

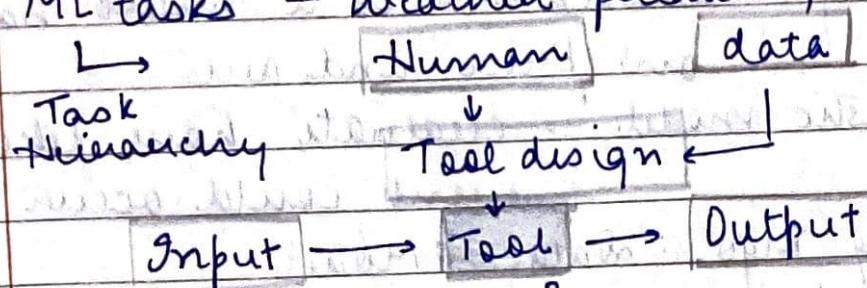
MLF

CLASSTIME	Page No.
Date	/ /

Week-1

L1.1 What is ML?

- ML is the study of comp. algo that improves automatically through experience and by-the use of data.
- ML tasks - Weather predict^n, Face detect^n



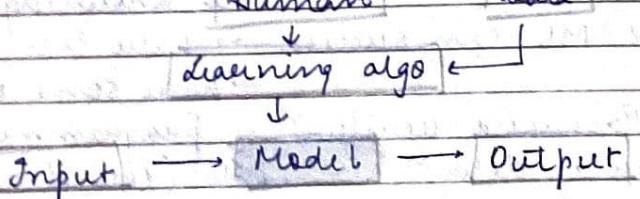
- Why & when ML?

1. Human fails → scale / speed / cost
→ inability to express
→ no knowledge of rules
 2. ML can succeed → have lots of common data
→ have some struc. idea
- Password verifica^n - programming
 - Face detect^n - ML (lots of image available)
 - Weather predict^n - ML
 - ML in inbox - spam classifier
 - ML in shopping cart - frequently recommended system - Netflix, YT
 - Smart Assistant - Alexa, Siri, etc.
 - ML in Robot AIs, Games (Chess, Checkers, Go), Marketing (budgeting)

L1.2 Data, Models and ML Tasks

- data - it is a collect^n of vectors
- metadata - it is info. on the data
- model - it is a mathematical simplification of reality
- In ML - Predictive Model → Regression
Probabilistic → Classification

- Regression model → model the price of a house based on meters dist. & area
output price = $0.5 \times \text{area} - \text{distance}$
- Classification → model whether a house is closer (discrete values) than 2 kms to a metro based on price and area
- probabilistic model → evaluate how likely an event could occur
- Learning algo : Data → Model
↳ choose from a collection of models, with same structure but diff. parameters
→ use data to get the 'best' parameters
- ML Human Data



L1.2 Supervised Learning: Regression

- supervised learning is curve-fitting
- Given $\{(x^1, y^1), (x^2, y^2), \dots, (x^n, y^n)\}$
↳ find a model f such that $f(x^i)$ is n to y^i
- Regression - e.g. predict house price from room, area, distance

$$f: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\text{Loss}(f) = \frac{1}{n} \sum_{i=1}^n (f(x^i) - y^i)^2$$

$$f(x) = w^T x + b = \sum_{j=1}^d w_j x_j + b \\ (w_1 x_1 + w_2 x_2 + \dots)$$

e.g. Room A dist. price

$$\begin{aligned} f &= 2 \times \text{rooms} \\ &\quad - 0.5 \times d \\ g &= \text{rooms} + \dots \end{aligned}$$

→ calculate f & g , & check for the lower loss

L1.4 SL : Classification

e.g. Predict if rooms > 3 from area & price
 $x^i \in \mathbb{R}^d, y^i \in \{+1, -1\}$
 $f: \mathbb{R}^d \rightarrow \{+1, -1\}$ Loss = $\frac{1}{n} \sum_{i=1}^n \text{I}\{f(x^i) \neq y^i\}$

$$f(x) = \text{sign}(w^T x + b)$$

↳ linear separator

- learning algo uses training data to get model f . But validation should n't be on the training data. Use test data which is not in the training data.
- Learning algo finds the best model in collection of models. This is c/d model selection and it is done by using another dataset c/d validation data that is distinct from train & test data

L1.5 Unsupervised Learning: Dimensionality Reduction

- UL - understanding data. Data : $\{x^1, x^2, \dots, x^n\}$
 $x^i \in \mathbb{R}^d$

- it build models that compress, explain and group data
- e.g. Coke tweet (e.g. → groups of 1D)
- Dimensionality reduction - compression and simplification
e.g. represent 1M gene expression of a million people

$$\text{Data} : \{x^1, x^2, \dots, x^n\} \quad x^i \in \mathbb{R}^d$$

Encoder: $R^d \rightarrow R^{d'}$ Decoding: $R^{d'} \rightarrow R^d$
 $d' < d$

Goal: $g(f(x_i)) \approx x_i$

$$\text{Loss} = \frac{1}{n} \sum_{i=1}^n \|g(f(x_i)) - x_i\|^2$$

L1-6 VL: Density Estimation

- a probabilistic model
- create a robot accnt. which generates more such tweets
- Eg. wisdomofcrowds.com
- DE model takes in several samples from a random source, and outputs a model that assigns a probability score to every possible sentence.
- Data: $\{x^1, x^2 \dots x^n\}$ - $x^i \in R^d$
- Probability mapping $P: R^d \rightarrow R$, that sum to 1.
- Goal: $P(x)$ is large if $x \in \text{Data}$, and low otherwise
- Loss: $\frac{1}{n} \sum_{i=1}^n -\log(P(x^i))$
use log likelihood

T1.1 To ML or not to ML

- data/samples
- do we need to have a ML approach to solve a given ques? \rightarrow If ML - find the metadata about the data.
- a, b, c are the parameters or wts. of the model. The best values for the parameters will

be learning from data

T1.2

Illustrate with a real-world dataset

- breast cancer classifier
- x_j $1 \leq j \leq 30 \Rightarrow$ columns (features)
- x^i $1 \leq i \leq 569 \Rightarrow$ rows (cases)
- $y^i = f(x^i) - b$
- $\text{Loss} = \frac{1}{n} \sum (y^i \neq y^i)$
- which func is suitable?
- Train Set, Validation Set, Test Set
80% of Total 20% of Train 20% of Total

T1.3

VL with applications

- how can we compress 30 features to 3 features w/o losing much information
↳ dimensionality reduction
- (only data with no labels)

$$u = f(x) = Wx \Rightarrow \text{encoder}$$

$$x = g(u) = W^T u \Rightarrow \text{decoder}$$

T1.4

- Density estimation \rightarrow all samples are assumed to be from some probability distribution. 569 samples are insufficient

Generate new samples by estimating the PDF for each features.

- new face generated by deep learning modi
- StyleGAN (Generative Adversarial networks)
↳ density estimation

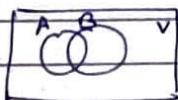
Week 2

CLASSTIME	Page No.
Date	/ /

CLASSTIME	Page No.
Date	/ /

1.2.1 Sets & Functions

- \mathbb{R} (set of real nos.)
- \mathbb{R}_+ (set of +ve real nos including 0)
- \mathbb{Z} (set of integers)
- \mathbb{Z}_+ (set of +ve integers with 0)
- $[a, b] = \{x \in \mathbb{R} : a \leq x \leq b\}$
- $(a, b) = \{x \in \mathbb{R} : a < x < b\}$
- \mathbb{R}^d set of d-dimensional vectors
- $\mathbb{R}^d : D(x, y) = \|x - y\| = \sqrt{(x_d - y_d)^2}$
- (open) $B(x, \epsilon) = \{y \in \mathbb{R}^d : D(x, y) < \epsilon\}$
- (closed) $\bar{B}(x, \epsilon) = \{y \in \mathbb{R}^d : D(x, y) \leq \epsilon\}$
- V is the universe



$$A \cup B, A \cap B, A^c, B^c$$

$$\begin{aligned} (A \cup B)^c &= A^c \cap B^c \\ (A \cap B)^c &= A^c \cup B^c \end{aligned} \quad \left. \begin{array}{l} \text{De Morgan's laws} \\ \text{laws} \end{array} \right\}$$

- \forall for all \exists there exist \Rightarrow implies
 \Leftrightarrow equivalent to

$$\begin{aligned} &x_1, x_2, \dots \text{ where } x_i \in \mathbb{R}^d \\ &\lim_{i \rightarrow \infty} x_i = x^* \\ &\forall \epsilon > 0, \exists N \end{aligned} \quad \left. \begin{array}{l} \text{convergent sequence} \\ \text{if } x_n \in B(x^*, \epsilon) \forall n \geq N \end{array} \right\}$$

- (i) $x_i \in \mathbb{R}$ $x_n = 1 + n$
- (ii) $x_i \in \mathbb{R}^2$ $x_n = \left(\frac{1}{2^n} \cos\left(\frac{\pi n}{2}\right), \frac{1}{2^n} \sin\left(\frac{\pi n}{2}\right) \right)$

- If V is a vector space, $v \in V, u \in V$
- $\alpha v + \beta w \in V$ (most crucial prop.)
- $\Rightarrow \mathbb{R}^d$ is a vector space
- $\rightarrow x \cdot y = x^T y = \sum_{i=1}^d x_i y_i$ (dot prod.)

$\rightarrow \|x\|^2 = x \cdot x = x^T x = \sum_{i=1}^d x_i^2$ (norm)
 $\rightarrow x \& y$ are \perp to each other if their dot product is 0.

- $f : A \rightarrow B$ $f : \mathbb{R} \rightarrow \mathbb{R}$
 A (domain) B (co-domain)

$$G_f \subseteq \mathbb{R}^{d+1}$$

$$G_f = \{ (x, f(x)) : x \in \mathbb{R}^d \}$$

- contour plots of 2 dimensional functions
 \hookrightarrow heat map (when fully coloured)
 \hookrightarrow it has ∞ contours. so more information

1.2.2 Univariate Calc: Continuity & Differentiability

$$f : \mathbb{R} \rightarrow \mathbb{R}$$

is continuous at $x^* \in \mathbb{R}$ if for all sequences $x = x_1, x_2, \dots$ converges to x^* . we have that $f(x_i)$ converges to $f(x^*)$

$$\lim_{i \rightarrow \infty} x_i = x^* \Rightarrow \lim_{i \rightarrow \infty} f(x_i) = f(x^*)$$

$$\lim_{x \rightarrow x^*} f(x) = f(x^*)$$

g.

$$f(x) = \operatorname{sign}(x) = \begin{cases} -1 & \text{if } x < 0 \\ 0 & \text{if } x = 0 \\ +1 & \text{if } x > 0 \end{cases}$$

- a funcⁿ is said to be continuous if it is continuous at all points in its domain

- a funcⁿ $f : \mathbb{R} \rightarrow \mathbb{R}$ is differentiable at $x^* \in \mathbb{R}$ if $\lim_{x \rightarrow x^*} \frac{f(x) - f(x^*)}{x - x^*}$ exists.

$$f'(x^*)$$

- if f is not continuous at x^* $\Rightarrow f$ is not differentiable at x^*

$$f'(x^*) = \lim_{x \rightarrow x^*} \frac{f(x) - f(x^*)}{x - x^*}$$

$$= \lim_{h \rightarrow 0} \frac{f(x^* + h) - f(x^*)}{h}$$

L2.3 Derivatives and linear approx.

