2E[f(x)]^2 + E[f(x)]^2$$

$$= E[f(x)^2] - 2E[f(x)]^2$$

$$= (E[x]^2) - 2(E[x])^2$$

Covariance:

relation

- Define between 2 different random variables

$$\text{cov}[x, y] = E_{x,y} [\{x - E[x]\} \{y - E[y]\}]$$

$$= E_{x,y} [xy - xE[y] - yE[x] + E[x]E[y]]$$

$$= E_{x,y} [xy] + E_{x,y} [-xE[y] - yE[x]] + E_{x,y} [E[x]E[y]]$$

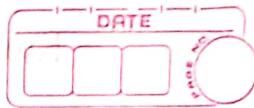
$$= E_{x,y} [xy] - E[x]E[y] - E[x]E[y] + E[x]E[y]$$

$$= E_{x,y} [xy] - E[x]E[y]$$

$$E_{x,y} [-xE[y]] \Rightarrow -E_{x,y} [x] E_{x,y} [E[y]]$$

$E[y] \Rightarrow$ since $E[y]$ is an expectation

$-E[x] E[y]$ & expectation on expectation
remains same & stays same



Bayesian Probability

- Probability is quantification of uncertainty
- Likelihood - observe the data

Priors - Having predetermined knowledge

Posterior - final answer which is correct

Properties of probability

- No probability can be > 1 i.e. $P(\text{sample space}) = 1$
- $P(\text{Exclude Any event}) \geq 0 \rightarrow P(E_1 \cup E_2) = P(E_1) + P(E_2) - P(E_1 \cap E_2)$
- Bayes Posterior \propto Likelihood \times Prior \Rightarrow Bayes theorem
- Sum rule

$$P(D) = \int_{\text{continuous}} P(D|w) P(w) dw$$

continuous

Since $P(D|w)$ & $P(w)$ are continuous, hence \int

Gaussian Distribution

$$* N(x|\mu, \sigma^2) = \frac{1}{\sqrt{(2\pi\sigma^2)^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} \quad \mu \Rightarrow \text{Mean}$$

$$\Downarrow \quad \sigma^2 \Rightarrow \text{variance}$$

single variate or x has 1 feature

$\sigma \Rightarrow$ S.D.

Expectation of x under Gaussian

$$E[x] = \int_{-\infty}^{\infty} N(x|\mu, \sigma^2) x dx$$

$$E[x^2] = \int N(x|\mu, \sigma^2) x^2 dx = \mu^2 + \sigma^2$$

$$\text{var}[x] = E[x]^2 - E[x^2]$$

- * If x is multi dimensional

$$N(x|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \times e^{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)}$$

$\mu \rightarrow$ Mean

$\Sigma \rightarrow$ covariance matrix $= D \times D$

D-dimensional vector



$$p(x|u, \sigma^2) = \prod_{n=1}^N N(x_n|u, \sigma^2)$$

Prove: $\log p(x|u, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - u)^2 - \frac{N}{2} \log \sigma^2 - \frac{N}{2} \log(2\pi)$

$$p(x|u, \sigma^2) = \prod_{n=1}^N \frac{1}{(2\pi\sigma^2)^{1/2}} \times e^{\left(\frac{-1}{2\sigma^2}(x-u)^2\right)}$$

Taking log on both sides

$$\log p(x|u, \sigma^2) = \log \left(\prod_{n=1}^N \frac{1}{(2\pi\sigma^2)^{1/2}} \times e^{\left(\frac{-1}{2\sigma^2}(x-u)^2\right)} \right)$$

$$= -\frac{1}{2\sigma^2} \sum_{n=1}^N (x-u)^2$$

$$\log p(x|u, \sigma^2) = \log \left(\prod_{n=1}^N \frac{1}{(2\pi\sigma^2)^{1/2}} \times e^{\left(\frac{-1}{2\sigma^2}(x-u)^2\right)} \right)$$

By property $\log(ab) = \log a + \log b$

$$= \log \left(\prod_{n=1}^N \frac{1}{(2\pi\sigma^2)^{1/2}} \right) + \left(-\frac{1}{2\sigma^2} \sum_{n=1}^N (x-u)^2 \right)$$

$$= -\frac{1}{2\sigma^2} \left(\sum_{n=1}^N (x_n - u)^2 \right) - \log \left(\prod_{n=1}^N (2\pi\sigma^2)^{1/2} \right)$$

(By property of $\log(\frac{1}{a}) = -\log a$)

$$= -\frac{1}{2\sigma^2} \left(\sum_{n=1}^N (x_n - u)^2 \right) - \frac{N}{2} \log(2\pi) - \frac{N}{2} \log \sigma^2$$



Independent events

$$P(E_1 \cap E_2) = P(E_1) \cdot P(E_2) = P(E_1, E_2)$$

$$\boxed{P(E_2 | E_1) = P(E_1) \cdot P(E_2) / P(E_1) = P(E_2)}$$

- 2 mutually exclusive events are not independent

$$E_1 \cap E_2 = \emptyset$$

$$P(E_1 \cap E_2) = 0$$

$$\text{But } P(E_1 \cap E_2) = P(E_1) \cdot P(E_2) \neq 0 \Rightarrow P(E_1 \cup E_2) = P(E_1) + P(E_2)$$

Independent event condition

- 2 identical events are not independent

$$P(E_1 \cap E_2) = P(E_1)$$

But by independent event defn,

$$P(E_1 \cap E_2) = P(E_1) \cdot P(E_2) \quad \left\{ P(E_1 \cap E_2) = P(E_1, E_2) \right\}$$

$$\text{But } P(E_1) \cdot P(E_2) \neq P(E_1)$$

$$\text{Since } P(E_2) \neq 1 \quad P(E_2) < 1$$

- If E_1 & E_2 are independent, are their complements too?

