

ML FOUNDATIONS



WEEK 0

Stats 1

$$\textcircled{1} \quad P(A \text{ or } B) = \frac{1}{2} \quad \text{for 5 question} \rightarrow (y_2)^5 = 0.03125$$

$$\textcircled{2} \quad P(A) = 0.3 \quad P(B) = 0.4$$

$$P(A \cup B) = 0.6$$

$$P(A \cap B) = P(A) + P(B) - P(A \cup B) = 0.1$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{0.1}{0.4} = 0.25$$

③	category	prob ^{<50}	freq	tot. prob.
	1	0.6	0.3	0.18
	2	0.5	0.2	0.1
	3	0.2	0.1	0.02
	4	0.8	0.4	0.32 sum = 0.62

$P(4|<50) = \frac{0.32}{0.62} = 0.51613$

$$\textcircled{4} \quad M_1 = 0.3 - 0.02 = 0.2449$$

$$M_2 = 0.45 - 0.03 = 0.5510$$

$$M_3 = 0.25 - 0.02 = 0.2041$$

$$\text{total probability} = 0.006 + 0.0185 + 0.005 = 0.0245$$

$$\textcircled{5} \quad \binom{2}{6} = \frac{6!}{2!4!} = \frac{6 \times 5}{2} = 15$$

$$\textcircled{6} \quad \text{Var}(x+y) = \text{Var}(x) + \text{Var}(y)$$

$$E(x^2) = 15.167 \quad (E(x))^2 = 12.25$$

$$\text{Var}(x) = 2.917 + 0.25$$

$$E(y^2) = 0.5 \quad (E(y))^2 = 0.25$$

$$= 3.167$$

$$\textcircled{7} \quad 60, 40, 40, 60$$

$$\frac{1}{3} \rightarrow 60 \quad E(x^2) = 2266.67, \quad (E(x))^2 = 2177.78$$

$$\frac{2}{3} \rightarrow 40 \quad \text{Var}(x) = 88.89, \quad \text{SD}(x) = 9.428$$

M	T	W	Th	F	Sa	Su
-	-	-	-	-	-	$\frac{3}{2}$

$$\textcircled{8} \quad \frac{3}{4} \int_0^{\infty} (0.25)^k dk \Rightarrow \lim_{n \rightarrow \infty} \int_0^n (0.25)^k dk \Rightarrow \lim_{n \rightarrow \infty} \left[\frac{(0.25)^k}{\ln(0.25)} \right]_0^\infty \Rightarrow \left(\lim_{n \rightarrow \infty} \frac{(0.25)^n}{\ln(0.25)} - \frac{1}{\ln(0.25)} \right) \cdot \frac{3}{4}$$

$$= 0 - \frac{1}{\ln(0.25)} \Rightarrow \frac{-3}{4(\ln(0.25))}$$

Stats 2

$$\textcircled{1} \quad \text{Var}(x) = E(x^2) - (E(x))^2$$

$$E(x^2) = \int_0^1 4x^5 dx \Rightarrow \frac{4}{6} [x^6]_0^1 = \frac{2}{3}$$

$$E(x) = \int_0^1 4x^4 dx \Rightarrow \frac{4}{5} [x^5]_0^1 = \frac{4}{5}$$

$$\text{Var}(x) = \frac{2}{3} - \frac{16}{25} = \frac{2}{75}$$

$$\textcircled{2} \quad k \cdot \int_{y=0}^{\infty} \int_{x=0}^{\infty} e^{-x} \cdot e^{-y} dx dy = 1 \Rightarrow k \cdot \int_0^{\infty} e^{-y} \left[-e^{-x} \right]_0^{\infty} dy = 1 \Rightarrow k \cdot \int_0^{\infty} e^{-y} dy = 1 \Rightarrow k = 1$$

$$\int_{y=0}^{\infty} \int_{x=0}^{\infty} e^{-x} \cdot e^{-y} dx dy \Rightarrow \int_0^{\infty} e^{-y} \left[-e^{-x} \right]_0^{\infty} dy = e^{-y} \int_0^{\infty} e^{-y} dy \Rightarrow e^{-y} (-e^{-y} + 1)$$

$$\textcircled{3} \quad X \sim \text{Bin}(1000, 1/6) \quad E[X_{1000}] = 1/6 \quad \text{Var}(X_{1000}) = \frac{1}{1000} \cdot 5/6 = \frac{5}{36000}$$

chebyshov's = $P(|X - \mu| \geq k) \leq \frac{\text{Var}(X)}{k^2} \quad k = 0.2$

$$P(|X_{1000} - 1/6| \geq 0.2) \leq \frac{5}{36000} \cdot 25 = \boxed{\frac{5}{1440}}$$

$$\textcircled{4} \quad L(u_1, u_2, \dots, u_n) = \left(\frac{1}{\sigma \sqrt{2\pi}} \right)^n \cdot \exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (u_i - \mu)^2 \right) \quad \hat{\mu} = 0$$

$$\log(L) = -n \log(\sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^n (u_i - \mu)^2 \Rightarrow \frac{\partial \log(L)}{\partial \sigma} = -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (u_i - \hat{\mu})^2$$

$$0 = -n + \frac{1}{\sigma^2} \sum_{i=1}^n (u_i - \hat{\mu})^2 \Rightarrow \hat{\sigma}^2 = \frac{\sum_{i=1}^n (u_i - \hat{\mu})^2}{n} = \boxed{2/3}$$