- let $f: \mathbb{R} \rightarrow \mathbb{R}$ be a diff. func.

$$f'(x^*) = \lim_{x \rightarrow x^*} \frac{f(x) - f(x^*)}{x - x^*}$$

$$f'(x^*) = \frac{f(x) - f(x^*)}{x - x^*} \quad (\text{around } x = x^*)$$

$$f(x) = f(x^*) + f'(x^*) \cdot (x - x^*)$$

$$\boxed{f(x) \approx L_{x^*}[f](x)}$$

$$\begin{aligned} \text{Q. } f(x) &= x^2 \quad L_{x^*}[f] = f(x^*) + f'(x^*)(x - x^*) \\ &\quad x=1 \\ &= 1^2 + 2 \cdot (x-1) \\ &= 2x-1 \end{aligned}$$

- $L_{x^*}[f]$ is a tangent to the graph of f at the point $(x^*, f(x^*))$.

$$\begin{aligned} \text{Eq. } f(x) &= \sin(x) \quad x^* = 0 \\ f'(x) &= \cos(x) \quad f(x) = 0 + 1(x-0) \\ f'(x^*) &= 1 \quad = x \\ f'(x^*) &= 0 \quad \sin x \approx x \quad \text{if } x=0 \end{aligned}$$

Eq.

$$\begin{aligned} f(x) &= e^x \quad x^* = 0 \\ e^x &= e^0 + (x-0) \cdot 1 \\ &\approx 1 + x \quad \text{around } x=0 \end{aligned}$$

Eq.

$$\ln(1+x) \quad \text{around } x^* = 0$$

Eq.

$$f(x) = (1+x)^n \quad \text{around } x^* = 0$$

L2.4

Multivariable Calculus: app & adv rules

Linear App. $\rightarrow f(x) \approx L_{x^*}[f](x)$

Quadratic App. $\rightarrow f(x) \approx L_{x^*}[f](x) + \frac{1}{2} f''(x^*)$

(higher order apprx. are better than linear apprx.)

$$\begin{aligned} f(x) &= x^2 \\ x^2 &= (x^*)^2 + 2x^*(x-x^*) + \frac{1}{2} \cdot 2 \cdot (x-x^*)^2 \end{aligned}$$

$$= x^2 \quad (\text{only for quadratic eqn})$$

Eq.

higher order approx. of e^x around $x^* = 0$

$$x \approx 1 + x + \frac{x^2}{2}$$

-

$$f(x) = g(x) \cdot h(x) \quad \text{prod. rule}$$

$$f'(0) = g'(0)(h(0)) + h'(0)g(0)$$

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

$$f'(0) = g'(h(0))h'(0)$$

Chain rule

Eq. $\frac{e^{3x}}{\sqrt{1+x}} \text{ around } x=0 \Rightarrow (x)$

$$\approx (1+3x) \left(1 - \frac{x}{2} \right)$$

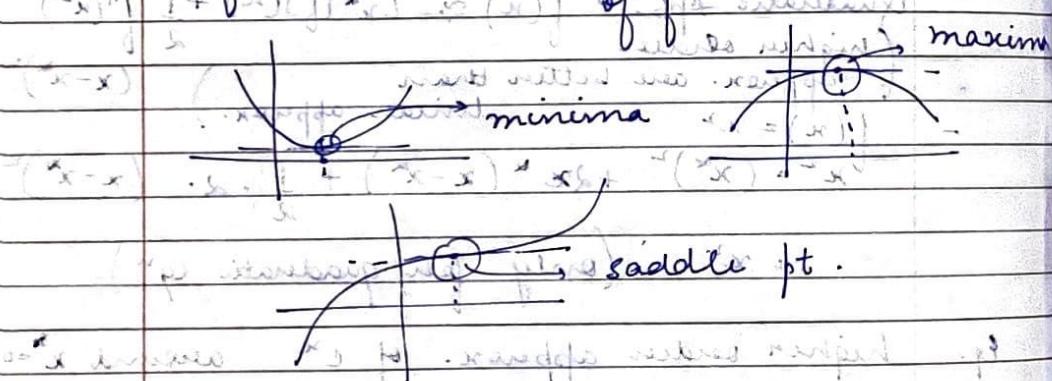
$$= 1 + 5x$$

$$= e^{3x} \text{ around } 2 \quad (x+1) \approx (x)$$

Eq. $e^{\sqrt{1+x}} \text{ around } x=1$

$$\approx e^{\sqrt{2}} + \frac{e^{\sqrt{2}}}{2\sqrt{2}} (x-1)$$

- $L(x) = f(x) + f'(x)(x-x_0)$
 $f'(x_0) = 0 \Leftrightarrow x_0 \text{ is a critical point}$



L2.5 Multi. Cal. - Lines & planes in higher dim.

(i) A line in \mathbb{R}^d : $x = u + \lambda v$ plane

(ii) A line thru the pt. $u \in \mathbb{R}^d$ along the

(iii) vector $w \in \mathbb{R}^d$: $x = u + \alpha w$

$$= \{x \in \mathbb{R}^d : x = u + \alpha w \text{ for } \alpha \in \mathbb{R}\}$$

similarly line thru (u, u') $\in \mathbb{R}^d$
 u along $u' - u$
 u along $u - u'$

- A $(d-1)$ dimensional hyperplane \mathbb{C}^{d-1}
- A hyperplane normal to the vector $w \in \mathbb{R}^d$ with value $b \in \mathbb{R}$

$$= \{x \in \mathbb{R}^d : w^T x = b\}$$

Eq. Line thru (1) along (2)

- Tuples vs Points vs Vectors
 (w, x) : (location) + (w direction)

$$f: \mathbb{R}^2 \rightarrow \mathbb{R} \quad f(x, x_2) = x_1^2 + x_2^2$$

$$\frac{df}{dx_1}(u) = \lim_{\alpha \rightarrow 0} \frac{f(u + [\alpha]) - f(u)}{\alpha}$$

$$= f(u_1 + \alpha, u_2) - f(u_1, u_2)$$

$$(u_1, u_2) \cdot (u) \approx f(u_1, u_2) + \frac{df}{dx_1}(u) \alpha (u_1, u_2)$$

$$\frac{df}{dx_2}(u) = \lim_{\alpha \rightarrow 0} \frac{f(u + \alpha, u_2) - f(u)}{\alpha}$$

- Gradients

$$f: \mathbb{R}^d \rightarrow \mathbb{R} \quad \nabla f(u) = (f'_1(u), f'_2(u), \dots, f'_d(u))$$

$$\nabla f(u) = \left[\frac{df}{dx_i}(u) \right]_{i=1}^d$$

$$df = 2x_1 dx_1 + 2x_2 dx_2 \quad \frac{df}{dx_1} = 2x_1, \quad \frac{df}{dx_2} = 2x_2$$

$$\frac{df}{dx_1}(u) = 2u_1 \quad \frac{df}{dx_2}(u) = 2u_2$$

$$f(x) = x_1 + 2x_2 + x_3 \quad \Rightarrow \text{linear func} \quad \nabla f(u) = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}$$

! subsequent gradients - previous

$$\|u\|_1 = |u_1| + |u_2| + \dots + |u_d| \geq \|u\|_2 \geq \|u\|_0 = |u_1| + |u_2| + \dots + |u_d|$$

L2.1 MC: Linear approximation & app:

$$f(x) \approx f(x^*) + f'(x^*)(x - x^*)$$

$$(x) = f(x^*) + f'(x^*)(x - x^*)$$

$$\text{Let } g: \mathbb{R}^d \rightarrow \mathbb{R} \text{ with } x \in \mathbb{R}^d \text{ which is valid}$$

$$f(x) \approx f(u) + (\nabla f(u))^T (x - u)$$

$$(x) = f(u) + L_u[f](x)$$

- $f: \mathbb{R}^2 \rightarrow \mathbb{R}$

$$f(y_1, u_2) \approx f(u_1, u_2) + \frac{\partial f}{\partial x_1}(u) \cdot (y_1 - u_1)$$

$$\text{Ex. } f(x_1, x_2) = x_1^2 + x_2^2 \quad \nabla f(x) = \begin{bmatrix} 2x_1 \\ 2x_2 \end{bmatrix}$$

(i) App. of around $(6, 2)$

$$f(u) = 40, \quad \nabla f(u) = \begin{bmatrix} 12 \\ 4 \end{bmatrix}$$

$$f(x) = 40 + [12, 4] \begin{bmatrix} x_1 - 6 \\ x_2 - 2 \end{bmatrix}$$

$$= 12x_1 + 4x_2 - 40$$

- The graph of $L_u[f]$ is tangent to the graph of f at the pt. $(u, f(u))$

- Directional derivatives

$$D_u[f](u) = \lim_{\alpha \rightarrow 0} \frac{f(u + \alpha u) - f(u)}{\alpha}$$

$$= \nabla f(u)^T u$$

- Cauchy-Schwarz Inequality.

$$- \|a\| \cdot \|b\| \leq a^T b \leq \|a\| \cdot \|b\|$$

$$\downarrow \quad \downarrow$$

$$a = \alpha b \quad a = \alpha b$$

$$\alpha \leq 0 \quad \alpha \geq 0$$

- dirⁿ of Steepest ascent

$$\text{maximising } D_u[f](u)$$

$$= \nabla f(u)^T u$$

- descent dirⁿ

$f: \mathbb{R}^d \rightarrow \mathbb{R}$
what are valid dirⁿ, $u \in$

$$u: D_u[f](u) \leq 0$$

$$\text{descent dir}^n: \{u \in \mathbb{R}^d : \nabla f(u)^T u \leq 0\}$$

- If $f(x)$ is minimised at u $\Rightarrow \nabla f(u) = 0$

$$\{u : \nabla f(u) = 0\} \rightarrow \text{critical point}$$

Week - 3

1.3.1 4 Fundamental Subspaces

 - Column space $C(A)$

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 1 & 3 \\ 3 & 1 & 4 \\ 4 & 1 & 5 \end{bmatrix} \quad C(A) = \text{span}(u_1, \dots, u_n)$$

 solving $Ax = b$:
 $b \in C(A)$ if and only if

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 1 & 3 \\ 3 & 1 & 4 \\ 4 & 1 & 5 \end{bmatrix} \quad 4 \times 3 \text{ with 2 unknowns}$$

col 3 = col 1 + col 2

$$C(A) = \text{span}(\text{col 1, col 2})$$

\Rightarrow 2D subspace of R^4

 Null space $N(A)$

$$N(A) = \{x \mid Ax = 0\}$$

 A as before \Rightarrow LC of columns of A
 Should give '0' vector

$$\begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix} \in N(A)$$

 $N(A)$ is a line in R^3
 \rightarrow If A is invertible, $N(A)$ has '0' only
 $\& C(A)$ is the whole space

 $Ax = b$ has a unique soln, $x = A^{-1}b$

 Else, $N(A) = x_n \neq 0$
 then, more solns

 \rightarrow Gaussian elimination to get $N(A)$

$$\begin{bmatrix} 1 & 2 & 2 & 2 \\ 2 & 4 & 6 & 8 \\ 3 & 6 & 8 & 10 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 2 & 2 \\ 0 & 0 & 2 & 4 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad \begin{aligned} x_2 &= 1 \\ x_4 &= 0 \\ x_3 &= 0 \\ x_1 &= -2 \end{aligned}$$

$$\begin{bmatrix} -2 \\ 1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 2 \\ 0 \\ -2 \\ 1 \end{bmatrix} \in N(A)$$

Rank = no. of pivot col = 2

Nullity = no. of free var. = 2

 Rank = $\dim(C(A))$, Nullity = $\dim(N(A))$
 A has n col. \Rightarrow Rank + nullity = n

 - Row space: colⁿ space of $A^T \Rightarrow$ span of rows of A

 col. rank = $\dim(C(A))$, $R(A) = C(A^T)$

 row. rank = $\dim(R(A))$

col. rank = row. rank

 - Left Null space $N(A^T)$

$$m \times n \Rightarrow \underbrace{\begin{bmatrix} y_1 & \dots & y_n \end{bmatrix}}_{N(A^T)} \begin{bmatrix} A^T \end{bmatrix} = \begin{bmatrix} 0 & \dots & 0 \end{bmatrix}$$