$$\begin{aligned} P(E_1' \cap E_2') &= 1 - P(E_1 \cup E_2) \\ &= 1 - P(E_1) - P(E_2) + P(E_1 \cap E_2) \\ &= 1 - P(E_1) - P(E_2) + P(E_1) \cdot P(E_2) \\ &= (1 - P(E_1))(1 - P(E_2)) \\ &= P(E_1') \cdot P(E_2') \end{aligned}$$

which means they're independent

Example

$$P(\text{clear, rising}) = P(\text{rising} | \text{clear}) \times P(\text{clear})$$

$$P(z_1, z_2) = P(z_1 | z_2) \times P(z_2)$$

$$P(\text{clear, rising}) = P(\text{clear} | \text{rising}) \times P(\text{rising})$$

$$P(\text{dry, clear, rising}) = \frac{P(\text{dry, clear, rising})}{P(\text{clear, rising})} = \frac{0.40}{0.47} \approx 0.85$$

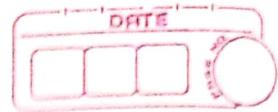
$$P(z_1) = \sum P(z_1 | z_3) \times P(z_3) + P(z_1)$$

$$0.07 + 0.1 + 0.11 = 0.28$$

$$\frac{0.11}{0.12} = \underline{0.33}$$

$$0.4 + 0.08 + 0.09 + 0.03 = 0.5$$

$$0.48 = \underline{0.12}$$



$P(\text{cloudy}) = P(\text{cloudy} | \text{dry}) P(\text{dry}) + P(\text{cloudy} | \text{wet}) P(\text{wet})$

≈ 0.28

Cloudy, wet $P(\text{wet})$

$$\frac{P(\text{wet})}{P(\text{wet})}$$

* $P(E_1 | E_2) \neq P(E_2 | E_1) \Rightarrow$ Inverse of probability

0.47

$$P(\text{clear} | \text{rising}) = P(\text{rising, clear}) / P(\text{rising})$$

$$= \frac{0.47}{0.67} = 0.7$$

0.11

0.09

0.2

* $P(\text{rising} \cup \text{falling}) = P(\text{rising}) + P(\text{falling}) = 1$

$$P(\text{Rising} | \text{clear}) + P(\text{Falling} | \text{clear}) = 1$$

$$P(\text{clear} | \text{falling}) + P(\text{clear} | \text{rising}) \neq 1$$

* $P(\cdot | E_1) = P(E_2 \cup E_3 \cup E_4 | E_1)$ (since there are 4 events E_1, E_2, E_3, E_4)

Any event

$$P(E_2 \cup E_3 \cup E_4 | E_1) = P(E_2 | E_1) + P(E_3 | E_1) + P(E_4 | E_1)$$

$E_2, E_3, E_4 \Rightarrow$ disjoint

* $P(E_1 | E_2 \cup E_3 \cup E_4) \neq P(E_1 | E_2) + P(E_1 | E_3) + P(E_1 | E_4)$

\Downarrow

$$P(E_1 | \cdot)$$

* If \bar{E} is complement of E

$$P(\bar{E}_1 | E_2) = 1 - P(E_1 | E_2)$$

but

$$P(E_1 | \bar{E}_2) \neq 1 - P(E_1 | E_2)$$

$$\frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2} - \frac{1}{2} (\log(\frac{1}{2} \cdot \frac{1}{2})) = -\log \left(\frac{1}{2}\right) = \log_2 e^{0.9}$$

$$P(\text{Positive} | \text{sick}) = P(\text{positive, sick}) / P(\text{sick}) = 0.009 / 0.01 = 0.9$$

$$P(- | \text{healthy}) = 0.891 / 0.99 = 0.9$$

$$P(\text{sick} | +) = 0.009 / 0.108 = 0.0833$$

$$P(\text{Healthy} | -) = 0.891 / 0.891 = 1$$

Entropy

- Entropy is measure of randomness of random variable Z .
- Represented by $H(Z)$, for single event random variable

$$H(Z) = - \sum_{z \in Z} p_z \log_2 p_z$$

- Entropy is maximum if $p(z) = 1$ ie equiprobable for all values of Z

Eg: Sun rises can rise from any direction / Toss of coin it can be Head or Tail

- It is minimum : $p(z) = 1$ for a value of Z (others are null)

Eg: Sun rises in east, there is no uncertainty

- If there several random variables it is called as ensemble of entropy

$$\sum_i p_i = 1 \rightarrow \text{Marginal probability}$$

$$H = -1 = - \sum_i p_i \log_2 p_i$$

- Entropy of an ensemble is simply average of all elements in it it can be computed by their average entropy by weighing each of $\log p_i$ contributions by its probability p_i occurring.

- In case of several events entropy is multiplied by weight/probability in ensemble but not in case of entropy with Example:

Let Z be an alphabet from A to Z but all prob of all alphabets is equiprobable

$$H = - \left[\frac{1}{26} \times 26 \log_2 \frac{1}{26} \right]$$

$$= 4.7$$



$$\begin{aligned}
 H &= - \left[0.165 \log_2 0.165 + 0.072 \log_2 0.072 + 0.066 \log_2 0.066 + \right. \\
 &\quad 0.063 \log_2 0.063 + 0.059 \log_2 0.059 + 0.055 \log_2 0.055 + \\
 &\quad 0.054 \log_2 0.054 + 0.052 \log_2 0.052 + 0.047 \log_2 0.047 + \\
 &\quad \left. 0.035 \log_2 0.035 + 0.029 \log_2 0.029 + 0.023 \log_2 0.023 \right] \\
 &= - \left[-0.105 \times 3.2515 - 0.072(3.8) - 0.066(3.92) - \right. \\
 &\quad 0.063(3.9) - 0.059(4.08) - 0.055(4.16) - \\
 &\quad 0.054(4.21) - 0.052(4.26) - 0.047(4.4) - \\
 &\quad \left. 0.035(4.8) - 0.029(5.1) - 0.023(5.42) \right] \\
 H &= 2.69
 \end{aligned}$$

Conditional Entropy:

Joint Entropy:

- $H = - \sum p(x,y) \log_2 p(x,y)$

$$H = \sum_{x,y} p(x,y) \log_2 \frac{1}{p(x,y)}$$

- Joint entropy is additive

$$H(X,Y) = H(X) + H(Y) \quad p(x,y) = p(x) \cdot p(y)$$

Given: $H(X,Y) = \sum_{x,y} p(x) \cdot p(y) \log_2 \left(\frac{1}{p(x) \cdot p(y)} \right)$

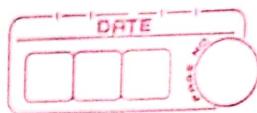
$$= - \sum_{x,y} p(x) \cdot p(y) \log_2 (p(x) \cdot p(y))$$

$$= - \sum_{x,y} p(x) \cdot p(y) [\log_2 p(x) + \log_2 p(y)]$$

$$= - \sum_{x,y} p(x) \cdot p(y) \log_2 p(x) + p(x) \cdot p(y) \log_2 p(y)$$

$$H(X) = - \sum_x p(x) \log_2 p(x) \quad H(Y) = - \sum_y p(y) \log_2 p(y)$$

$$H(X) + H(Y) = - \left[\sum_x p(x) \log_2 p(x) + \sum_y p(y) \log_2 p(y) \right]$$



$$\begin{aligned}
 &= - \sum_x P(x) \log_2 P(x) \sum_y P(y) - \sum_y P(y) \log_2 P(y) \sum_x P(x) \\
 \sum_x P(x) = 1 \quad &\& \sum_y P(y) = 1 \\
 &= - \left[\sum_x P(x) \log_2 P(x) + \sum_y P(y) \log_2 P(y) \right] \\
 &= H(X) + H(Y)
 \end{aligned}$$

Conditional Entropy

Measures uncertainty about random variable X when r.v Y has taken

$$H(X|Y=b_j) = \sum P(x|Y=b_j) \log_2 \frac{1}{P(x|Y=b_j)}$$

some random value
↳ This is for single event $P(x|Y=b_j)$

$$\begin{aligned}
 H(X|Y) &= \sum_y P(y) \left[\sum_x P(x|Y) \log_2 \frac{1}{P(x|Y)} \right] \\
 &= \sum_{x,y} P(y) \cdot P(x|Y) \cdot \cancel{\log_2 \frac{1}{P(x|Y)}} \\
 &= \sum_{x,y} P(x,y) \log_2 \frac{1}{P(x|y)}
 \end{aligned}$$

$$P(x|y) = \frac{P(x,y)}{P(y)} \quad \text{or} \quad P(y|x) \cdot P(x) = P(x,y) \quad \Rightarrow \quad P(x) = \frac{P(x,y)}{P(y|x)}$$

Chain Rule for Entropy

$$\text{Prove: } H(X,Y) = H(X) + H(Y|X) = H(Y) + H(X|Y)$$

$$\begin{aligned}
 H(X) + H(Y|X) &= - \left[\sum_x P(x) \log_2 P(x) + \sum_{x,y} P(x,y) \log_2 \frac{P(y|x)}{P(x)} \right] \\
 &= - \cancel{\sum_x P(x)} \cdot \cancel{\sum_{x,y} P(x,y) \log_2 P(x)} \\
 &= - \left[\sum_x P(x) \log_2 P(x) \right]
 \end{aligned}$$

$$\begin{aligned}
 H(X, Y) &= - \sum_{x,y} P(x, y) \log_2 P(x, y) \\
 &= - \sum_{x,y} P(Y|x) P(x) \log_2 [P(Y|x) P(x)] \\
 &= - \sum_{x,y} P(Y|x) P(x) [\log_2 P(Y|x) + \log_2 P(x)] \\
 &= - \left[\sum_{x,y} P(Y|x) \cdot P(x) \log_2 P(Y|x) + \sum_{x,y} P(Y|x) P(x) \cdot \log_2 P(x) \right] \\
 &= - \left[\sum_{x,y} P(Y|x) \cdot P(x) \log_2 P(Y|x) - \sum_x P(x) \log_2 P(x) \sum_y P(Y|x) \right] \\
 &= - \sum_{x,y} P(x, y) \log_2 P(Y|x) + H(X) \sum_y P(Y|x) \\
 &\quad \text{1 since its summation of } P(Y|x) \text{ & } x \text{ is already given} \\
 &= - \sum_{x,y} P(x, y) \log_2 P(Y|x) + H(X) \\
 &\quad - H(Y|X) + H(X)
 \end{aligned}$$

Axiomatic



Prove

$$H(x_1, x_2, \dots, x_n) \leq \sum_{i=1}^n H(x_i)$$

$$H(X) = -\sum_x p(x) \log p(x)$$

$$\sum_{i=1}^n H(x_i) = \sum_{i=1}^n -\sum_x p(x_i) \log p(x_i)$$

$$H(x_1, x_2, \dots, x_n) = H(x_1) + H(x_2|x_1) + H(x_3|x_1, x_2) + \dots +$$

$H(x_n|x_1, x_2, \dots, x_{n-1}) \Rightarrow$ By chain rule for entropy

$$\text{By formula, } H(x, y) = -\sum_{x,y} p(x, y) \log p(x, y)$$

$$\therefore H(x_1, x_2, \dots, x_n) = -\sum_x p(x_1, \dots, x_n) \log p(x_1, \dots, x_n)$$

Each conditional entropy $H(x_i|x_1, \dots, x_{i-1})$ is non negative, indicating
 ~~$H(x_1, x_2, \dots, x_n) \leq H(x_i)$~~ that joint entropy will not
 exceed marginal total entropy of individual components. This
 follows the rule that $H(x|y) \leq H(x)$ as existence of known value