\textcircled{5} Likelihood $\sim N(\bar{y}, 25/10)$ prior $\sim N(50, 25)$

$$\begin{aligned} \text{Post. mean} &= \frac{1}{\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}} \left(\frac{\mu_0}{\sigma_0^2} + \frac{\sum u_i}{\sigma^2} \right) \quad \mu_0 = 50 \quad \sigma_0^2 = 25 \quad \sigma = 25 \\ &= \frac{1}{\frac{1}{25} + \frac{10}{25}} \left(\frac{50}{25} + \frac{63.636}{25} \right) = \frac{25}{11} \cdot \frac{700}{25} = 63.636 \end{aligned}$$

Maths 2

$$\textcircled{1} \quad \lim_{u \rightarrow 0^+} f(u) = 0 \quad \lim_{u \rightarrow 0^-} f(u) = 0 \quad f(0) = 0 \quad \checkmark$$

$$\lim_{h \rightarrow 0^-} \frac{f(0+h) - f(0)}{h} = \frac{f(h)}{h} = 1 \quad ; \quad \lim_{h \rightarrow 0^+} \frac{f(0+h) - f(0)}{h} = \lim_{h \rightarrow 0^+} \frac{h^2}{h} = 0 \quad \times$$

$$\lim_{u \rightarrow 1^+} f(u) = \lim_{u \rightarrow 1^+} u^3 = 1 \quad \lim_{u \rightarrow 1^-} f(u) = \lim_{u \rightarrow 1^-} u^2 = 1 \quad f(1) = 1 \quad \checkmark$$

$$\lim_{h \rightarrow 0^-} \frac{f(1+h) - f(1)}{h} = \frac{(1+h)^2 - 1}{h} = \frac{1+2h+h^2-1}{h} = \frac{h(2+h)}{h} = \lim_{h \rightarrow 0} 2+h = 2$$

$$\lim_{h \rightarrow 0^+} \frac{f(1+h) - f(1)}{h} = \frac{(1+h)^3 - 1}{h} = \frac{1+3h+3h^2+h^3-1}{h} = \frac{h(3+3h+h^2)}{h} = \lim_{h \rightarrow 0} 3+3h+h^2 = 3 \quad \times$$

$$\textcircled{2} \quad f'(u) = \begin{cases} -2u+2 & 0 \leq u \leq 50 \\ 3u^2 & -50 \leq u < 0 \end{cases}$$

$$0 = -2u+2 \Rightarrow u=1 \quad f(1) = 4 \quad \text{glob max} \quad f(50) = -2397 \quad \text{loc min}$$

$$0 = 3u^2 \Rightarrow u=0 \quad f(-50) \approx -(50)^3 \quad \text{glob min}$$

$$\textcircled{3} \quad f'(u) = 3u^2 - b$$

$$0 = 3u^2 - b \Rightarrow 2 = u^2 \Rightarrow u = \sqrt{2}, -\sqrt{2}$$

$$f''(u) = 6u \quad ; \quad f''(\sqrt{2}) = 6\sqrt{2} \quad ; \quad f''(-\sqrt{2}) = -6\sqrt{2}$$

$$\textcircled{4} \quad A^* = \begin{bmatrix} 2 & 2 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 2 & 2 \\ 0 & 2 \end{bmatrix} = \begin{bmatrix} 4 & 8 \\ 0 & 4 \end{bmatrix} \begin{bmatrix} 2 & 2 \\ 0 & 2 \end{bmatrix} = \begin{bmatrix} 8 & 24 \\ 0 & 8 \end{bmatrix}$$

$$AB = \begin{bmatrix} a_1 & a_2 \\ a_3 & a_4 \end{bmatrix} \begin{bmatrix} b_1 & b_2 \\ b_3 & b_4 \end{bmatrix} = \begin{bmatrix} a_1 b_1 + a_2 b_2 & a_1 b_2 + a_2 b_4 \\ a_3 b_1 + a_4 b_3 & a_3 b_2 + a_4 b_4 \end{bmatrix}$$

$$A = \begin{bmatrix} a & b \\ c & d \\ e & f \end{bmatrix} \quad B = \begin{bmatrix} g & h & i \\ j & k & l \end{bmatrix}$$

$$\begin{aligned} (ab+gi) & \quad (ab+hk) \quad (ab+il) \\ (cd+gj) & \quad (cd+hk) \quad (cd+il) \\ (ef+gj) & \quad (ef+hk) \quad (ef+il) \end{aligned}$$