 - A is $m \times n$ matrix

$$\dim(C(A)) + \dim(N(A)) = \text{no. of col} = n$$

$$\dim(C(A^T)) + \dim(N(A^T)) = m = \text{no. of rows}$$

$$A \Rightarrow N(A^T) = \{(1, 1, -1)\}$$

$$\text{Eq. } A = \begin{bmatrix} 1 & 2 \\ 3 & 1 \end{bmatrix} \quad m=n=2$$

$$C(A) = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$$

rank = 1

$$N(A) = \begin{bmatrix} -2 \\ 1 \end{bmatrix} \quad R(A) = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad \dim(R(A)) = 1$$

$$\text{nullity} = 1 \quad N(A^T) = \begin{bmatrix} -3 \\ 1 \end{bmatrix} \quad \dim(N(A^T)) = 1$$

- L 3.2 Orthogonal Vectors & Subspaces
- length of a vector $\|x\|^2 = x_1^2 + x_2^2$
 - same for general n -dim. space $\|x\|^2 = x_1^2 + x_2^2 + \dots + x_n^2$
 - Orthogonal $\rightarrow x \perp y \text{ if } x^T y = 0$
 - $x^T y = \sum_{i=1}^n x_i y_i$ inner prod.
 - 0 is orthogonal to every x
 - if (u, \dots, u_n) are mutually orthogonal, then this is a linearly indep. set
 - Orthonormal \rightarrow orthogonal & their unit lengths are 1
 $(u, u) = 1$ if $u^T u = 1$ & $\|u\| = 1$
 - Orthogonal subspaces
 U, V if $A|_{U \times V} = 0$
if $x^T y = 0 \forall x \in U, y \in V$
 - Orthogonality w.r.t A has fundamental spaces
 - ① $R(A) \perp N(A)$ since $0 \in R(A)$
 - ② $C(A) \perp N(A^T)$ since $0 \in C(A)$
- L 3.3 Projections $(A, b) \rightarrow (A)u \leftarrow A$
- projection onto a line

$$P = \hat{x} \hat{x}^T$$

$$\hat{x} = \hat{a}^T b$$

$$\hat{x} = \frac{\hat{a}^T b}{\hat{a}^T \hat{a}} \hat{a}$$
 - Cauchy-Schwarz Inequality
 $|a^T b| \leq \|a\| \|b\|$

L 3.4 Least sq. & projections onto a subspace

Eq. $a = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, A = P = \frac{a a^T}{a^T a} = \frac{1}{3} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix}$

- P is symmetric
- $P^T = P$
- $C(P) \Rightarrow a = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$
- $N(P) \Rightarrow$ plane orthogonal to a
- Rank $\text{r}(P) = 1$

L 3.4 Least sq. & projections onto a subspace

$2x = b_1$
 $3x = b_2$
 $4x = b_3$: sys. is solvable if b is on line $\begin{bmatrix} 2 \\ 3 \\ 4 \end{bmatrix}$

minimize aug. error
 $E^2 = (2x - b_1)^2 + (3x - b_2)^2 + (4x - b_3)^2$

$$\frac{dE^2}{dx} = 0$$

$$\Rightarrow \hat{x} = \frac{2b_1 + 3b_2 + 4b_3}{2+3+4} = \frac{a^T b}{a^T a}$$

\rightarrow Taking deviation & finding the minimum turns out to be as same as performing an projection

$Ax = b$, A is $m \times n$, $m \geq n$

\Rightarrow projection of b onto $C(A)$

$|A^T A \hat{x} = A^T b| \Rightarrow$ this leads to \hat{x}
 which minimizes $\|Ax - b\|^2$
 \rightarrow columns of A are linearly indep.
 Then, $A^T A$ is invertible

$\rightarrow b \in C(A)$ i.e., $b = Ax$

$$\phi = b$$

$\rightarrow b \in N(A^T)$

$\rightarrow A$ is eq. & invertible $\Leftrightarrow C(A) = \mathbb{R}^n$

$$\phi = b$$

$\rightarrow A$ is rank 1

$$\hat{x} = \frac{a^T b}{a^T a}$$

- Projecⁿ matrix,

$$P = A(A^T A)^{-1} A^T$$

symmetrce

$$P^T = P$$

$$P^2 = P$$

So, projecⁿ matrix is sym. & satisfies $P^2 = P$,
as the converse is also true.

L3.5 Example of Least Squares

- 1D

$$(x_1, b_1), \dots, (x_m, b_m)$$

$$b_i = \alpha x_i + \theta' \quad i = 1, \dots, m \quad \text{--- (1)}$$

$$A\theta = b \quad \text{where } \theta = \begin{bmatrix} \theta' \\ \theta'' \end{bmatrix}$$

$A\theta = b$ may be inconsistent

\therefore minimize $E^2 = \|b - A\theta\|^2$.

$$(\hat{\theta}', \hat{\theta}'') = \text{arg min } \|b - A\theta\|^2$$

$$4. \quad \begin{bmatrix} -1 \\ 1 \\ 2 \\ 1 \end{bmatrix} \begin{bmatrix} \theta' \\ \theta'' \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 3 \end{bmatrix} \quad \text{not consistent}$$

since
 $b \notin C(A)$

$$A^T A \hat{\theta} = A^T b$$

$$A^T A = \begin{bmatrix} 6 & 2 \\ 2 & 3 \end{bmatrix}$$

$$A^T A \hat{\theta} = A^T b$$

$$\theta' = \frac{4}{7} \quad \& \quad \theta'' = \frac{9}{7}$$

$$\hat{\theta} = \begin{bmatrix} 4/7 \\ 9/7 \end{bmatrix}$$

$$\text{line } \Rightarrow \frac{4}{7}x + \frac{9}{7}$$

$$\text{Projections: } P_1 = \frac{4}{7}(+1) + \frac{9}{7} = \frac{5}{7},$$

$$\text{distance w.r.t } P_1 = \frac{12}{7} \quad (A^T A)^{-1} = \frac{17}{7}$$

original data is not on line, $E^2 > 0$

$$E^2 = \|b - A\hat{\theta}\|^2, c = \left(\frac{-21}{7}, \frac{6}{7}, \frac{4}{7} \right)$$

$$\text{new distance w.r.t } P_2 = \frac{1}{7} \quad (A^T A)^{-1} = 1$$

parallel line

$$(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2$$

minimum distance

when θ is

: wanted line equation

$$L_\theta = (x_1, \theta), (x_2, \theta), \dots, (x_n, \theta) = (x, \theta)$$

$$(x, \theta)^T S = (x, \theta)^T P$$

newly measured, wanted with prior

$$(A^T A)^{-1} = \frac{1}{17} I_3$$

Week 4

14.1 Linear & Polynomial Regression

- $(x_1, y_1) \dots (x_n, y_n)$ $x_i \in \mathbb{R}^d$ $y_i \in \mathbb{R}$

$$L(\theta) = \frac{1}{2} \sum_{i=1}^n (x_i^\top \theta - y_i)^2$$

Minimizing L : define $A = \begin{bmatrix} x_1^\top \\ x_2^\top \\ \vdots \\ x_n^\top \end{bmatrix}$ ← Feature matrix

$$\text{so, } L(\theta) = \frac{1}{2} \|A\theta - y\|^2$$

$$\Rightarrow \theta = \text{least eq. soln} \Rightarrow (A^\top A)\theta = A^\top y$$

$$\theta = (A^\top A)^{-1} A^\top y \Rightarrow \text{if } A \text{ is full rank}$$

then $A^\top A$ is invertible

- Maximum likelihood & least eq -

\rightarrow Model $\rightarrow y = \theta^\top x + \epsilon$ ← zero-mean noise

$D = \{(x_i, y_i), i=1 \dots n\}$ Err gaussian with mean iid θ var $\frac{1}{\sigma^2}$

Max. likelihood

$$L(\theta) = \prod_{i=1}^n \frac{\sqrt{\sigma}}{\sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma^2} (y_i - \theta^\top x_i)^2\right)$$

- Polynomial regression

1 D data

Transformed features:

$$\phi(x) = \sum_{j=0}^m \phi_j(x), \quad \phi_j(x) = x^j$$

$$\hat{y}(x) = \theta^\top \phi(x)$$

using these features, perform linear regression

$$(A^\top A)\theta = A^\top y$$

$$L(\theta) = \frac{1}{2} \sum_{i=1}^n (\theta^\top x_i - y_i)^2 + \lambda \|\theta\|^2$$

minimized = least squares for above "fit"

14.2 Eigenvalues & Eigenvectors

- ordinary differentiated matrices

$$\text{single eq. } \frac{du}{dt} = Au \quad \begin{cases} 1 & t=0 \\ 1 & t=1 \end{cases}$$

$$\text{soln } \Rightarrow u(t) = e^{ut} u(0)$$

$\begin{cases} \text{is unstable if } \lambda > 0 \\ \text{is neutral stable if } \lambda = 0 \\ \text{is stable if } \lambda < 0 \end{cases}$

$$\begin{aligned} & - 2D - \text{parameters not written} \\ & \text{top row } u(t) = e^{\lambda t} y, \text{ middle row } u(t) = e^{\lambda t} \begin{bmatrix} 1 & 2 \end{bmatrix}, \\ & \text{bottom row } u(t) = e^{\lambda t} x, \text{ where } x = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \\ & \boxed{Ax = \lambda x} \Rightarrow \text{eigenvalue eqn} \end{aligned}$$

- For a matrix A , λ is an eigenvalue & $x \neq 0$ is an eigenvector if $Ax = \lambda x$

if x is an eigenvector, then A is called eigenvectors / shrinks it, but A

does not \rightarrow the direction of x is not changed

$$\text{if } \lambda = 0: Ax = 0 \Rightarrow x \in N(A)$$

$$\text{Eq. } B = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad x = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

- Finding the eigenvalues:

$$\text{Find } Ax = \lambda x \Rightarrow (A - \lambda I)x = 0 \text{ has non-trivial solns} \Rightarrow \det(A - \lambda I) = 0$$

charac. poly. of matrix A

If A is $n \times n$, the dgf (charac. poly.) λ

" n " roots of charac. poly. = eigenvalues

For eigenvectors its linearly independent -

$$\text{Ex. } A = \begin{bmatrix} 3 & 1 \\ 1 & 3 \end{bmatrix} \quad \det(A - \lambda I) = 0 \quad \text{piv 2} \\ \lambda^2 - (3+1)\lambda + 8 = 0 \quad \text{piv 2} \\ (\lambda - 5)(\lambda - 1) = 0 \quad \lambda_1 = 5, \lambda_2 = 1$$

$$x_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, x_2 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

- Sym. matrix has real eigenvalues
- a 2×2 matrix where we don't get 2 independent eigenvalues

$$\text{Ex. } A = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

Diagonalization of matrix

A matrix A is diagonalizable if there is a invertible matrix S such that

$$S^{-1}AS = D \text{ where } D \text{ is a diag. matrix}$$

- This is tried by having enough independent eigenvectors x_1, x_2, \dots, x_n : $x_i = \lambda_i p_i$
- if λ_1, λ_2 are the eigenvalues with corresp. eigenvectors x_1, x_2 & $\lambda_1 \neq \lambda_2$, then $\{x_1, x_2\}$ is a linearly indep. set

and extensions (if λ_1, λ_2 are distinct, then eigenvectors are linearly indep.)

$$\text{Ex. } A = \begin{bmatrix} 1 & 2 & 4 \\ 2 & 3 & 1 \end{bmatrix} \quad A - \lambda I = \begin{bmatrix} 1-\lambda & 2 & 4 \\ 2 & 3-\lambda & 1 \end{bmatrix}$$

Is A diagonalizable? $\lambda^2 - 4\lambda - 5 = 0$
Yes, it has 2 distinct eigenvalues

$$\lambda_1 = 5 \quad A - \lambda_1 I = \begin{bmatrix} -4 & 2 \\ 2 & -2 \end{bmatrix} \quad x_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\lambda_2 = -1 \quad A - \lambda_2 I = \begin{bmatrix} 2 & 4 \\ 2 & 4 \end{bmatrix} \quad x_2 = \begin{bmatrix} -2 \\ 1 \end{bmatrix}$$

$$S = \begin{bmatrix} 1 & -2 \\ 1 & 1 \end{bmatrix} \quad \lambda = \begin{bmatrix} 5 & 0 \\ 0 & -1 \end{bmatrix}$$