Hence,

$$H(x_1, x_2, \dots, x_n) \leq \sum_{i=1}^n H(x_i)$$

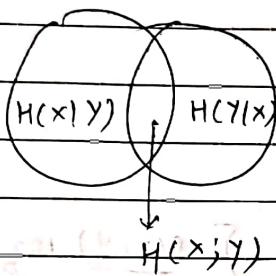
for y tend to reduce randomness
in x



Mutual Information between X & Y

Mutual info between 2 r.v measures amount of information that one conveys about other. Equivalently measures avg reduction in uncertainty about X that results from learning Y

$$H(X;Y) = \sum_{x,y} P(x,y) \log_2 \frac{P(x,y)}{P(x) \cdot P(y)}$$



- When X & Y are independent, mutual entropy is 0

$$I(X;Y) = \sum_{x,y} P(x,y) \log_2 \frac{P(x,y)}{P(x) \cdot P(y)}$$

↓ 1

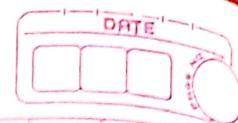
$$= 0 \quad \{ \log_2 1 = 0 \}$$

Reactive Entropy Mutual Info:

- It is always non negative or ≥ 0

$$P(X,X) = P(X)$$

$$\begin{aligned} I(X;X) &= \sum_x P(x) \log_2 \frac{P(x)}{P(x) \cdot P(x)} \\ &= \sum_x P(x) \log_2 \frac{1}{P(x)} \\ &= - \sum_x P(x) \log_2 P(x) \\ &= H(X) \end{aligned}$$



①

$$\begin{aligned} I(X;Y) &= H(X) - H(X|Y) \\ &= I(Y;X) \\ &= H(X) + H(Y) - H(X,Y) \rightarrow ② \end{aligned}$$

Prove:

$$\begin{aligned} ① \quad I(X;Y) &= H(X) - H(X|Y) \\ &= - \sum_x P(x) \log_2 P(x) - \left[- \sum_{x,y} P(x|y) \log_2 P(x|y) \right] \end{aligned}$$

By Bayes theorem,

$$\begin{aligned} P(x|y) &= \frac{P(x,y)}{P(y)} \\ &= + \sum_x P(x) \log_2 \frac{1}{P(x)} + \sum_{x,y} \frac{P(x,y)}{P(x)P(y)} \log_2 \frac{P(x|y)}{\frac{P(x,y)}{P(x)P(y)}} \\ &= \sum_x P(x) \log_2 \frac{1}{P(x)} + \sum_{x,y} P(x,y) \log_2 \left(\frac{P(x,y)}{P(y)} \right) \\ &= \sum_x P(x) \sum_{x,y} \log_2 \frac{1}{P(x)} + \sum_{x,y} P(x,y) \log_2 \left(\frac{P(x,y)}{P(y)} \right) \\ &= \sum_x P(x) \log_2 \frac{1}{P(x)} + \sum_{x,y} P(x,y) \log_2 \left(\frac{P(x,y)}{P(x)} \right) \end{aligned}$$

$$\begin{aligned} I(X;Y) &= \sum_{x,y} P(x,y) \log_2 \frac{P(x,y)}{P(x)P(y)} \\ &= \sum P(x|y) \cdot P(y) \log_2 \frac{P(x|y) \cdot P(y)}{P(x) \cdot P(y)} \\ &= \sum P(x|y) \cdot P(y) \log_2 \frac{P(x|y)}{P(x)} \\ &\stackrel{2}{=} \sum P(x|y) \cdot P(y) \log_2 P(x|y) - \sum P(x|y) P(y) \log_2 P(x) \\ &\stackrel{3}{=} \left[\sum P(x,y) \log_2 P(x|y) \right] + \left[\sum P(x|y) P(y) \log_2 P(x) \right] \end{aligned}$$

$$= \left[-H(X|Y) + H(X) \right]$$

$$= H(X) - \underline{H(X|Y)}$$

$$= H(X) - H(X|Y)$$

$$= H(X) - \left[-H(Y) + H(X,Y) \right]$$

$$= \underline{H(X) + H(Y) - H(X,Y)}$$

additivity of entropy result

Spanning of $\{A, B\}$

(1, 2) (A, B) are independent

$(X, Y) = (Y, X)H = (Y, X)I$

mutual information in example 5 $\leq (X, Y)$

0

$I(X; Y) = 0$

$I(X; Y)$

mutual info. $I(X; Y) \leq I(X, Y)$

value of mutual information $I(X; Y) \geq I(X, Y) \geq (X, Y)$

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

1 1 1 1 1 1 1 1 1 1

Cross Entropy

If there are 2 distributions $p(x)$ & $q(x)$ over same set of outcomes

$$H(p, q) = - \sum_x p(x) \log q(x)$$

If $p(x) = q(x)$

$$H(p) = \sum_x p(x) \log p(x)$$

- Cross entropy is asymmetric

$$H(p, q) \neq H(q, p)$$

Distance between x & y (R.V)

$$D(x, y) = H(x, y) - I(x; y)$$

$$D(x, y) \geq 0 \quad \because \text{distances are non-negative}$$

$$D(x, x) = 0$$

$$D(x, y) = D(y, x) \Rightarrow \text{symmetric}$$

$$D(x, z) \leq D(x, y) + D(y, z) \quad (\text{Triangle Inequality rule})$$

KL divergence

$$D_{KL}(p || q) = \sum_x p(x) \log \frac{p(x)}{q(x)} \quad \text{Here } p(x) \& q(x) \text{ are diff probability distributions defined over } x$$

Application:

A model say $x \rightarrow y_x$ is a training data has a model, this model cannot be applied on $z \rightarrow y_z$ if $p(x)$ & $q(z)$ are very different. They can be determined as different using KL distance between them & we can use KL divergence

$$D_{KL}(p || q) = H(p, q) - H(p)$$

If variables are independent, they are always uncorrelated

If variables are uncorrelated, not necessary are independent under condition variables

Correlation & covariance Matrix & Wiener Filter

follows Gaussian joint distribution

1. Hermitian Matrix

$$C^H = C$$

? what is conjugate?