$$\textcircled{5} \quad A = \left[\begin{array}{cccc|c} 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 \end{array} \right]$$

$$\textcircled{6} \quad A = \left[\begin{array}{ccc|c} -1 & 1 & 2 & 1 \\ 2 & 1 & -2 & -1 \\ 0 & 3 & c & d \end{array} \right] \xrightarrow{-R_1} \left[\begin{array}{ccc|c} 1 & -1 & -2 & -1 \\ 0 & 3 & 2 & 1 \\ 0 & 3 & c & d \end{array} \right] \xrightarrow{R_2/3} \left[\begin{array}{ccc|c} 1 & -1 & -2 & -1 \\ 0 & 1 & 2/3 & 1/3 \\ 0 & 3 & c & d \end{array} \right] \xrightarrow{R_3-3R_2} \left[\begin{array}{ccc|c} 1 & 0 & -4/3 & -2/3 \\ 0 & 1 & 2/3 & 1/3 \\ 0 & 0 & c-2 & d-1 \end{array} \right]$$

$$\textcircled{7} \quad \left[\begin{array}{ccc|c} c & 1 & -1 \\ -1 & 0 & -3 \\ -2 & -1 & c \end{array} \right] \quad \det(A) = 3c + (-c+6) - (1) \\ 0 = 2c + 5 \Rightarrow c = 5/2 = 2.5$$

$$\textcircled{8} \quad A = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \quad A^T = \begin{bmatrix} a & c \\ b & d \end{bmatrix} \quad \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

WEEK 1

→ Why and When ML?

- ↳ Programming / Human labour fails?
 - ① scale / speed / cost of human labour
 - ② inability to express rules using language
 - ③ don't know exact rules
- ↳ ML can succeed when
 - ① Have lots of example data
 - ② Have some structural idea on the rules

→ Machine learning applications

- ↳ e-mail inbox spam classification
- ↳ recommender systems (e-commerce, social media, streaming services, etc.)
- ↳ smart assistants
- ↳ robot AIs
- ↳ games (chess, go)
- ↳ social media marketing

→ Types of Models in ML

- ↳ Predictive Model
 - ↳ Regression ; e.g. → model the price of a house based on its area
 - ↳ real-valued output
 - ↳ Classification ; e.g. → whether a house is closer than 2kms to metro
 - ↳ discrete valued output

↳ Probabilistic Model

- ↳ goal is not to predict future
- ↳ to evaluate how likely an event is

→ Learning Algorithms

- ↳ learning algorithms : Data → Models
 - ↳ choose from a collection of models, with same structure but different parameters

Notation

- \mathbb{R} : real numbers, \mathbb{R}_+ : Positive reals, \mathbb{R}^d : d-dimensional vector of reals.
- \mathbf{x} : vector. x_j : j^{th} co-ordinate. $\|\mathbf{x}\|$: Length of vector \mathbf{x} .
- $\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^n$: Collection of n vectors.
- x_j^i : j^{th} co-ordinate of i^{th} vector.
- $(x_1)^2$: Square of the first co-ordinate of the vector \mathbf{x}
- $\underbrace{1(2 \text{ is even})}_{\text{true is } 1} = 1, \underbrace{1(2 \text{ is odd})}_{\text{false is } 0} = 0$.

→ Supervised Learning

↳ curve-fitting at its core

↳ Given $\{(x^1, y^1), (x^2, y^2), \dots, (x^n, y^n)\}$, find a model f such that $f(x^i)$ is 'close' to y^i

↳ Regression

→ e.g. Predict house price from room, area, distance

→ takes training data: $\{(x^1, y^1), (x^2, y^2), \dots, (x^n, y^n)\}$

→ algorithm outputs a model $f: \mathbb{R}^d \rightarrow \mathbb{R}$

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

$$\rightarrow f(x) = w^T x + b = \sum_{j=1}^d w_j x_j + b$$

↳ Classification

→ e.g. Predict if room > 3 from area and price

→ training data

$$x^i \in \mathbb{R}^d, y^i \in \{+1, -1\}$$

→ Algorithm outputs a model $f: \mathbb{R}^d \rightarrow \{+1, -1\}$

$$\rightarrow \text{Loss} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(f(x^i) \neq y^i) \rightarrow \text{fraction of misclassified instances}$$

$$\rightarrow f(x) = \text{sign}(w^T x + b) \quad \text{a.k.a. linear separator}$$

↳ uses **training data** to get model f

↳ " **test data** not in the training data for model evaluation

↳ " **validation data** for model selection

→ Unsupervised Learning

↳ much more vague than supervised learning

↳ usually used to 'understand data'