Diagonalization is applied
eigenvectors are matched

$$\Rightarrow S^{-1}AS = \lambda \quad \checkmark$$

$$S^{-1}AS = \lambda \quad \text{or} \quad A = SAS^{-1} \quad \text{Is } S \text{ unique? No}$$

$$A = SAS^{-1} \quad \text{for example, } A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Powers of A : use A is diagonalizable

$$A^2 = A(Ax) = \lambda Ax = \lambda^2 x$$

$$S^{-1}A^2S = \lambda^2 \quad \Rightarrow \text{Yes}$$

The easiest form a general $k \geq 1$

- Not all matrices are diagonalizable

L4.4 Solving Fibonacci Seq. using "diagonalization"

Fibonacci

$$F_{k+2} = F_{k+1} + F_k \quad \text{what is } F_{100}?$$

$$u_k = A^k u_0 \quad \text{set } u_0 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad u_{k+1} = A u_k$$

$$A = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$$

$$\lambda^2 - \lambda - 1 = 0$$

$$\lambda_1 = \frac{1+\sqrt{5}}{2}, \lambda_2 = \frac{1-\sqrt{5}}{2}$$

$$x_1 = \begin{bmatrix} \lambda_1 \\ 1 \end{bmatrix}, x_2 = \begin{bmatrix} \lambda_2 \\ 1 \end{bmatrix}$$

$$F_k \approx \frac{1}{\sqrt{5}} \left(\frac{1+\sqrt{5}}{2} \right)^k, F_{100} \approx \frac{1}{\sqrt{5}} \left(\frac{1+\sqrt{5}}{2} \right)^{100}$$

helps to understand linear recurrence relation — diagonalization

L4.5.1 Orthogonally diagonalizable matrices

- A is a real sym. matrix, then $\exists Q \in \mathbb{R}^{n \times n}$
- ① eigenvalues of A are real
- ② eigenvectors λ_i corresponds to diff. eigenvalues are linearly indep.

③ A is orthogonally diagonalizable

$$A = Q \Lambda Q^T, Q^T Q = I$$

- any real matrix is not necessarily diagonalizable, but a real sym. is diagonalizable

$$A \in \mathbb{R}^{n \times n} \text{ at } A^T A = I$$

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

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

Week 5

L5.1 Complex Matrices

- $i \in \mathbb{C}^n$: complex counterpart of \mathbb{R}^n
- $x_1, \dots, x_m \in \mathbb{C}^n$
- Add \mathbb{C}^n : $(a+ib) + (c+id) = (a+c) + i(b+d)$
- Multiplication: $(a+ib)(c+id) = (ac-bd) + i(bc+ad)$
- complex conjugate of $(a+ib)$ is $(a-ib)$
- Linear combination: $c_1 x_1 + c_2 \dots + c_n x_n = 0$
- In \mathbb{R}^n ,

$$x \cdot y = \overline{x}^T y = \overline{x}_1 y_1 + \dots + \overline{x}_n y_n$$

Also, $\overline{x}^T y \neq y^T x$

- Length of a complex vector in \mathbb{C}^n
- For $x \in \mathbb{C}^n$, $\|x\|^2 = \overline{x}^T x$
- $x \cdot y = \overline{y} \cdot x$ from defn of inner prod.
- $x \cdot (cy) = c(x \cdot y)$ using prod.
- $(cx) \cdot y = \overline{c}(x \cdot y)$
- $x \cdot y = \overline{\overline{x}^T y} = \sum_{i=1}^n \overline{x}_i y_i$

Conjugate transpose:

$$\begin{aligned} 1. & A^* = \text{conj. transpose of } A \\ 2. & A^{**} = \overline{A^T} = iA^T \\ 3. & (A^*)^* = A \\ 4. & (AB)^* = B^* A^* \end{aligned}$$

Real matrix: $A \in \mathbb{R}^{n \times n}$ case

$$\begin{aligned} \rightarrow (A^*)^* &= A & \rightarrow (A^T)^T &= A \\ \rightarrow (AB)^* &= B^* A^* & \rightarrow (AB)^T &= B^T A^T \end{aligned}$$

$$\boxed{x \cdot y = \overline{x}^* y}$$

- A matrix is Hermitian if $A^* = A$
- These matrices are similar as similar matrices in complex vector space

L5.2 Hermitian matrices

Ex. $A = \begin{bmatrix} 2 & 3-3i \\ 3+3i & 5 \end{bmatrix}$ is Hermitian $\Rightarrow A^* = \overline{A^T} = \begin{bmatrix} 2 & 3+3i \\ 3+3i & 5 \end{bmatrix}^T$

$$(2+i) + (3-i) = (2-i) + (3+i) = A$$

$$(\text{Hermitian}) + (\text{Hermitian}) = (\text{Hermitian}) + (\text{Hermitian})$$

$$(di - s) \text{ is } (di + s) \text{ for transpose relation}$$

- diagonal entries of hermitian matrix are real \Rightarrow eigenvalues are real

- Properties -

① If A is hermitian, all eigenvalues are real.

Ex. $|A - \lambda I| = \begin{vmatrix} 2-\lambda & 3-3i \\ 3+3i & 5-\lambda \end{vmatrix} \Rightarrow \lambda^2 - 7\lambda - 8$

$$\lambda = 8, -1$$

② If A is hermitian, then eigenvectors corresponding to diff. eigenvalues, are orthogonal \Rightarrow thus, linearly independent.

$$\text{Let } A \vec{x} = \lambda_1 \vec{x}, \vec{x} \neq 0 \Rightarrow (\lambda_1 - \lambda_2) \vec{x} = \vec{x}^T \vec{y} = 0$$

$$\therefore A \vec{y} = \lambda_2 \vec{y}$$

Ex. $(A - 8I) \vec{x} = \begin{bmatrix} -6 & 3-3i \\ 3+3i & -3 \end{bmatrix} \begin{bmatrix} \vec{x}_1 \\ \vec{x}_2 \end{bmatrix} = 0 \quad \vec{x} = \begin{bmatrix} 1 \\ i+1 \end{bmatrix}$

$$(A + iI) \vec{y} \Rightarrow \vec{y} = \begin{bmatrix} 1-i \\ -1 \end{bmatrix}$$

$$A \vec{x} \cdot \vec{y} = (1)(1-i) \begin{bmatrix} 1 \\ i+1 \end{bmatrix} = 0$$

- The equivalent of hermitian matrix is real symmetric matrix

All "real" matrices are hermitian

- If eigenvalues are distinct, the A is diagonalizable

L5.3

Unitary Matrices \Rightarrow inverse is itself

A matrix is unitary, if its square & has orthonormal columns, i.e.,

$$\text{Real: } Q^T Q = I \Rightarrow Q \text{ is orthogonal case}$$

$$\text{Complex: } U^* U = I \Rightarrow U \text{ is unitary}$$

$$U = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \sqrt{2}/2 \end{bmatrix}$$

Ex. $U_h = \begin{bmatrix} \frac{1}{\sqrt{2}} + \frac{i}{\sqrt{2}} \\ \frac{i}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} - \frac{i}{\sqrt{2}} \end{bmatrix} - h = (h)$

- Properties - $U \rightarrow$ Unitary matrix

① "length unchanged" $\|Ux\| = \|x\|$

② Eigenvalues of U have absolute value 1
then, $\|\lambda\| = 1$

③ Eigenvectors corresponding to different eigenvalues of a unitary matrix U are orthogonal (because $A = UT$ satisfies orthonormality condition)

- For a hermitian matrix A , we can find unitary matrix U

$$A = U \lambda U^*, \quad \lambda \Rightarrow \text{diagonal matrix with eigenvalues of } A$$

$$T = U \lambda^* U = U^* \lambda^* U = \lambda^*$$

orthogonal matrix

Diagonalisation of Hermitian Matrix - I

$$A^* = \overline{A^T A} = \overline{A} = A$$

Matrix A is hermitian $\Rightarrow A^* = A$

Matrix U is unitary $\Rightarrow U^* U = I$ & U is Sq.

- A matrix 'A' is unitarily diagonalizable if there exists a unitary matrix U

$$A = U \lambda U^*$$

- Schur's Theorem - Any non-singular matrix is similar to upper triangular matrix T , there exists unitary matrix U such that $T = U^* A U$ (here $A = T^* T$)

position of U as $U = U^*$ (natural)

$$\text{Eq. } A = U \begin{bmatrix} 5 & 0 & 0 \\ 0 & 9 & 0 \\ 0 & 0 & -10 \end{bmatrix}$$

$$p(\lambda) = -(\lambda - 1)(\lambda + 3)^2 \quad \lambda_1 = 1, \lambda_2 = -3$$

$$\text{Eigenvector } z_1 = \begin{bmatrix} -2 \\ 1 \\ 1 \end{bmatrix}$$

another position $\Rightarrow U$

$$\text{Hilf. } T = U^* A U = \begin{bmatrix} 1 & -1/\sqrt{2} & -12/\sqrt{3} \\ 0 & 9 & 0 \\ 0 & 0 & -10 \end{bmatrix}$$

$$\text{euler standard } \begin{bmatrix} 1 & -2 & \sqrt{6} \\ 0 & 4 & 0 \\ 0 & 0 & -3 \end{bmatrix}$$

inverses of eigenvectors are unique

- " L.B.S. " Part 2 meeting 2 for A is Hermitian

- Spectral Th. - A hermitian matrix "A" is unitarily diagonalizable being unitary $U^* A U = D$ maintains \Rightarrow U is unitary \Rightarrow $U^* = U^{-1}$ \Rightarrow $D = U^* A U$ \Rightarrow $U^* D U = A$ \Rightarrow $U^* U = I \Rightarrow I = U U^* = I$ \Rightarrow $U^* = U^{-1}$ \Rightarrow U is unitary

$$T^* = U^* A^* U = U^* A U = T$$

Lower triangular

Upper triangular

$$\text{Eq. } A = \begin{bmatrix} 2 & 1-i \\ 1+i & 3 \end{bmatrix} \quad A^* = A^T = A$$

$$p(\lambda) = (\lambda - 1)(\lambda - 4) \Rightarrow \lambda_1 = 1, \lambda_2 = 4$$

$$z_1 = \begin{bmatrix} -1+i \\ 1 \end{bmatrix} \quad z_2 = \begin{bmatrix} 1-i \\ 2 \end{bmatrix}$$

Spectral Th.

- A real symmetric matrix A is orthogonally diagonalizable, there exists a matrix Q such that $Q^T A Q = D$ where $Q^T Q = I$ (i)

$$\text{diag. } Q^T A Q = D \text{ where } Q^T Q = I \text{ (ii)}$$

\hookrightarrow diag. matrix with real nos, eigenvalues and it has either positive or negative entries

- Hermitian \Rightarrow unitarily diagonalizable

$$\cancel{\Leftrightarrow} \quad \beta \circ \beta = A$$

but, the reverse is not true

counterexample when $\beta \neq I$

$$\begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} = \beta^2 \text{ where } \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \neq \beta$$

$$\beta \circ \beta = A$$

$$\beta \circ \beta = \beta \circ \beta = \beta \circ \beta = A$$

but $\beta \circ \beta = A$ for non-diagonalizable matrix

$$\beta \circ \beta = \beta \circ \beta = A$$

but $\beta \circ \beta = A$ for non-diagonalizable matrix

$$\left\{ \text{diag. } \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = A \right\}$$

but $\beta \circ \beta = A$ for non-diagonalizable matrix

(counterexample) $\beta \circ \beta = A$ for non-diagonalizable matrix

L6.1 SVD (Singular Value Decomposition)

A is an $n \times n$ real sym. matrix w.r.t.