$H \rightarrow$ Hermitian operator

\Leftrightarrow conjugate + transpose

$C \rightarrow$ conjugate matrix

$$y = a + jb \Rightarrow y^* = a - jb$$

2. Covariance Matrix

$$\checkmark \quad C_{p \times p} = \begin{bmatrix} (x_1 - \bar{x}_x) & [(x_1 - \bar{x}_x), \dots, (x_p - \bar{x}_x)]^* \\ (x_2 - \bar{x}_x) & [H(x_1 - \bar{x}_x), \dots, (x_p - \bar{x}_x)]^* \\ \vdots & \vdots \\ (x_p - \bar{x}_x) & [H(x_1 - \bar{x}_x), \dots, (x_p - \bar{x}_x)]^* \end{bmatrix}_{p \times p}$$

Auto covariance
matrix

$$= \begin{bmatrix} c_{xx_1} & c_{xx_2} & \dots & c_{xp} \\ c_{xx_2} & c_{xx_2} & \dots & c_{xp} \\ \vdots & \vdots & \ddots & \vdots \\ c_{xp} & c_{xp} & \dots & c_{xp} \end{bmatrix}_{p \times p}$$

Prove : covariance matrix is always symmetric

$$c = E[(\underline{x} - \bar{x}_x)(\underline{x} - \bar{x}_x)^H]$$

Given $c^H = \bar{c} \Rightarrow$ To prove symmetric

$$c = \begin{bmatrix} c_{xx_1} & c_{xx_2} & \dots & c_{xp} \\ c_{xx_2} & c_{xx_2} & \dots & c_{xp} \\ \vdots & \vdots & \ddots & \vdots \\ c_{xp} & c_{xp} & \dots & c_{xp} \end{bmatrix}_{p \times p}$$

$$c^H = \begin{bmatrix} c_{xx_1} & c_{xp} & \dots & c_{xp} \\ c_{xx_2} & c_{xx_2} & \dots & c_{xp} \\ \vdots & \vdots & \ddots & \vdots \\ c_{xp} & c_{xp} & \dots & c_{xp} \end{bmatrix}_{p \times p}$$

$$(AB)^H = [(AB)^T]^* \\ = [B^T A^T]^* \\ = [B^H A^H]^*$$

$$c^H = E[(\underline{x} - \bar{x}_x)(\underline{x} - \bar{x}_x)^H]$$
$$= E[(\underline{x} - \bar{x}_x)^H (\underline{x} - \bar{x}_x)^H]$$
$$= \bar{c}$$



covariance matrix is always hermitian matrix

If mean is 0 Then

$$\rightarrow E[(\underline{x})(\underline{x})^H]$$



correlation / covariance / auto-covariance

$$\underline{\mu_x} = E[\underline{x}]$$

$$\sigma_x^2 = E[(\underline{x} - \underline{\mu_x})(\underline{x} - \underline{\mu_x})^H]$$

$$c_{x,y} = E[(\underline{x} - \underline{\mu_x})(\underline{y} - \underline{\mu_y})^H]$$

Here \underline{x} & \underline{y} are vectors

$$c_{xy} = E[(\underline{x} - \underline{\mu_x})(\underline{y} - \underline{\mu_y})^H]$$

Here x & y are scalar vector

$$= E[(\underline{x} - \underline{\mu_x})(\underline{y} - \underline{\mu_y})^*]$$

$$= E[xy^*] - E[\mu_x y^*] - E[x \mu_y^*] + E[\mu_x \mu_y^*]$$

$\therefore \mu_x = E[\underline{x}]$ expectation & expectation on

expectation is always same ($\mu_x = E[x]$)

$$= E[xy^*] - \mu_x E[y^*] - \mu_y^* E[x] + \mu_x \mu_y^*$$

$$\Rightarrow E[xy^*] - \mu_x \mu_y^* - \cancel{\mu_y^* \mu_x} + \mu_x \mu_y^*$$

$$= E[xy^*] - \mu_x \mu_y^*$$

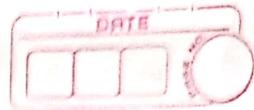
If $\mu_x, \mu_y > 0$ ie

$$E[xy^*] = \mu_x \mu_y^* \Rightarrow \text{uncorrelated}$$

[correlation between 2 random variables equals to product of means]



Random process \Rightarrow Random variable becomes
fn of time



Stationarity

- The data is now wst time

First order stationary

$$E[x(n)] = \bar{x}_n = \bar{x}$$

$x(n) \rightarrow$ Discrete time random

Avg Power = 1 \Rightarrow First order

process

Eg: So if there are 500 samples are taken & it results into mean \bar{x} ,
now again another 500 samples are taken & mean is taken again
which comes out to \bar{x} & so on

Then the series is stationary

Second order stationary / wide sense stationarity (WSS)

1. $E[x(n)] = \bar{x}$

2. $E[x(n) x^*(n+k)] = r(k)$ k: lag/gap

$$\begin{aligned} r^*(k) &= E[x(n+k) x^*(n)] & n+k = m \\ &= E[x(m) x^*(m-k)] & m = m-k \\ &= r(k) \end{aligned}$$

Hence, it is Hermitian/conjugate symmetry

If $k=0$,

$$r(0) = E[|x(n)|^2] \geq 0 \Rightarrow \text{Average power of signal}$$

Auto correlation:

- Correlation with same random variable

- $\bar{x} = 0$

$$R_{pp} = E\left[\frac{x}{-} \frac{x^*}{-}\right]$$

$$= \left[E[|x(n)|^2] - E[x(n)x^*(n-1)] - \dots - E[x(n)x^*(n-p)] \right]$$
$$= E[|x(n-p)|^2]$$

To get this value $\Rightarrow n - cn - k$



$$= \begin{bmatrix} r(0) & r(1) & \cdots & r(p) \\ r(p) & r(1-p) & \cdots & r(0) \end{bmatrix}$$

$$\therefore r^*(-k) = r(k)$$

$$= \begin{bmatrix} r(0) & r(1) & \cdots & r(p) \\ r^*(p) & r^*(p-1) & \cdots & r(0) \end{bmatrix}$$

Toeplitz Matrix \Rightarrow If above diagonal values are known then we can compute lower diagonal values.