↳ build models that compress, explain and group data

↳ Dimensionality Reduction

→ compression and simplification

$$\rightarrow \text{Data: } \{x^1, x^2, \dots, x^n\}$$

$$x^i \in \mathbb{R}^d$$

$$\text{Encoder } f: \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$$

$$\text{Decoder } g: \mathbb{R}^{d'} \rightarrow \mathbb{R}^d$$

$$\text{Goal: } g(f(x^i)) \approx x^i$$

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

↳ Density Estimation

→ probabilistic model

$$\rightarrow \text{Data: } \{x^1, x^2, \dots, x^n\}$$

$$x^i \in \mathbb{R}^d$$

Probability mapping $P: \mathbb{R}^d \rightarrow \mathbb{R}_+$ that 'sums' to one

Goal: $P(x)$ is large if $x \in \text{Data}$, and low otherwise

$$\text{Loss} = \frac{1}{n} \sum_{i=1}^n -\log(P(x^i)) \quad \text{a.k.a. negative log likelihood}$$

WEEK 2

→ Sets

- ↪ \mathbb{R} → real numbers
- ↪ \mathbb{R}_+ → the real numbers
- ↪ \mathbb{Z} → integers
- ↪ \mathbb{Z}_+ → the integers

↪ $[a, b] = \{x \in \mathbb{R} : a \leq x \leq b\}$
 $(a, b) = \{x \in \mathbb{R} : a < x < b\}$

- ↪ \mathbb{R}^d → set of d -dimensional vectors = $\mathbb{R} \times \mathbb{R} \times \dots \times \mathbb{R}$ $\leftarrow d$ times
- ↪ $[a, b]^d \rightarrow \{x \in \mathbb{R}^d : x_i \in [a, b], i \in \{1, 2, \dots, d\}\}$

→ Metric Space

- ↪ A set with a distance function associated with it
- ↪ Default metric space = \mathbb{R}^d : $D(x, y) = \|x - y\| = \sqrt{(x_1 - y_1)^2 + \dots + (x_d - y_d)^2}$
- ↪ An open ball with radius ϵ with center x
 $\rightarrow B(x, \epsilon) = \{y \in \mathbb{R}^d : D(x, y) < \epsilon\}$
- ↪ closed ball $\rightarrow \bar{B}(x, \epsilon) = \{y \in \mathbb{R}^d : D(x, y) \leq \epsilon\}$
 center radius

→ Logic Statements/Modifiers

- ↪ \forall : for all
- ↪ \Rightarrow : implies
- ↪ \exists : there exists
- ↪ \Leftrightarrow : equivalent

→ Sequence

- ↪ ordered collection of elements
- ↪ example: $x_n = \left(1 + \frac{4}{2^n}, 3 - \frac{4}{2^n}\right)$ → each element is in \mathbb{R}^2 and as n increases, it will approach $(1, 3)$
- ↪ x_1, x_2, \dots where $x_i \in \mathbb{R}^d$
- ↪ $\lim_{i \rightarrow \infty} x_i = x^* \Leftrightarrow \forall \epsilon > 0, \exists N \text{ s.t. } x_n \in B(x^*, \epsilon) \forall n \geq N$
- ↪ interpretation: after a point N , all the terms of sequence stay within a ball. whatever ϵ , there exists some integer N , such that after N , x_n is going to stay within a ball of radius ϵ centered x^* .

→ Vector Space

- ↪ collection of vectors that must satisfy certain properties
- ↪ most crucial property: if V is a vector space, $u \in V, v \in V, \alpha, \beta \in \mathbb{R}$
 then $\alpha u + \beta v \in V$
- ↪ $u, v \in V ; u \cdot v = u^T v = \sum_{i=1}^d u_i v_i \rightarrow$ (dot product)
- ↪ $\|u\|^2 = u \cdot u = u^T u = \sum_{i=1}^d (u_i)^2 \rightarrow$ (norm)
- ↪ u & v are perpendicular (orthogonal) if $u \cdot v = u^T v = \sum_{i=1}^d u_i v_i = 0$

→ Functions and Graphs

↪ graph of a function in n -dimensional space is $n+1$ dimensions

example of 1-dimensional function $f: \mathbb{R} \rightarrow \mathbb{R}$ $f(u) = u^2$ $G_f = \{(u, u^2) : u \in \mathbb{R}\} \subseteq \mathbb{R}^2$

↪ contour plots of 2-dimensional function

e.g. $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ $f(u) = u_1 + u_2$ values it can take = $\{-1, 0, 1\}$ (maybe more)

$f(u) = -1$, plot all possible u_1, u_2 st. $u_1 + u_2 = -1$

$f(u) = 0$, " " " " " " $u_1 + u_2 = 0$

→ Continuity

↪ $f: \mathbb{R} \rightarrow \mathbb{R}$ is continuous at $u^* \in \mathbb{R}$ if & sequences converging to u^* , $f(u_i)$ converges to $f(u^*)$

$$\lim_{i \rightarrow \infty} u_i = u^* \Rightarrow \lim_{i \rightarrow \infty} f(u_i) = f(u^*)$$