(i) All eigenvalues are real

(ii) A is orthogonally diagonalizable

$$A = Q_1 \Sigma Q_2^T$$

Every matrix can't be diagonalized, but, any real $n \times n$ matrix can be decomposed into SVD form.

$$A = Q_1 \Sigma Q_2^T$$

$m \times n$ $m \times m$

Q_1, Q_2 are orthogonal

$$\Sigma = \begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix}, \text{ where } D = \begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_m \end{bmatrix}$$

①

$$A = Q_1 \Sigma Q_2^T$$

$$A A^T = Q_1 \Sigma \Sigma^T Q_2^T$$

eigen decomposition of $A A^T$ into Q_1 ,

$$② A^T A = Q_2 \Sigma^T \Sigma Q_2^T$$

eigen decomposition of $A^T A$ into Q_2

L6.2 Eg. of SVD

$$\text{Ex. } A = \begin{bmatrix} \sqrt{2} & 1 \\ 0 & \sqrt{2} \end{bmatrix}$$

Find SVD?

Is matrix A diag.? No

Eigenvalue = $\sqrt{2}, \sqrt{2}$ (repeated)

$$A - \sqrt{2} I = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 0 \end{bmatrix} \text{ eigenvector}$$

$\lambda = 0 \neq \sqrt{2}$. & they aren't linearly indep.

Find SVD of $\begin{bmatrix} \sqrt{2} & 1 \\ 0 & \sqrt{2} \end{bmatrix}$.

(i) Eigenvalues & eigenv. of $A^T A$

$$(ii) \sigma_1 y_1 = A x_1, \sigma_2 y_2 = A x_2$$

to find y_1, y_2

$$(iii) Q_1 = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}, Q_2 = \begin{bmatrix} 1 & 1 \\ \sqrt{2} & \sqrt{2} \end{bmatrix}, \Sigma = \begin{bmatrix} \sqrt{\lambda_1} & 0 \\ 0 & \sqrt{\lambda_2} \end{bmatrix}$$

$$(i) A^T A = \begin{bmatrix} 2 & \sqrt{2} \\ \sqrt{2} & 3 \end{bmatrix}, \lambda_1 = 4, \lambda_2 = 1$$

$$\sigma_1 = \sqrt{\lambda_1} = 2, \sigma_2 = \sqrt{\lambda_2} = 1$$

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

$$\text{Eigenvectors of } A^T A \Rightarrow \begin{bmatrix} -2, \sqrt{2} \\ \sqrt{2}, -1 \end{bmatrix}, \begin{bmatrix} 1, \sqrt{2} \\ \sqrt{2}, 2 \end{bmatrix}, \begin{bmatrix} \sqrt{2} \\ -1 \end{bmatrix}$$

$$\text{Normalize, } x_1 = \frac{1}{\sqrt{3}} \begin{bmatrix} 1 \\ \sqrt{2} \end{bmatrix}, x_2 = \frac{1}{\sqrt{3}} \begin{bmatrix} \sqrt{2} \\ 1 \end{bmatrix}$$

$$Q_2 = \frac{1}{\sqrt{3}} \begin{bmatrix} 1 & \sqrt{2} \\ \sqrt{2} & -1 \end{bmatrix}$$

$$(ii) y_1 = \frac{1}{\sqrt{3}} \begin{bmatrix} \sqrt{2} \\ 1 \end{bmatrix}, y_2 = \frac{1}{\sqrt{3}} \begin{bmatrix} 1 \\ \sqrt{2} \end{bmatrix}$$

(iii)

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

- L6.3 +ve definiteness: $f = ax^2 + bxy + cy^2$
- $f(x,y) = 2x^2 + 4xy + y^2$
 - $\frac{\partial f}{\partial x} = 4x + 4y = 0$ and $\frac{\partial f}{\partial y} = 4x + 2y = 0$
 - $\begin{cases} 4x + 4y = 0 \\ 4x + 2y = 0 \end{cases}$ gives $\frac{\partial^2 f}{\partial x^2} = 4$, $\frac{\partial^2 f}{\partial x \partial y} = 4$, $\frac{\partial^2 f}{\partial y^2} = 2$
 - $(x,y) = (0,0)$: both $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}$ are 0

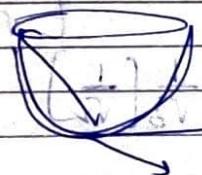
Hence, $(0,0)$ is a stationary point of

$$\frac{\partial^2 f}{\partial x^2} = 4, \quad \frac{\partial^2 f}{\partial x \partial y} = 4, \quad \frac{\partial^2 f}{\partial y^2} = 2$$

$\Rightarrow f$ has a minima at $(0,0)$

- Every quadratic funcⁿ of the form $(ax^2 + 2bxy + cy^2)$ has a stationary point at $(0,0)$

- A function f that vanishes at $(0,0)$ and is strictly +ve at other pts is called positive definite. $\Rightarrow f > 0$



$$(ax^2 + 2bxy + cy^2)$$

Conditions -

- ① if $f > 0$, then $a > 0$ (i)
- ② if $f > 0$, then $c > 0$ (ii)

Further,

$$f(x,y) = a \left(x + \frac{b}{a} y \right)^2 + \left(c - \frac{b^2}{a} \right) y^2 \geq 0$$

- ③ if $f > 0$, $ac > b^2$

Combining all these 3 conditions, $f(x,y) = ax^2 + 2bxy + cy^2$ is +ve definite iff. 2. $a > 0$ & $ac > b^2$

- if $ac = b^2$, then it is semi-def. $a > 0$
- we have a saddle point at $(0,0)$ if $ac < b^2$

$$ax^2 + 2bxy + cy^2 = [x \ y] \begin{bmatrix} a & b \\ b & c \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

$$u = \begin{bmatrix} x \\ y \end{bmatrix}, \quad A = \begin{bmatrix} a & b \\ b & c \end{bmatrix} \Rightarrow ax^2 + 2bxy + cy^2 = u^T A u$$

$$\textcircled{1} \quad f(x,y) = 2x^2 + 4xy + y^2 \Rightarrow \text{saddle pt. at origin}$$

$$\textcircled{2} \quad f(x,y) = 2xy \Rightarrow \text{saddle at origin}$$

L6.4 +ve definite matrices.

- $A = \begin{bmatrix} a & b \\ b & c \end{bmatrix}$ is +ve def. if $a > 0$, $ac - b^2 > 0$

concs. funcⁿ in $f(u) = u^T A u$ $u = \begin{bmatrix} x \\ y \end{bmatrix}$

- if $a > 0$, $ac - b^2 > 0$, then both eigenvalues are +ve

$$ac - b^2 = \det(A) = \lambda_1 \lambda_2 > 0$$

$$\text{Trace}(A) = \lambda_1 + \lambda_2 = a + c > 0$$

- a real sym. matrix is +ve def. if
 - (i) $u^T A u > 0 \quad \forall u \in \mathbb{R}^n, u \neq 0$
 - (ii) all eigenvalues of A are > 0

Week-7

L7.1 Principal Comp. Analysis (PCA)

- PCA (7 dim) - Project given data onto a lower dimensional subspace such that -
 - reconstruction error is minimized
 - covariance of proj. data is minimized

$$(x_i)^T = \sum_{j=1}^m (x_i^T u_j) u_j \text{ and } w =$$

$$[x_i^T] = \frac{1}{n} \sum_{i=1}^n = \|x_i - \bar{x}\|_F^2 + \dots$$

$$\sim [w] = A$$

$$x_i = \sum_{j=1}^m z_{ij} u_j + \sum_{j=m+1}^n (x_i^T u_j + \beta_j) u_j$$

$$\text{Lagrange} \leftarrow \mu + \mu \lambda + \omega \omega = (\mu, \omega)$$

(cont.)

$$L(u, \lambda) = u^T (u + \lambda (I - u^T u))$$

$$D_u L(u, \lambda) = 0 \Rightarrow Cu = \lambda u \quad (\text{as } u)$$

$$C = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$

where PCA used mult. $\lambda < 0$, $\lambda > 0$ if $\lambda < 0$

$$\textcircled{1} \text{ data : } \{x_1, \dots, x_n\}$$

$$\textcircled{2} \text{ let } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \text{ tel. } \mathbb{R}^d$$

$$C = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$

$$\textcircled{3} \text{ find eigenvalues with scores (i)} \\ \text{or eigenvectors with right (ii)}$$

L7.2 (4) projected data

$$x_i = \sum_{j=1}^m (x_i^T u_j) u_j + \sum_{j=m+1}^n x_i^T u_j$$

L7.3 PCA is maximizing variance
 x_i projection along (u) is $(x_i^T u) u$

$$D_x = \{x_1, \dots, x_n\}, \text{ mean } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

variance between x_i and \bar{x} is $(x_i^T u) u$, mean $= (\bar{x}^T u) u$
 variance is $(x_i^T u - \bar{x}^T u)^2$

$$\text{maximize } \frac{1}{n} \sum_{i=1}^n (x_i^T u - \bar{x}^T u)^2$$

A calculus argument,

$$\max_u \frac{u^T C u}{u^T u}$$

In vector form,

$$u^T C u = (u^T C u) u \Rightarrow C u = \left(\frac{u^T C u}{u^T u} \right) u$$

$$\text{or } C u = \lambda u$$

To maximize $\frac{u^T C u}{u^T u}$, choose λ to be the

largest eigenvalue of C & let w be the corresponding eigenvector.

L7.4 PCA in higher dimensions

- PCA eqn: finding the eigenvectors of

$$C = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T,$$

a $d \times d$ matrix

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

rank (C) $\leq n$ if

λ_i is an eigenvalue of $\frac{1}{n} A^T A$

so, instead of working with $d \times d$ matrix
 $C = \frac{1}{n} A^T A$ to find its eigenvalues/vectors,
it is enough to find eigenvalues/vectors

of $A^T A$, which is a $n \times n$ matrix.

transposed called A^T

new matrix

$A^T A$

new matrix

$$(w_1^T w_1) = w_1 \cdot w_1 = w_1 (w_1^T w_1) = w_1^T w_1$$

add at k points, will minimize at

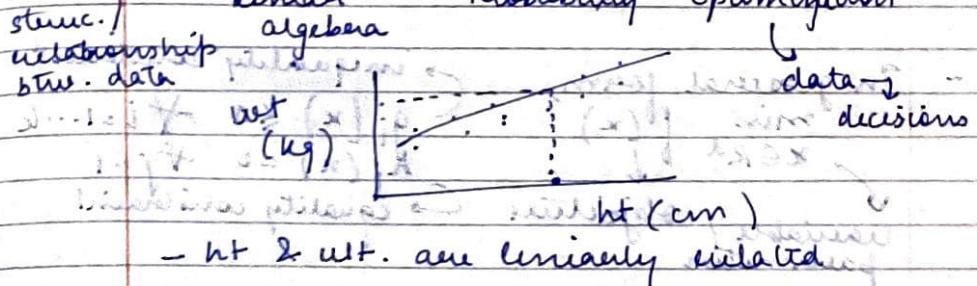
and want to minimize the sum of squares of distances from a central point

Week 8

L8.1 Pillars of ML

noise / variance / uncertainty

P_1 linear algebra
 P_2 probability
 P_3 optimization

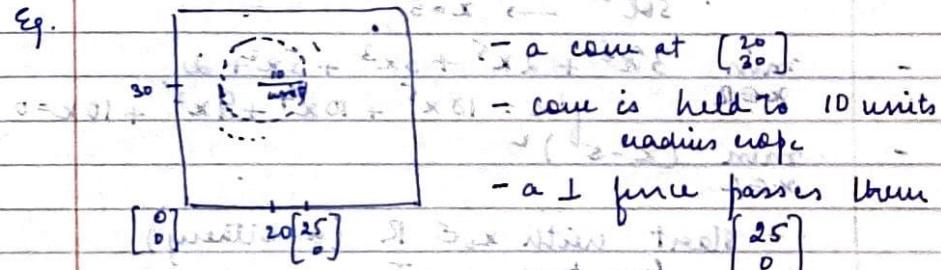


- ht & wt. are linearly related

L8.2 Intro to Optimization

we care about finding the 'best' classifier
 or least 'loss'
 or maximum 'reward'

Eg.



- grass on the field at $[40, 40]$

"How close the cow get to the grass?"