$$A = \begin{bmatrix} r(0) & r(1) & \cdots & r(p) \\ r(p) & r(1-p) & \cdots & r(0) \end{bmatrix}$$

$$A^{-1} = \begin{bmatrix} (r(0))^* & (r(1))^* & \cdots & (r(p))^* \end{bmatrix}$$

$$A^{-1} = \begin{bmatrix} (r(0))^* & (r(1))^* & \cdots & (r(p))^* \end{bmatrix}$$

$$A^{-1} = \begin{bmatrix} (r(0))^* & (r(1))^* & \cdots & (r(p))^* \end{bmatrix}$$

Werner Föhrer

→ If $T_{p \times p}$ is a Hermitian matrix or a Toeplitz matrix $\left\{ \begin{array}{l} T_{p \times p} = T_{p \times p}^H \\ \underline{x} \neq 0 \end{array} \right.$

$$\xrightarrow{\text{Matrix}} \xrightarrow{\text{Eigen vector}} T\underline{x} = \lambda \underline{x} \rightarrow \text{Eigen value decomposition}$$

$$(T\underline{x})^H = \lambda^* \underline{x}^H \quad \{ \text{order doesn't matter in } \lambda \underline{x} \text{ since } \lambda \}$$

$$\underline{x}^H T^H = \lambda^* \underline{x}^H \quad \{ \text{is a scalar} \}$$

$$\underline{x}^H T = \lambda^* \underline{x}^H$$

$$\underline{x}^H T \underline{x} = \lambda^* \underline{x}^H \underline{x} \quad \{ \text{(Multiplying } \underline{x} \text{ on both the sides)} \}$$

$$\underline{x}^H T \underline{x} = \lambda^* \underline{x}^H \underline{x} \quad \{ \text{we multiply after } \underline{x}^H T \text{ since otherwise multiplication might get merged up) } \}$$

$$\lambda \underline{x}^H \underline{x} = \lambda^* \underline{x}^H \underline{x} \quad \{ \text{multiplication might get merged up) } \}$$

$$(\lambda - \lambda^*) \underline{x}^H \underline{x} = 0$$

By property,

$$\underline{x}^H \underline{x} = \sum_{i=1}^p |x_i|^2$$

$$\therefore (\lambda - \lambda^*) = 0$$

$$\therefore \lambda = \lambda^* \quad \{ \text{eigen values are always real} \}$$

Hence, eigen values are always real $\therefore \underline{x}^H \underline{x} \neq 0$

→ Prove $\underline{x}_1^H \underline{x}_2$ are orthogonal if $\lambda_1 \neq \lambda_2$

$$T \underline{x}_1 = \lambda_1 \underline{x}_1 \quad \{ \text{---(1)} \}$$

$$T \underline{x}_2 = \lambda_2 \underline{x}_2 \quad \{ \text{---(2)} \}$$

Apply Hermitian on eqn (1)

$$(T \underline{x}_1)^H = (\lambda_1 \underline{x}_1)^H$$

$$\underline{x}_1^H T^H = \lambda_1^* \underline{x}_1^H \quad \{ \text{---(3)} \}$$

Multiply \underline{x}_2 on eqn (3)

$$\underline{x}_1^H T^H \underline{x}_2 = \lambda_1^* \underline{x}_1^H \underline{x}_2$$

$$\underline{x}_1^H T \underline{x}_2 = \lambda_1^* \underline{x}_1^H \underline{x}_2$$

$$\underline{x}_1^H \lambda_2 \underline{x}_2 = \lambda_1^* \underline{x}_1^H \underline{x}_2$$

$$\lambda_2 \underline{x}_1^H \underline{x}_2 = \lambda_1^* \underline{x}_1^H \underline{x}_2$$

$$(\lambda_2 - \lambda_1^*) \underline{x}_1^H \underline{x}_2 = 0 \quad \{ \because \lambda_1 = \lambda_1^* \}$$

since $\lambda_1 \neq \lambda_2 \therefore \underline{x}_1^H \underline{x}_2 = 0 \Rightarrow \text{orthogonal}$

Example:

$$A = \begin{bmatrix} 1+2j & 2 \\ 3 & 4+2j \end{bmatrix} \quad A^H = \begin{bmatrix} (1+2j)^* & 3^* \\ 2^* & (4+2j)^* \end{bmatrix}$$

$$A^H A = \begin{bmatrix} (1+2j)^* & 3 \\ 2 & (4+2j)^* \end{bmatrix} \begin{bmatrix} 1+2j & 2 \\ 3 & 4+2j \end{bmatrix}$$
$$= [(1+2j)^*(1+2j) + 9 + 4 + (4+2j)^*(4+2j)]$$

$$= (1-4j^2)(1+2j) + 13 + (4-2j)(4+2j)$$
$$= (1-(4j^2)) + 13 + 16 - 4j^2$$
$$= 30 - 8j^2$$
$$= 38$$

$$= \begin{bmatrix} 1+4+9 & 2-2j+12+6j \\ 2+4j+12-6j & 4+16+4 \end{bmatrix}$$
$$= \begin{bmatrix} 14 & 14+4j \\ 14-2j & 24 \end{bmatrix}$$

\Rightarrow If all eigen values & vectors are combined into matrix

$$\underline{E} = [e_1 \ e_2 \ e_3 \cdots e_p] = [\lambda_1 e_1 \ \lambda_2 e_2 \ \cdots \ \lambda_p e_p]$$

lets call as E

$$\underline{E} \underline{E} = [e_1 \ e_2 \ \cdots \ e_p] \begin{bmatrix} \lambda_1 & 0 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & & & & \lambda_p \end{bmatrix}$$

$$\underline{E} \underline{E} = \underline{E} \underline{D}$$

$$\underline{E}^H \underline{E} = \begin{bmatrix} e_1^H \\ e_2^H \\ \vdots \\ e_p^H \end{bmatrix} [e_1 \ e_2 \ \cdots \ e_p] \quad \text{--- (4)}$$