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

→ Differentiability

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

$$\lim_{h \rightarrow 0} \frac{f(u+h) - f(u)}{h}$$

→ Derivatives and Linear Approximation

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

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

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

$$f(u) \approx \underbrace{f(u^*) + f'(u^*)(u - u^*)}_{L_{u^*}[f](u)} \quad \text{—— around } u = u^*$$

→ Higher Order Approximation

$$f(u) \approx f(u^*) + f'(u^*)(u - u^*) + \frac{1}{2} f''(u^*)(u - u^*)^2 \quad \text{—— quadratic approximation}$$

ex. give lin. approx. of $f(u) = e^{\sqrt{1+u}}$ around $u = 1$

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

$$f(u) = f(1) + f'(1)(u - 1)$$

$$= \boxed{e^{\sqrt{2}} + \frac{e^{\sqrt{2}}}{2\sqrt{2}}(u - 1)}$$

$$f'(u) = e^{\sqrt{1+u}} \cdot \frac{1}{2\sqrt{1+u}}$$

$$f'(1) = \frac{e^{\sqrt{2}}}{2\sqrt{2}}$$

→ Geometry of Lines

↪ A line in $\mathbb{R}^d \subseteq \mathbb{R}^d$

↪ A line through the point $u \in \mathbb{R}^d$ along the vector $v \in \mathbb{R}^d$

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

A line through the points $u, u' \in \mathbb{R}^d = \{u \in \mathbb{R}^d : u = u + \alpha(u' - u) \text{ for } \alpha \in \mathbb{R}\}$

→ Geometry of (Hyper)planes

↪ A $(d-1)$ -dimensional hyperplane $\in \mathbb{R}^d$

↪ A hyperplane normal to the vector $w \in \mathbb{R}^d$ with $b \in \mathbb{R}$

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

$$\{u \in \mathbb{R}^d : \sum_{i=1}^d w_i u_i = b\}$$

→ Gradient and Linear Approximations

↪ $v, u \in \mathbb{R}^d$ and u is close to v

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

$$= f(v) + \sum_{i=1}^d \frac{\partial f(v)}{\partial u_i} \cdot (u_i - v_i)$$

$$\hookrightarrow f: \mathbb{R}^2 \rightarrow \mathbb{R} \quad f(y_1, y_2) - f(v_1, v_2) \approx \frac{\partial f}{\partial u_i}(v)(y_i - v_i)$$

↪ e.g. → approx of $f(u_1, u_2) = u_1^2 + u_2^2$ around $(3, 1)$?

$$f(u) = f(3, 1) + [6, 2] \begin{bmatrix} u_1 - 3 \\ u_2 - 1 \end{bmatrix}$$

$$\nabla f(u) = (2u_1, 2u_2)$$

$$\nabla f(v) = (6, 2)$$

$$= 10 + 6u_1 - 18 + 2u_2 - 2$$

$$f(v) = 10$$

$$= 6u_1 + 2u_2 - 10$$

→ Directional Derivative

↪ $D_u[f](v) \rightarrow$ directional derivative of f at the point v , along u

$$D_u[f](v) = \nabla f(v)^T(u)$$

→ Cauchy-Schwarz Inequality

↪ two d -dimensional vectors a & b

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

$$\begin{array}{c} \downarrow \\ a = \alpha b \\ \alpha < 0 \end{array}$$

$$\begin{array}{c} \downarrow \\ a = \alpha b \\ \alpha > 0 \end{array} \rightarrow \text{to maximize } a^T b$$

→ Direction of Steepest Ascent

↪ Find a direction u , that maximises the rate of change of f as you move from v along u .

$$\text{Maximise } D_u[f](v)$$

$$f(u, y) = u e^{uy}$$

$$\nabla f(u, y) = (e^{uy} + uye^{uy}, u^2 e^{uy})$$

$$f(1, 0) = 1$$

$$\nabla f(1, 0) = (1, 1)$$

$$L(x, y) = f(1, 0) + \frac{\partial f(1, 0)}{\partial u}(x-1) + \frac{\partial f(1, 0)}{\partial y}(y-0)$$

$$= 1 + 1(x-1) + y = u + y$$

$$f(u) = \sqrt{u+4}$$

$$f'(u) = \frac{1}{2\sqrt{u+4}}$$

$$f(5) = 3$$

$$f'(5) = 1/6$$

$$\begin{aligned} L_a[f](u) &= f(a) + f'(a)(u-a) \\ &= 3 + \frac{1}{6}(u-5) \end{aligned}$$

$$= \frac{19}{6}$$

WEEK 3

→ Four Fundamental Subspaces

↳ Column Space

$$A = \begin{bmatrix} | & | & | \\ u_1 & u_2 & \dots & u_n \\ | & | & | \end{bmatrix}$$

$C(A) = \text{span}(u_1, u_2, \dots, u_n) = \{\text{linear combinations of the vectors } u_1, u_2, \dots, u_n\}$

Solving $Ax=b$

↳ For what b does $Ax=b$ have a solution?