- dist. from grass $\left[\begin{array}{c} 40 \\ 40 \end{array} \right]$ } $d = (x_1 - 40)^2 + (x_2 - 40)^2$
 say the cow is $\left[\begin{array}{c} x_1 \\ x_2 \end{array} \right]$

minimize d or make this as small as possible

$$\text{min} - (x_1 - 40)^2 + (x_2 - 40)^2$$

b. i. f. $\left[\begin{array}{c} x_1 \\ x_2 \end{array} \right]$ based on d "min" or
 "close" to a ad. fence

$$\text{Rope constraint} \quad (x_1 - 20)^2 + (x_2 - 30)^2 \leq 10^2$$

Fenus nescie"
wait (approximately) x 15 25 min?

- In general form, $\min_{x \in \mathbb{R}^d} f(x)$ \rightarrow inequality constraints $g_i(x) \leq 0 \quad \forall i=1 \dots k$
 $\quad\quad\quad$ equality constraints $h_j(x) = 0 \quad \forall j=1 \dots l$

L8.3 Solving an unconstrained opt. problem - I

$$\rightarrow \text{Unconstrained Opt.}$$

$$\text{LHS} = \text{RHS} \rightarrow x = 5$$

$$\min_{x_1, x_2} \quad 3x_1^6 + 2x_1^5 + 3x_1^3 + 5x_1^2 + 2$$

$$f(x) = 18x^5 + 10x^4 + 9x^2 + 10x = 0$$

$$x \in \mathbb{R} \text{ where } t = -$$

Start with $x_0 \in R$ (arbitrarily)
for $t = 1, \dots, T$

- update $\stackrel{O(2)}{\rightarrow}$ to ~~list~~ \Rightarrow decision -

! change soft | $x_{t+1} = x_t + d$ | now most "good direction"

Given x , what is a good direction?

if $x > 5$ then $\leftarrow d < 0$

\Rightarrow direction d will depend [on] x , i.e., d must be a funcⁿ of x .

$$f'(x) = 2(x-5) \Rightarrow 2(8-5) = 6$$

$$\begin{aligned} f'(x) &= 2(x-5) \Rightarrow f'(x) > 0 \\ \Leftrightarrow x &> 5 \\ \dots \Rightarrow x &< 5 \Rightarrow f'(x) < 0 \end{aligned}$$

choose, $d = -f'(x)$
 $f(x) = (x-5)^2$ $f'(x) = 2(x-5)$ $d = -f'(x)$
 but, in the value oscillates from 0 to 10 but atleast the direcⁿ is giving good answers towards minima

L8.4 Part -2

$$x_{t+1} = x_t - \eta_t f'(x_t)$$

\rightarrow step size (scalar quantity +ve)

Now, have to choose a step size

$$1^{\text{st}} \text{ attempt } n_t = 1, m_1 = 1 \Rightarrow \dots \rightarrow n_t = \frac{1}{2^t} + \frac{1}{2^{t+1}} X$$

$$2^{\text{nd}} \text{ attempt} \cdot \text{ takes } 1, \frac{1}{2}, \frac{1}{3}, 1 \frac{1}{4}, \dots$$

$$\dots + (x_1^{t+1})^2 \dots + (x_n^{t+1})^2 \Bigg) \left| \frac{\partial y_t}{\partial x_i} + \frac{1}{t+1} \right| \rightarrow (\text{good step}) \text{ or } (\text{bad step})$$

IND ~~in~~ ~~in~~ ~~in~~ ~~in~~

卷之三

Basic Alg: Gradient Descent

minimum f(x) at zero

Alg - Gradient Descent (1^{st} Order Alg.)

Initializ at $x_0 \in \mathbb{R}$

$f'(x)$ for $t = 1, 2, \dots$

$$x_{t+1} = x_t - \eta_t f'(x_t) \quad \text{where } \eta_t = \frac{1}{t+1}$$

Properties

① if $\eta_t = 1/t$, the alg converges

principally at $t+1$ and $t+2$

② gradient descent converges to local minima

- there are funcs in which local minima
- convex functions
- global minima

18.6 Gradient descent & Taylor series

$$x_{t+1} = x_t - \eta_t (-f'(x_t))$$

what is so spcl. abt. this?

Taylor series

$$\begin{aligned} f(x + \eta d) &= f(x) + \eta d | f'(x) + \frac{\eta^2 d^2}{2} f''(x) + \dots \\ x &= x + \eta d \end{aligned}$$

The local info. gives evaluations at all at x .

$f(x + \eta d)$
↓ small dirce
the further higher
the step size

$$\begin{aligned} f(x + \eta d) &= f(x) + \eta d f'(x) \\ \Rightarrow f(x + \eta d) - f(x) &\approx \eta d f'(x) \end{aligned}$$

func evalua at func evalua at current
updated pt. along. pt.
dirce d

want to choose a
dirce s.t. we want this to be negative

$$\begin{aligned} \eta d f'(x) &< 0 \Rightarrow d f'(x) < 0 \\ \text{small } \eta & \text{ const.} \\ d f'(x) &= -(f'(x))^2 < 0 \end{aligned}$$

18.7 Gradient descent for multivariate func

- let's go to higher dimensions

$$f(x_1, x_2) = x_1^2 + 4x_2 + 8x_2^2$$

Derivatives \Leftrightarrow Gradient

(vector of partial derivatives)

$$\nabla f \begin{pmatrix} a \\ b \end{pmatrix} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \Big|_{x_1=a} \\ \frac{\partial f}{\partial x_2} \Big|_{x_2=b} \end{bmatrix}$$

$$\text{Ex. } \nabla f \begin{pmatrix} 1 \\ 3 \end{pmatrix} = \begin{bmatrix} 2x_1 \\ 4+16x_2 \end{bmatrix} \Big| \begin{array}{l} x_1=1 \\ x_2=3 \end{array} = \begin{bmatrix} 2 \\ 52 \end{bmatrix}$$

$$\begin{aligned} \text{Parr. eq. } d(x_1, x_2) &= (x_1 - 40)^2 + (x_2 - 40)^2 \\ \nabla d \begin{pmatrix} x_1, x_2 \end{pmatrix} &= \begin{bmatrix} 2(x_1 - 40) \\ 2(x_2 - 40) \end{bmatrix} = \begin{bmatrix} -70 \\ -76 \end{bmatrix} \end{aligned}$$

$$-\nabla f \begin{pmatrix} 5 \\ 2 \end{pmatrix} = \begin{bmatrix} 70 \\ 76 \end{bmatrix}$$

- Gradient descent : $\vec{x}_{t+1} = \vec{x}_t + \gamma(-\nabla f(\vec{x}_t))$

General form of "steepest descent"

L 8.8 Taylor series in higher dimensions

- Higher order taylor series,

$$f(\vec{x} + \gamma d) = f(\vec{x}) + \gamma d^T \nabla f(\vec{x}) + \dots$$

$$f(\vec{x} + \gamma d) - f(\vec{x}) \approx \gamma d^T \nabla f(\vec{x})$$

$$\Rightarrow d = -\nabla f(\vec{x})$$

- $-\nabla f(\vec{x})$ gives the steepest descent

- G.D is also called the steepest descent.

$$\min_{\vec{x} \in \mathbb{R}^d} f(\vec{x})$$

$\boxed{g(\vec{x}) \leq 0}$

Week 9

19.1 Constrained Optimization - I

$$\min_{\vec{x}} f(\vec{x}) \text{ st. } g(\vec{x}) \leq 0$$

\vec{x}^* solves this, how to check whether it is indeed optimal

$$(1) \quad g(\vec{x}^*) \leq 0$$

(2) no "descent direc" should be a "feasible direc" \rightarrow any direc that induces the

any direc that takes to func value point that is feasible (for same step-size) satisfying constraint $g(\vec{x}) \leq 0$

- no descent direc is feasible direc

- $d^T \nabla f(\vec{x}) \leq 0 \Rightarrow d$ is a descent direc

- if \vec{x}^* satisfies $g(\vec{x}^*) \leq 0$, then any descent direc for g is also a "feasible direc"

- \vec{x}^* can't be optimal for the $\nabla f(\vec{x}^*)$, $\nabla g(\vec{x}^*)$ deviates above

19.2

Part - 2

Necessary condition for optimality of \vec{x}^*

$$\nabla f(\vec{x}^*) = -\lambda \nabla g(\vec{x}^*)$$

↳ Lagrange multiplier

can be

any arbitrary scalar

L9.3 Method of Lagrange Multiplier, Projected Gradient descent

- method of Lagrange multipliers

$$\textcircled{1} \quad g(x^*) = 0$$

$$\textcircled{2} \quad \nabla f(x^*) = -\lambda \nabla g(x^*) \quad \text{for some } \lambda$$

$$\text{Eq. } f(x_1, x_2) = x_1^2 + 2x_2 + 4x_2^2$$

$$g(x_1, x_2) = x_1^2 + x_2^2 - 1$$

$$\nabla f([x_1 \atop x_2]) = \begin{bmatrix} 2x_1 \\ 2+8x_2 \end{bmatrix}; \quad \nabla g([x_1 \atop x_2]) = \begin{bmatrix} 2x_1 \\ 2x_2 \end{bmatrix}$$

$$\begin{bmatrix} 2x_1 \\ 2+8x_2 \end{bmatrix} = -\lambda \begin{bmatrix} 2x_1 \\ 2x_2 \end{bmatrix}$$

$$2x_1 = -\lambda 2x_1 \quad \textcircled{1}$$

$$2+8x_2 = -\lambda (2x_2) \quad \textcircled{2}$$

from \textcircled{1} \Rightarrow either $x_1 = 0$ or $\lambda = -1$

Case 1: $\lambda = -1$

$$2+8x_2 = -(-1)2x_2 \Rightarrow x_2 = -\frac{1}{3}$$

$$\text{then } x_1 = \left\{ \frac{+\sqrt{8}}{3}, \frac{-\sqrt{8}}{3} \right\}$$

Case 2: $x_1 = 0$

$$\text{then } x_2 = 1 \text{ or } -1$$

To find $\min f(x)$ substitute each pt.
 $g(x) = 0$ into

$$f(x) = x_1^2 + 2x_2 + 4x_2^2 \quad \left| \begin{array}{l} f[0] = 6 \\ f[-1] = 2 \end{array} \right. \quad \left| \begin{array}{l} f[\frac{1}{3}] = \frac{19}{3} \\ f[-\frac{1}{3}] = \frac{2}{3} \end{array} \right.$$

- It may not be true always possible to solve the sys. of eq^{ns} that satisfy the Lagrange eqⁿ.

- Project Gradient descent

$$x_0 \quad \text{for } t=1, \dots T$$

$$x_{t+1} = x_t - \gamma \nabla f(x_t)$$

Projection operator

$$\Pi(x) = \min_{y \in S} \|x-y\|^2$$

$$\{y : g(y) \leq 0\}$$

$$x_{t+1} = \Pi(x_t - \gamma \nabla f(x_t))$$

gradient

projection start

another x ends

L9.4 Intuition to Convexity

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \rightarrow \begin{bmatrix} a & b \\ \lambda a + (1-\lambda)b & \lambda c + (1-\lambda)d \end{bmatrix}$$

for more genrl case

- A set $S \subseteq \mathbb{R}^d$ is a convex set if $\forall x_1, x_2 \in S$, then $\lambda x_1 + (1-\lambda)x_2 \in S$ $\forall \lambda \in [0, 1]$

$$\begin{cases} x_1 \\ x_2 \end{cases} \rightarrow \begin{cases} \lambda x_1 + (1-\lambda)x_2 \\ \lambda \geq 0 \\ \lambda \leq 1 \end{cases}$$

- Half space hyperplanes with convex sets

$$S \subseteq \mathbb{R}^d \quad \left\{ x : w^T x \geq b \right\}$$

- Half-spaces are convex

$$S \subseteq \mathbb{R}^d \text{ if } \{x \in \mathbb{R}^d : w^T x \leq b\}$$

$$\|p - x\|_{\infty} = \|x\|_1$$

L 9.5 Prop. of Convex Sets

- Intersection of convex sets

$$S_{1,2} = S_1 \cap S_2 = \{x : x \in S_1, x \in S_2\}$$

$\Rightarrow S_{1,2}$ is also convex

$$\text{Ex. } \{x : Ax = b\} \quad A \in \mathbb{R}^{m \times d}, b \in \mathbb{R}^m$$

is it convex?