Apply E^H both the sides on eqn (1)

$$I = E^H E \Rightarrow I = D E^H$$

As we proved in previous theorem, $\underline{x}_1^H \underline{x}_1$ are orthogonal

$$\begin{array}{l} I = E D E^H \\ \boxed{I = E D^H E} \end{array}$$

decomposition of Toeplitz matrix

$\left. \begin{array}{l} S \\ EE^H = I \\ E^H = E^H \end{array} \right\}$ We use co-hermitian since it avoids inverse matrix computation

- * computation of inverse matrix is difficult for non-square matrix
- * If $|A| = 0$ then it might be problematic

Weiner Filters

- Important for optimization in image processing
- Uses supervised learning type
- Also called Minimum Mean Squares or least square error filter
- Minimizing error

$$e(n) = y(n) - d(n)$$

↓ ↓ ↓

error o/p actual value

$$E^2 = E[e^2(n)] \Rightarrow \text{We apply statistical operation since there will be randomness in signals}$$

Mean

squared error

$$\begin{aligned} &= E[e(n)e^H(n)] \\ &= E[e(n)e^T(n)] \quad \left\{ \because e^H(n) is e^T(n) as e(n) has real values \right\} \\ &= E[(d(n) - y(n))e^T(n)] \\ &= E[(d(n) - \underline{w}^T \underline{x}(n))e^T(n)] \quad \left\{ \because y(n) = \underline{w}^T \underline{x}(n) \right\} \end{aligned}$$

$$= E[(d(n) - \underline{w}^T \underline{x}(n))(d(n) - \underline{w}^T \underline{x}(n))^T]$$

↓
Weiner filter is linear
estimation of original image



$$\begin{aligned}
 & R_{ww}(n) \begin{pmatrix} x(n) \\ x(n-1) \end{pmatrix} = w^T x(n) + w_0 \\
 & x^T(n) = [x(n) \ x(n-1)] \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \\
 & \Rightarrow E[d^2(n) - w^T x(n)d(n) - w^T x(n)x(n)^T + w^T x(n)(w^T x(n))^T] \\
 & = E[d^2(n) - w^T x(n)d(n) - d(n)(w^T x(n))^T + w^T x(n)(w^T x(n))^T] \\
 & = E[d^2(n) - w^T x(n)d(n) - d(n)x^T(n)w + \\
 & \quad w^T x(n)(x^T(n)w)] \\
 & = E[d^2(n) - w^T x(n)d(n) - d(n)w^T x(n) + \\
 & \quad w^T x(n)x^T(n)w] \\
 & = E[d^2(n) - 2w^T x(n)d(n) + w^T x(n)x^T(n)w] \\
 & = E[d^2(n)] - 2E[w^T x(n)d(n)] + E[w^T x(n)x^T(n)w] \\
 & = \sigma^2 - 2w^T E[x(n)d(n)] + w^T E[x(n)x^T(n)]w \\
 & E^2 = \frac{\sigma^2}{d} - \frac{2w^T p}{d} + \frac{w^T R w}{d} \quad (1) \quad \left. \begin{array}{l} \text{we don't apply expectation of } w \\ w \text{ since its static in nature} \end{array} \right\} \text{cross correlation b/w } x(n) \text{ and } d(n) \quad \left. \begin{array}{l} \text{doesn't change with time} \end{array} \right\}
 \end{aligned}$$

Gradient descent optimization

objective: w should be optimized ie coefficients should be optimized

Take first order derivative of w for eqn (1)

$$\nabla_w E^2 = 0 \quad (\because \text{we want minima})$$

Since E^2 is a vector so we need to equate with 0 vector

$$\nabla_w E^2 = \frac{\partial}{\partial w_k} \left(\frac{\sigma^2}{d} - \frac{2w^T p}{d} + \frac{w^T R w}{d} \right)$$

$$\text{let } A = w^T p$$

second order derivative \rightarrow +ve \rightarrow Minima
 second order derivative \rightarrow -ve \rightarrow Maxima



$$\frac{\partial A}{\partial w_k} =$$

$$\frac{\partial A}{\partial w_k}$$

$$A = W_0 P(0) + W_1 P(1) + \dots + W_k P(k) + \dots + W_N P(N)$$

$$\frac{\partial A}{\partial w_k} = P(k)$$

$$\frac{\partial A}{\partial w_k}$$

$$\nabla_w A = \begin{bmatrix} \frac{\partial A}{\partial w_0} \\ \frac{\partial A}{\partial w_1} \\ \vdots \\ \frac{\partial A}{\partial w_N} \end{bmatrix} = \begin{bmatrix} P(0) \\ P(1) \\ \vdots \\ P(N) \end{bmatrix} \Rightarrow P$$

$$\text{let } B = \underbrace{W^T R w}_{\substack{\rightarrow (N+1) \times 1 \\ \rightarrow (N+1) \times (N+1)}} \quad (N+1) \times 1$$

$$B = \sum_{i=0}^N w_i \sum_{j=0}^N R_{ij} w_j$$

$$\nabla_w B = \begin{bmatrix} \frac{\partial B}{\partial w_0} \\ \frac{\partial B}{\partial w_1} \\ \vdots \\ \frac{\partial B}{\partial w_N} \end{bmatrix} = \begin{bmatrix} w_0 + R_{00}w_0 + R_{01}w_1 + \dots + R_{0N}w_N \\ w_1 + R_{00}w_0 + R_{11}w_1 + \dots + R_{1N}w_N \\ \vdots \\ w_N + R_{00}w_0 + R_{11}w_1 + \dots + R_{NN}w_N \end{bmatrix}$$

$$\frac{\partial B}{\partial w_k} B = \sum_{i=0}^N w_i \sum_{j=0}^N R_{ij} w_j + w_k \sum_{j=0}^N R_{kj} w_j$$

if k

\downarrow
(j does not contain k^{th} term)

$$\frac{\partial B}{\partial w_k} = \sum_{i=0}^N w_i R_{ik} + \frac{\partial}{\partial w_k} \left[\sum_{j=0}^N R_{kj} w_j + w_k R_{kk} w_k \right]$$