↳ $\forall b \in C(A)$

e.g.

$$A = \begin{bmatrix} 1 & 1 & 2 \\ 2 & 1 & 3 \\ 3 & 1 & 4 \\ 4 & 1 & 5 \end{bmatrix}$$

→ Some vectors in \mathbb{R}^4 are not in $C(A)$ because

$Ax=b$ is "4 equations in 3 unknowns"

→ $\text{col } 3 = \text{col } 1 + \text{col } 2$. So, $C(A) = \text{span}(\text{col } 1, \text{col } 2)$

$C(A)$ is 2-dim. subspace of \mathbb{R}^4 .

→ Pivot columns in RREF

↳ Null Space

$$N(A) = \{u \mid Au=0\}$$

e.g.

$$A = \begin{bmatrix} 1 & 1 & 2 \\ 2 & 1 & 3 \\ 3 & 1 & 4 \\ 4 & 1 & 5 \end{bmatrix}$$

→ column 1 + column 2 = column 3. So the vector $(1, 1, -1)$ is in $N(A)$.

Note: If A is invertible, then $N(A)$ has "zero" only,

and $C(A)$ is the whole space. $Ax=b$ has a unique solution.

Else, $N(A)$ has $x_n \neq 0$ and $Ax=b$ solutions are of the form

$$x = x_p + x_n, \text{ s.t. } Ax_p=b, Ax_n=0$$

Note: Gaussian elimination to find $N(A)$. Solve $Ax=0$. Take the pivot variables and set them to basis entry

Rank = # of pivot columns

Rank = $\dim(C(A))$

Nullity = # of free variables

Nullity = $\dim(N(A))$

↳ Row Space

↳ column space of $A^T \Leftrightarrow$ span of rows of A

↳ col rank = $\dim(C(A))$ row rank = $\dim(R(A))$

col rank = row rank

↳ Left Null Space

$$\hookrightarrow N(A^T) = \{y \mid A^T y = 0\} = \{y \mid y^T A = 0\}$$

↳ linear combination of rows leading to zero vector

↳ A is $m \times n$ matrix

$$\dim(C(A)) + \dim(N(A)) = \# \text{ of col. in } A = n$$

$$\dim(C(A^T)) + \dim(N(A^T)) = \# \text{ of rows in } A = m$$

$$\dim(C(A^T)) = r \Rightarrow r + \dim(N(A^T)) = m \Rightarrow \dim(N(A^T)) = m - r$$

↳ e.g.:

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 6 \end{bmatrix}_{2 \times 2}$$

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

$$A^T = \begin{bmatrix} 1 & 3 \\ 2 & 6 \end{bmatrix}$$

$$N(A) = \text{line through } \begin{bmatrix} -2 \\ 1 \end{bmatrix}$$

$$C(A^T) = \text{line through } \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$N(A^T) = \text{line through } \begin{bmatrix} -3 \\ 1 \end{bmatrix}$$

→ Orthogonality

↳ length of a vector $u = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix} \in \mathbb{R}^n$, $\|u\|^2 = u_1^2 + u_2^2 + \dots + u_n^2$

↳ Orthogonal vector

↳ $u \perp y$ if $\underbrace{u^T y = 0}$ inner product

↳ if $\{v_1, v_2, \dots, v_n\}$ are mutually orthogonal, then they are linearly independent. To check, solve for $Au=0$.

↳ Orthonormal vector

↳ $\{u, v\}$ are orthonormal if $u^T v = 0$ and $\|u\| = \|v\| = 1$

↳ Orthogonal subspaces

↳ U, V are orthogonal subspaces

if $x^T y = 0 \forall x \in U, y \in V$

↳ e.g.

$$U = \text{span} \left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix} \right\} \quad \text{and} \quad V = \text{span} \left\{ \begin{pmatrix} 0 \\ 0 \\ 1 \\ 2 \end{pmatrix} \right\}$$

↳ Orthogonality w.r.t four fundamental subspaces of a matrix A

$$\textcircled{1} \quad R(A) \perp N(A) \Leftrightarrow C(A^T) \perp N(A)$$

$$\rightarrow n \in N(A) \Rightarrow An = 0$$

$$\begin{bmatrix} \text{row 1} \\ \text{row 2} \\ \vdots \\ \text{row } m \end{bmatrix} \begin{bmatrix} n_1 \\ n_2 \\ \vdots \\ n_m \end{bmatrix} = 0$$

implies row 1 $\perp n$, row 2 $\perp n$, ...

$$\textcircled{2} \quad C(A) \perp N(A^T)$$

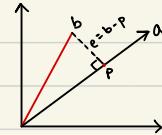
\rightarrow because $C(A^T) \perp N(A)$

→ Projections

↳ Why?