S is "convex [intersection]" of hyperplanes is also a convex

- Convex combination

Let $S = \{x_1, \dots, x_n\} \subseteq \mathbb{R}^d$. $z \in \mathbb{R}^d$ to be a convex combination of pts. in S if $\exists \lambda_1, \dots, \lambda_n$

$$[z = \lambda_1 x_1 + \lambda_2 x_2 + \dots + \lambda_n x_n]$$

and they sum up

$$[\sum \lambda_i = 1]$$

$$-\text{Convex Hull}(S) = \left\{ z : z = \sum_{i=1}^n \lambda_i x_i \text{ for some } \lambda_i \geq 0, \sum \lambda_i = 1 \right\}$$

$$\{x_1, \dots, x_n\} \rightarrow \left\{ \sum \lambda_i x_i \mid \lambda_i \geq 0, \sum \lambda_i = 1 \right\}$$

this in itself is a convex set

- Convex hull $\{x_1, \dots, x_n\}$ as the intersection of all convex sets that contain $\{x_1, \dots, x_n\}$

- Euclidean balls in \mathbb{R}^d

$$B = \{x : \|x\|_2 \leq r\} \rightarrow \text{This is also a convex set}$$

$$\text{is just } \sqrt{\sum_{i=1}^d x_i^2} \geq (x)$$

L 9.6 Convex Functions

$$f: \mathbb{R}^d \rightarrow \mathbb{R}$$

\hookrightarrow any convex set

$$\text{epi}(f) = \left\{ \begin{bmatrix} x \\ z \end{bmatrix} \in \mathbb{R}^{d+1} : z \geq f(x) \right\}$$

- A funcⁿ $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex funcⁿ if $\text{epi}(f)$ is a convex set.

- or,
A funcⁿ $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex iff $\forall x_1, x_2 \in \mathbb{R}^d$ and all $\lambda \in [0, 1]$

$$f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2)$$

- or,
Assume, $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is differentiable

$$f \text{ is convex iff } f(y) \geq f(x) + (y-x)^T \nabla f(x)$$

- or,
 f is twice differentiable
 $H \in \mathbb{R}^{d \times d}$

$$H_{ij} = \frac{\partial^2}{\partial x_i \partial x_j} f(x_1, \dots, x_d) \rightarrow \mathbb{R}$$

f is convex iff. H is p.s.d matrix
the semi-def. eigenvalues

E9.7 Prop. of Convex funcⁿ (with examples)

- If f is a convex funcⁿ, then all local minima of f are also global minimum

$f(z) < f(x^*) \rightarrow$ this is convex

local min
global min

Week - 0

(10.1) Prop. of Convex Funcⁿ

- it continues to last pt., there could be multiple global minimum.
- set of all global minima of a convex funcⁿ is a convex set.
- necessary & sufficient condⁿ for optimality of convex funcⁿ

$$\min_{\mathbf{x}} f(\mathbf{x})$$

convex,
differentiable

- If $f \rightarrow$ convex, differential funcⁿ $\mathbb{R}^d \rightarrow \mathbb{R}$
 x^* G.R. is a global minimum of f iff. $\nabla f(x^*) = 0$

and also, the converse is true.

- Addrd prop. of convex funcⁿ -

- ① If $f: \mathbb{R} \rightarrow \mathbb{R}$ & $g: \mathbb{R}^d \rightarrow \mathbb{R}$ are both convex
then, $h(x) = f(x) + g(x)$ is convex
⇒ sum of convex funcⁿ is convex

② Compositions

- $f: \mathbb{R} \rightarrow \mathbb{R}$ be convex & non-decreasing func
 $g: \mathbb{R}^d \rightarrow \mathbb{R}$ be convex
- $$h(x) = f(g(x)) \Rightarrow h := f \circ g$$
- h is also convex.

⇒ Composition of convex with (convex + nondec.) is convex

- ③ $f: \mathbb{R} \rightarrow \mathbb{R}$ convex $g: \mathbb{R}^d \rightarrow \mathbb{R}$ convex
 $h = f \circ g \Rightarrow h$ is convex

- In general, if f & g are convex, then $h = f \circ g$ may not be convex

L10.2 App. of Optimizaⁿ in ML

Linear regression: Minimizing data

Linear regression: Minimizing $\{x_1, \dots, x_n\}$

Linear regression: $y_i \in \mathbb{R}$

$x_i \in \mathbb{R}^d$ + linear model for y_i

Goal: $h: \mathbb{R}^d \rightarrow \mathbb{R}$ such that $\hat{y}_{\text{test}} = h(x_{\text{test}})$

(a) min data \rightarrow reg model $w \rightarrow \hat{y}_{\text{test}}$

h is linear, $h(x) = w^T x$ for some $w \in \mathbb{R}^d$

Performance: sum of squared errors

$$\min_{w \in \mathbb{R}^d} \frac{1}{2} \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$f(w) = \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$h_i(w) = (w^T x_i - y_i)^2$$

$$g(w) = w^T x - y$$

linear \Rightarrow convex

\Rightarrow hence, $f(w)$ is a convex funcⁿ

$$f(w) = \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$f(w) = \frac{1}{2} \|xw - y\|_2^2$$

$$f(w) = \frac{1}{2} (xw - y)^T (xw - y)$$

$$\nabla f(w) = (x^T x) w - x^T y$$

Advantages + Analytical solution

Issue: needs an "inverse" computaⁿ - $O(d^2)$

can use iterative procedures

$$w^{t+1} = w^t - \eta_t \nabla f(w^t)$$

$$(x^T x) w - x^T y \quad \text{doesn't have an inverse}$$

→ Approximation of gradient

by Stochastic Gradient Descent

samples a small set of data

points uniformly at random

new dataset & compute gradient

$$\frac{1}{T} \sum_{t=1}^T w_t \rightarrow w^*$$

L10.3 Revisiting Constrained Optimization

- unconstrained, f is convex, $\nabla f(x^*) = 0 \Rightarrow x^*$ is global opt.
- constrained

$$\begin{aligned} \min_w & f(w) \\ \text{s.t.} & h(w) \leq 0 \end{aligned}$$

- Lagrangian funcⁿ

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

Fix $x \in \mathbb{R}^n$

$$\max_{\lambda \geq 0} L(x, \lambda)$$

$$= \max_{\lambda \geq 0} f(x) + \lambda h(x) \begin{cases} f(x) & h(x) \leq 0 \\ \infty & h(x) > 0 \end{cases}$$

$$\begin{aligned} \min_x & f(x) \\ \text{s.t.} & h(x) \leq 0 \end{aligned} = \min_x \left[\max_{\lambda \geq 0} L(x, \lambda) \right] \quad \hookrightarrow \text{primal problem}$$

$$\max_{\lambda \geq 0} \left[\min_x L(x, y) \right] = \max_{\lambda \geq 0} \left[\min_x f(x) + \lambda h(x) \right]$$

\Rightarrow

$$= \max_{\lambda \geq 0} g(\lambda)$$

\downarrow concave funcⁿ

dual problem

L10.4. Relax^b btree. Periodic & dual push.; KKT cond.

$$\min_{\mathbf{x}} \max_{\lambda \geq 0} L(\mathbf{x}, \lambda) \quad \text{Dual}$$

$$T(x) = \begin{cases} f(x) & \text{if } x \leq 0 \\ \infty & \text{otherwise} \end{cases}$$

For any $\lambda \geq 0$

$$L(x, t_2) = f(x) + \lambda h(x)$$

Ein A30

$$L(x, \lambda) \leq J(x)$$

$$\max_{\lambda \geq 0} \min_{x \in \mathbb{R}^n} L(x, \lambda) \leq f(x^*)$$

$$\max_{\lambda \geq 0} [g(\lambda)] = g(\lambda^*) \leq f(x^*)$$

value at dual opt. \leq value at primal opt.
Weak - duality

If f & h are convex, then Strong-Duality holds

$$\begin{aligned} f(x^*) &= g(\lambda^*) \\ &= \min_{\lambda} \underbrace{f(x) + \lambda^* h(x)}_{\text{Convex combination}} \\ &\leq f(x^*) + \lambda^* h(x^*) \end{aligned}$$

- f, h are convex \Rightarrow strong duality

$x^{\frac{1}{n}}, y$ must satisfy

$$(a) \nabla f(x^*) + \lambda^* \nabla h(x^*) = 0$$

$$(b) \lambda^{\alpha} \ln(x^{\alpha}) = 0$$

$$(c) \quad h(x^*) \leq 0$$

$$(d) \quad \lambda^* \geq 0$$

If x^*, y^* satisfies these conditions \Rightarrow Local Optima

- KKT (Karush - Kuhn - Tucker) Conditions

$$\min f(x)$$

$$h_i(x) \leq 0 \quad \forall i = 1 \dots m$$

$$l_j(x) = 0 \quad \forall j = 1 \dots n$$

$$L(x, u, v) = f(x) + \sum_{i=1}^m u_i h_i(x) + \sum_{j=1}^n v_j l_j(x)$$

L 10.5 KKT cond \Rightarrow contd.

SVM (support vector machine)
Optimization problem

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 \rightarrow \text{quad.} \rightarrow \text{convex}$$

$$\text{s.t. } \mathbf{w}^T \mathbf{x}_i + b \geq 1 \quad \forall i$$

↳ convex

\Rightarrow Convex

Week-11

Contd. RV

- Sample Space $\rightarrow (\Omega, \mathcal{F}, \mathbb{P})$
- $\mathcal{F} \subseteq \{\emptyset, \Omega\}^{\Omega}$
 - (i) Ω GF
 - (ii) $A \in \mathcal{F} \Rightarrow A^c \in \mathcal{F}$
 - (iii) $A_1, A_2, \dots, A_n \in \mathcal{F} \Rightarrow \bigcup_{i=1}^n A_i \in \mathcal{F}$
- $X : \Omega \rightarrow \mathbb{R}$

Domain & range to be uncountable

Expt. - Waiting for bus, X : amt. of time you wait
 7:00, 7:15, 7:30, ..., Reach the bus stop 7:10, 7:20

$$f_x(x) = \frac{1}{15} \int_x^{x+15} dx \quad x \in \mathbb{R} \quad \text{PDF}$$

$$P(X=n) = 0$$

$$F_x(x) = P(X \leq x) \quad \text{CDF}$$

- ① $f_x(x) \geq 0$
- ② $\int_{-\infty}^{\infty} f_x(x) dx = 1$
- ③ $F_x(-\infty) = 0$
- ④ $F_x(\infty) = 1$
- ⑤ F_x is increasing

Ex.



$$f_x(x) = \begin{cases} \frac{1}{15} & \text{if } x \in [0, 15] \\ 0 & \text{otherwise} \end{cases}$$

$$F_x(x) = \begin{cases} 0 & x < 0 \\ x/15 & x \in [0, 15] \\ 1 & x > 15 \end{cases}$$

Ex. $f_x(x) = \begin{cases} x/2 & \text{if } x \in [0, 2] \\ 0 & \text{otherwise} \end{cases}$

$$F_x(x) = P(X \leq x) = \int_0^x f_x(x) dx = \frac{1}{2} \left[\frac{x^2}{2} \right]_0^x = \frac{x^2}{4}$$

Properties of cumulative distribution function:

L11.2 Conditional PDF (Ω, \mathcal{F}, P) $X, A \subseteq \Omega$

$$f_{X|A}(x) = \frac{P(X \in [x, x+dx] | A)}{\int_{x \in A} dx}$$

Ex. $f_x(x) = \begin{cases} x/2 & x \in [0, 2] \\ 0 & \text{others} \end{cases}$

$$f_{X|A}(x) = \begin{cases} 0 & \text{if } x < 1 \\ \frac{P(X \in [x, x+dx])}{P(A) dx} & \text{if } x \in [1, 2] \\ 0 & \text{if } x > 2 \end{cases}$$

$$= \begin{cases} 0 & \text{if } x < 1 \\ \frac{4x}{6} & \text{if } x \in [1, 2] \\ 0 & \text{if } x > 2 \end{cases}$$

- functions of RV:
 $x \rightarrow y = \frac{x}{2}$ or $y = x^2$ or $y = |x|$

Ex. $f_x(x) = \begin{cases} 1/2 & \text{if } x \in [-1, 1] \\ 0 & \text{otherwise} \end{cases}$ $y = \frac{x}{2}$

$$\begin{aligned} f_y(y) &= P(X \in [y, y+dy]) \\ &= P(X \in [2y, 2y+2dy]) \\ &\Rightarrow 0 \cdot y \text{ if } 2y \notin [-1, 1] \\ &\Rightarrow 2 \cdot \frac{1}{2} \text{ if } y \in [-\frac{1}{2}, \frac{1}{2}] \end{aligned}$$

Ex. $y = 1|x|$ X is uniform $[-1, 1]$

$$f_y(y) = \begin{cases} 1 & \text{if } y \in [0, 1] \\ 0 & \text{otherwise} \end{cases}$$

L11.3 Expectation

$$E[X] = \int_{-\infty}^{\infty} x \cdot f_x(x) dx$$

- Properties -

$$\begin{aligned} ① E[X+Y] &= E[X] + E[Y] \\ ② Y = g(x) & \quad E[Y] = \int_{-\infty}^{\infty} g(x) \cdot f_x(x) dx \end{aligned}$$

Ex. $X \sim \text{unit } [a, b]$

$$f_x(x) = \begin{cases} \frac{1}{(b-a)} & x \in [a, b] \\ 0 & \text{otherwise} \end{cases}$$

$$E[X] = \frac{(b+a)}{2}$$

Ex. $f_x(x) = \begin{cases} x/2 & \text{if } x \in [0, 2] \\ 0 & \text{otherwise} \end{cases}$

$$E[X] = \int_0^2 x \cdot \frac{x}{2} dx = \frac{4}{3}$$

$$\begin{aligned} - \text{Var}[X] &= E[(X - E(X))^2] \\ &= E(X^2) - (E(X))^2 \end{aligned}$$

$$\sqrt{\text{Var}[X]} = SD[X]$$

- Prop -

- ① $\text{Var}[X+Y] \neq \text{Var}[X] + \text{Var}[Y]$
- ② $\text{Var}[aX] = a^2 \text{Var}[X]$
- ③ $\text{Var}[X] \geq 0$

Ex. $f_x(x) = \begin{cases} \frac{1}{b-a} & x \in [a, b] \\ 0 & \text{others} \end{cases}$

$$E[X^2] = \frac{1}{3} [b^2 + a^2 + ab]$$

$$(E[X])^2 = \frac{(b+a)^2}{4}$$

$$\boxed{\text{Var}[X] = \frac{(b-a)^2}{12}}$$

- X, A

$$E[X|A] = \int_{-\infty}^{\infty} x \cdot f_{X|A}(x) dx$$

Ex. $f_{X|A}(x) = \begin{cases} 5 & \text{if } x \in [0, 5] \\ 0 & \text{otherwise} \end{cases}$ $E[X|A] = 2.5$

$$E[X|A^c] = 7.5$$

* $E[X] = E[X|A] \cdot P(A) + E[X|A^c] \cdot P(A^c)$

$$E[X] = \frac{50}{8}$$

11.4 Multiple RV

- $f_{XY}(x, y) = P(X \in [x, x+dx], Y \in [y, y+dy])$
 \hookrightarrow Joint Distribution / Density Function

$$F_{XY}(x, y) = P(X \leq x, Y \leq y)$$

Prop.

- ① $f_{XY}(x, y) \geq 0$
- ② $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{XY}(x, y) dx dy = 1$
- ③ $F_{XY}(-\infty, -\infty) = 0$
- ④ $F_{XY}(\infty, \infty) = 1$

- $f_X(x) = \int_{-\infty}^{\infty} f_{XY}(x, y) dy$ Marginal density

$$f_Y(y) = \int_{-\infty}^{\infty} f_{XY}(x, y) dx$$

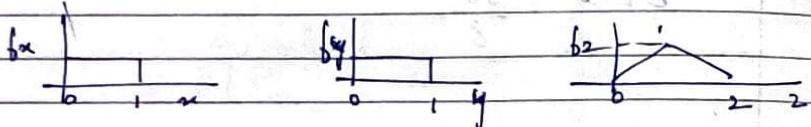
$$f_{X|Y}(x|y) = \frac{f_{XY}(x, y)}{f_Y(y)}$$

conditions
 $\downarrow \neq 0$

L11.5 Indep. RV.

$$X, Y \\ Z = X + Y$$

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(x) f_Y(z-x) dx$$



$$Z = \max(X, Y)$$

$$F_Z(z) = P(Z \leq z) \\ = P(X \leq z) \cdot P(Y \leq z) \\ = F_X(z) \cdot F_Y(z)$$

$$Z = \min(X, Y)$$

$$F_Z(z) = P(Z \leq z) \\ = P(X \leq z \cup Y \leq z)$$

$$X, Y$$

$$\text{Cov}[X, Y] = E[XY] - E[X]E[Y]$$

$$\rho(X, Y) = \frac{\text{Cov}[X, Y]}{\sqrt{\text{Var}[X] \cdot \text{Var}[Y]}}$$

Independent \implies Uncorrelated

$$\text{Ex. } X \sim \text{Unif. } [-1, 1]$$

$$E(XY) = E(X^2) = 0$$

$$E[X] \cdot E[Y] = 0$$

$$Y = X^2$$

$$E[XY] = 0$$

$$E[X] \cdot E[Y] = 0$$

L11.6 Transformed RV

$$W, X \\ Y = g(W, X)$$

$$z = h(W, X)$$

$$f_{YZ} = ?$$

$$f_{YZ}(y, z) = P(Y=y, Z=z) = P(Y=y, h(W, X)=z) = P(Y=y, W=w) = f_Y(y) f_W(w)$$

L11.7 Uniform Exponential, Normal

- $X \sim \text{Unit } (a, b)$.

$$f_X(x) = \begin{cases} \frac{1}{b-a} & x \in [a, b] \\ 0 & \text{otherwise} \end{cases}$$

$$E[X] = \frac{a+b}{2} \quad \text{Var}[X] = \frac{(b-a)^2}{12}$$

- Exponential

$$X \sim \exp(\lambda) \\ f_X(x) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$F_X(x) = 1 - e^{-\lambda x}$$

Memoeyless property

$$P(X \geq a | X \geq b) = P(X \geq a-b)$$

$$E[X] = \frac{1}{\lambda}$$

$$X \sim \exp(\lambda)$$

$$Y \sim \exp(\tau)$$

$$Z = \min(X, Y)$$

$$Z \sim \exp(\lambda + \tau)$$

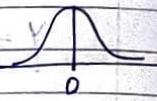
$$F_Z(z) = 1 - e^{-(\lambda + \tau)z}$$

$$f_Z(z) = \begin{cases} 0 & z < 0 \\ (\lambda + \tau) e^{-(\lambda + \tau)z} & z \geq 0 \end{cases}$$

- Normal

$$z \sim N(0, 1)$$

$$f_z(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right)$$



- Gaussian

$$x = \sigma z + \mu \quad z \sim N(0, 1)$$

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

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

$$x \sim N(\mu, \sigma^2)$$

$$E[x] = \mu$$

$$\text{Var}(x) = \sigma^2$$

Week-12

112.1 Bivariate & Multivariate Normal

$$z = \begin{bmatrix} z_1 \\ \vdots \\ z_d \end{bmatrix}$$

$$z_1 \sim N(0, 1) \dots z_d \sim N(0, 1)$$

$$f_z(z) = \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{1}{2} \|z\|^2\right)$$

$$x_1 = z_1; x_2 = \rho z_1 + \sqrt{1-\rho^2} z_2$$

$$x = \begin{bmatrix} 1 & 0 \\ 0 & \sqrt{1-\rho^2} \end{bmatrix} z \Rightarrow z = \begin{bmatrix} 1 & 0 \\ \frac{\rho}{\sqrt{1-\rho^2}} & \frac{1}{\sqrt{1-\rho^2}} \end{bmatrix} y$$

$$\det(A) = \sqrt{1-\rho^2} \quad \det(A^{-1}) = \frac{1}{\sqrt{1-\rho^2}}$$

$$E[x_1] = E[x_2] = 0$$

$$\text{cov}[x_1, x_2] = E[x_1 x_2] = E[\rho z_1^2 + \sqrt{1-\rho^2} z_2^2]$$

$$\text{var}[x_1] = 1 \quad \text{var}[x_2] = \rho^2 + 1 - \rho^2 = 1$$

$$\text{cov}[x] = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} = A A^T$$

$$\det(\Sigma) = 1 - \rho^2$$

$$[\text{cov}[x]]^{-1} = \frac{1}{1-\rho^2} \begin{bmatrix} 1 & -\rho \\ -\rho & 1 \end{bmatrix}$$

$$f_x(x) = f_z(A^{-1}x) \cdot |\det(A^{-1})|$$

$$\begin{aligned} &= \frac{1}{2\pi\sqrt{1-\rho^2}} \exp\left(-\frac{1}{2(1-\rho^2)} [x_1, x_2] \begin{bmatrix} 1 & -\rho \\ -\rho & 1 \end{bmatrix} [x_1, x_2]\right) \\ &= f_{x_1}(x_1) \cdot f_{x_2|x_1}(x_2|x_1) \\ &\sim N(x_1 | 0, 1) \quad \sim N(x_2 | \rho x_1, 1 - \rho^2) \end{aligned}$$

- $x = A^{-1}z$
 $f_x(x) = f_z(A^{-1}x) \cdot |\det(A^{-1})|$

L12.2 Estimation of Parameters using ML

- $P = \{P_\theta : \theta \in \Theta\}$

x_1, \dots, x_n drawn i.i.d from P_θ for $\theta \in \Theta$

- $L(\theta) = P(x_1 = x_1, \dots, x_n | \theta)$

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

$$= \prod_{i=1}^n P_\theta(x_i)$$

- $\log(L(\theta)) = \sum_{i=1}^n \log(P_\theta(x_i))$

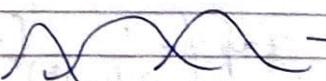
- $P = \{ \text{Bern}(\theta) : \theta \in [0, 1] \}$

$$P_\theta(x) = \begin{cases} 0 & \text{if } x=1 \\ 1-\theta & \text{if } x=0 \end{cases}$$

- $R(\theta) = \sum_{i=1}^n \log(P_\theta(x_i))$

- $\hat{\theta}_{ML} = \frac{a}{n} = \frac{\sum_{i=1}^n x_i}{n}$

L12.3 Gaussian mixture model & expectation maximization

 \rightarrow multimodal

$$f_x(x) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k)$$

L12.4 Law of large nos.

- Markov inequality -

$$P(X \geq t) \leq \frac{\mu}{t} \quad X \text{ is +ve RV} \quad E[X] = \mu$$

- Chebyshhev Inequality

$$E[X] = \mu, \text{Var}[X] = \sigma^2$$

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

- Hoeffding Inequality

$$P(|\bar{X}_n - \mu| \geq t) \leq \frac{\text{Var}[\bar{X}_n]}{t^2} = \frac{\sigma^2}{nt^2}$$

- Convergence in Probability

$$\bar{X}_n \xrightarrow{P} \mu \quad \text{if}$$

$$\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| \geq t) = 0 \quad \forall \epsilon > 0$$

- Law of large nos.

$$X_1, \dots, X_n \quad E[X_i] = \mu$$

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \quad \bar{X}_n \xrightarrow{P} \mu$$