$$= \sum_{i=0}^N w_i R_{ik} + \sum_{j=0, j \neq k}^N R_{kj} w_j + 2w_k R_{kk}$$

$$= \sum_{i=0, i \neq k}^N w_i R_{ik} + \sum_{i=0, i \neq k}^N R_{kj} w_j + w_k R_{kk} + w_k R_{kk}$$

$$= \sum_{i=0}^N R_{ik} w_i + \sum_{j=0}^N R_{kj} w_j$$

$\because R$ is hermitian, $R^T = R \therefore R_{ji} = R_{ij}$

$$= \sum_{i=0}^N w_i R_{ki} + \sum_{j=0}^N R_{kj} w_j$$

$$= \left[w_0 R_{k0} + w_1 R_{k1} + \dots + w_N R_{kN} \right] + \left[R_{k0} w_0 + R_{k1} w_1 + \dots + R_{kN} w_N \right]$$

$$= 2 \sum_{i=0}^N w_i R_{ki}$$

$$\frac{\partial B}{\partial w_0} = 2 \sum_{i=0}^N w_i R_{0i} = 2 \left[w_0 R_{00} + w_1 R_{01} + \dots + w_N R_{0N} \right]$$

$$\frac{\partial B}{\partial w_1} = 2 \sum_{i=0}^N w_i R_{1i} = 2 \left[w_0 R_{10} + w_1 R_{11} + \dots + w_N R_{1N} \right]$$

$$\therefore \nabla_{wB} = \begin{bmatrix} \frac{\partial B}{\partial w_0} \\ \frac{\partial B}{\partial w_1} \\ \vdots \\ \frac{\partial B}{\partial w_N} \end{bmatrix} = 2 R w$$

(R is matrix
 w is vector)

$$\nabla_w J E^2 = 0 \cancel{+ 2P} + 2Rw = 0$$

$$= -2P + 2Rw = 0$$

$$\therefore P = R w$$

w is optimal signal
of final guess

$$\boxed{w^* = R^{-1} P}$$



(When R is orthogonal then we can write $R^{-1} = R^H$)

$$\text{Ex: } R = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix} \quad P = \begin{bmatrix} 0.2 \\ 0.1 \end{bmatrix}$$

$w = ?$ using Gradient descent opti

$$|R| = ad - bc$$

$$= 1 - 0.81$$

$$= 0.19$$

$$R^H = \frac{1}{0.19} \begin{bmatrix} 1 & -0.9 \\ -0.9 & 1 \end{bmatrix} = \begin{bmatrix} 5.26 & -4.73 \\ -4.73 & 5.26 \end{bmatrix}$$

$$w = R^H P$$

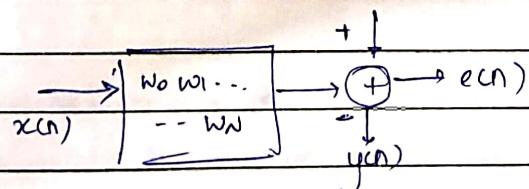
$$= \frac{1}{0.19} \begin{bmatrix} 1 & -0.9 \\ -0.9 & 1 \end{bmatrix} \begin{bmatrix} 0.2 \\ 0.1 \end{bmatrix}$$

$$= \begin{bmatrix} 5.26 & -4.73 \\ -4.73 & 5.26 \end{bmatrix} \begin{bmatrix} 0.2 \\ 0.1 \end{bmatrix}$$

$$= \begin{bmatrix} 0.579 \\ -0.42 \end{bmatrix}$$

$\underline{\underline{}}$

dn)



$$E^2 = E[e^2(n)]$$

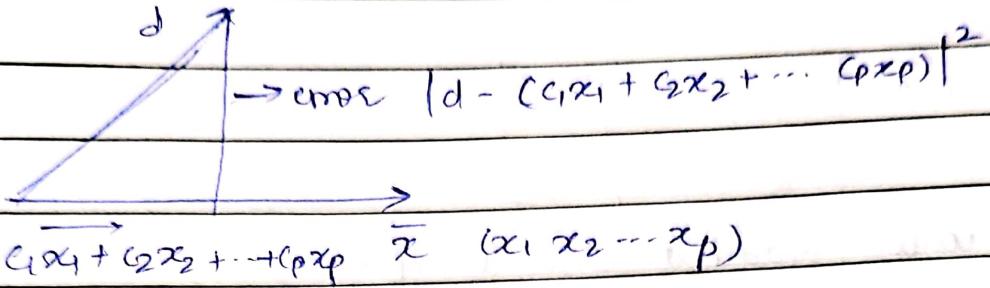
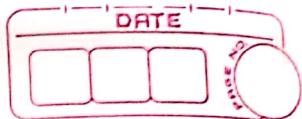
Objective: Minimize mean square error.

Let us consider a random variable x whose values are x_1, x_2, \dots, x_p .

Suppose we want to estimate another random variable who is linear combination of x & it is called linear estimation.

$$d^D = c_1 x_1 + c_2 x_2 + \dots + c_p x_p$$





if weights are optimal then error is minimum or is orthogonal
to \bar{x}

$$\begin{bmatrix} E_1 \\ E_2 \\ E_3 \\ E_4 \\ E_5 \end{bmatrix} = \begin{bmatrix} P_1 & Q_1 \\ P_2 & Q_2 \\ P_3 & Q_3 \\ P_4 & Q_4 \\ P_5 & Q_5 \end{bmatrix} \begin{bmatrix} S_1 \\ S_2 \end{bmatrix}$$

$$\begin{bmatrix} S_1 \\ S_2 \end{bmatrix} = \begin{bmatrix} P_1 & Q_1 \\ P_2 & Q_2 \\ P_3 & Q_3 \\ P_4 & Q_4 \\ P_5 & Q_5 \end{bmatrix}^{-1} \begin{bmatrix} E_1 \\ E_2 \\ E_3 \\ E_4 \\ E_5 \end{bmatrix}$$