Inconsistent system of linear equations, i.e. For a system $Ax = b$, $b \notin C(A)$. In such situations, it makes sense to project b onto $C(A)$

↳ Projection onto a line



$$p = \hat{n}a; e = b - p = b - \hat{n}a; e \perp a$$

$$(b - \hat{n}a) \perp a$$

$$a(b - \hat{n}a) = 0 \Rightarrow ab - \hat{n}a^2a = 0$$

$$\text{projection of } b \text{ onto } a = \left(\frac{a^T b}{a^T a} \right) a$$

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

Cauchy-Schwarz inequality

$$\|e\|^2 = \|b - p\|^2 \geq 0 \Rightarrow \left\| b - \frac{a^T b}{a^T a} a \right\|^2 = b^T b - \frac{2(a^T b)^2}{a^T a} + \left(\frac{a^T b}{a^T a} \right)^2 a^T a$$

$$= \frac{(b^T b)(a^T a) - (a^T b)^2}{(a^T a)} \geq 0 \Rightarrow (b^T b)(a^T a) \geq (a^T b)^2$$

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

↳ Projection matrix

$$p = \left(\frac{a^T b}{a^T a} \right) a = \left(\frac{a a^T}{a^T a} \right) b \quad \text{let } P = \frac{a a^T}{a^T a}$$

Then, projection of b onto a is $P \cdot b$

↳ just left multiply by the projection matrix.

e.g.

$$a = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \quad P = \begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{bmatrix}$$

① P is symmetric

② $P^2 = P$ i.e., $P^2 b = Pb$

③ $C(P) = \text{line through } a; N(P) = \text{plane orthogonal to } a$

$\text{Rank}(P) = 1$

→ Least Squares

$$\left. \begin{array}{l} 2\kappa = b_1 \\ 3\kappa = b_2 \\ 4\kappa = b_3 \end{array} \right\} \text{system is solvable if } b \text{ is on line through } \begin{bmatrix} 2 \\ 3 \\ 4 \end{bmatrix}$$

- if b is such that the system is inconsistent

↳ pick a subset of equations and solve. Problem: large errors in some inputs and no errors in others.

↳ solution: minimize average error

$$E^2 = (2\kappa - b_1)^2 + (3\kappa - b_2)^2 + (4\kappa - b_3)^2$$

↳ minimize E^2

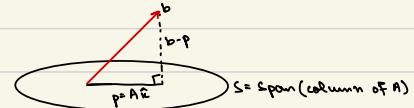
$$\frac{\partial E^2}{\partial \kappa} = 0 \Leftrightarrow 2[2(2\kappa - b_1) + 3(3\kappa - b_2) + 4(4\kappa - b_3)] = 0$$

$$\text{leads to } \hat{\kappa} = \frac{2b_1 + 3b_2 + 4b_3}{(2)^2 + (3)^2 + (4)^2} = \frac{\vec{a}^\top \vec{b}}{\vec{a}^\top \vec{a}}, \text{ with } \vec{a} = \begin{bmatrix} 2 \\ 3 \\ 4 \end{bmatrix}$$

→ taking a derivative and finding the minima is the same as taking a projection onto the column space

↳ Projection onto a subspace

→ project b onto $C(A)$



$$\text{orthogonal vector } e = b - p = b - A\hat{\kappa}$$

→ $e \perp \forall$ vectors in $C(A)$ and $C(A) \perp N(A^\top)$, therefore

$$e \in N(A^\top) \Rightarrow A^\top e = 0 \Rightarrow A^\top(b - A\hat{\kappa}) = 0$$

→ $A^\top A \hat{\kappa} = A^\top b$ → equation to solve for proj. of b onto $C(A)$.

↳ even if $A\kappa = b$ is inconsistent,

$A^\top A \hat{\kappa} = A^\top b$ has a solution.

→ if columns of A are lin. independent, then $A^\top A$ is invertible.

Solving $A^\top A \hat{\kappa} = A^\top b$ when $(A^\top A)$ is invertible $\Rightarrow \hat{\kappa} = (A^\top A)^{-1} A^\top b$

$$\text{Projection } P = A\hat{\kappa} = A(A^\top A)^{-1} A^\top b$$

→ if $b \in C(A)$, then $b = A\kappa$

$$\hookrightarrow \text{projection} = A(A^\top A)^{-1} A^\top A \kappa = A\kappa = b$$

→ if $b \in N(A^\top)$, then

$$P = A(A^\top A)^{-1} A^\top b = 0 \text{ because } A^\top b = 0$$

→ If A is square and invertible $\Leftrightarrow C(A) = \mathbb{R}^n$

$$P = A(A^\top A)^{-1} A^\top b = A A^{-1}(A^\top)^{-1} A^\top b = b$$

→ if $\text{Rank}(A) = 1$, then

$$\hat{\kappa} = \frac{\vec{a}^\top \vec{b}}{\vec{a}^\top \vec{a}}$$

↳ Projection matrix

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

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

$$\rightarrow P^2 = P ; P^2 = A(A^T A)^{-1} A^T A (A^T A)^{-1} A^T = A(A^T A)^{-1} A = P$$

→ projection matrix is symmetric and satisfies $P^2 = P$

↳ if a matrix A is symmetric and satisfies $A^T = A$, then A is a projection matrix.

→ Pb = projection of b onto the column space of P

→ Examples of least squares

↳ One-dimensional example:

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

$$b_i = \theta' x_i + \theta'' \text{ offset}$$

① system of equations:

$$\begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_m & 1 \end{bmatrix} \begin{bmatrix} \theta' \\ \theta'' \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix} \quad A\theta = b, \text{ where } \theta = \begin{bmatrix} \theta' \\ \theta'' \end{bmatrix}$$

② least squares approach: minimize $E^2 = \|b - A\theta\|^2 = (b_1 - \theta' x_1 - \theta'')^2 + \dots + (b_m - \theta' x_m - \theta'')^2$

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

e.g.

$$A\theta = b \quad A = \begin{bmatrix} -1 & 1 \\ 1 & 1 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} \theta' \\ \theta'' \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 3 \end{bmatrix} \quad \text{constant offset}$$

$$A^T A \hat{\theta} = A^T b, \text{ where } \hat{\theta} = \begin{bmatrix} \hat{\theta}' \\ \hat{\theta}'' \end{bmatrix}$$

$$A^T A = \begin{bmatrix} -1 & 1 & 2 \\ 1 & 1 & 1 \\ 2 & 1 & 1 \end{bmatrix}_{3 \times 3} \begin{bmatrix} -1 & 1 \\ 1 & 1 \\ 2 & 1 \end{bmatrix}_{3 \times 2} = \begin{bmatrix} 6 & 2 \\ 2 & 3 \\ 5 & 1 \end{bmatrix}; \quad \begin{bmatrix} -1 & 1 & 2 \\ 1 & 1 & 1 \\ 2 & 1 & 1 \end{bmatrix}_{3 \times 3} \begin{bmatrix} 1 \\ 1 \\ 3 \end{bmatrix}_{3 \times 1} = \begin{bmatrix} 6 \\ 5 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} 6 & 2 \\ 2 & 3 \\ 5 & 1 \end{bmatrix} \hat{\theta} = \begin{bmatrix} 6 \\ 5 \\ 1 \end{bmatrix}$$

$$\left[\begin{array}{cc|c} 6 & 2 & 6 \\ 2 & 3 & 5 \\ 5 & 1 & 1 \end{array} \right] \sim \left[\begin{array}{cc|c} 1 & 1/3 & 1 \\ 2 & 3 & 5 \\ 5 & 1 & 1 \end{array} \right] \sim \left[\begin{array}{cc|c} 1 & 1/3 & 1 \\ 0 & 7/3 & 3 \\ 5 & 1 & 1 \end{array} \right]$$

$$\hat{\theta}' = 4/7 \quad \hat{\theta}'' = 9/7$$



WEEK 4

→ Linear Regression

Given data $\{(u_1, y_1), \dots, (u_n, y_n)\}$, $u_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$, $i = 1, \dots, n$

$$L(\theta) = \frac{1}{2} \sum_{i=1}^n (u_i^\top \theta - y_i)^2 \rightarrow \text{Loss Function}$$

Minimise L : Define $A = \begin{bmatrix} u_1^\top \\ u_2^\top \\ \vdots \\ u_n^\top \end{bmatrix}$ ← feature matrix , $y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

$$A\theta - y = \begin{bmatrix} u_1^\top \theta - y_1 \\ \vdots \\ u_n^\top \theta - y_n \end{bmatrix} \quad (A\theta - y)^\top (A\theta - y) = \sum_{i=1}^n (u_i^\top \theta - y_i)^2$$

$$\begin{aligned} \text{Minimising } L : \quad \nabla_\theta L(\theta) &= 0 \Leftrightarrow \nabla_\theta ((A\theta - y)^\top (A\theta - y)) = 0 \\ &= A^\top (A\theta - y) = 0 \\ &= (A^\top A)\theta = A^\top y \rightarrow \text{Least squares solution} \end{aligned}$$

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

Maximum Likelihood and least squares regression end up solving the same equation, under a linear model.

→ Polynomial Regression

given one-dimensional data: $\{(u_1, y_1), \dots, (u_n, y_n)\}$ $u_i, y_i \in \mathbb{R} \forall i$
 polynomial of degree m

Transformed Features:

$$\begin{aligned} \hat{y}(u) &= \theta_0 + \theta_1 u + \theta_2 u^2 + \dots + \theta_m u^m \\ &= \sum_{j=0}^m \theta_j \phi_j(u) \text{, where } \phi_j(u) = u^j \end{aligned}$$

For given u , transformed feature vector $\phi(u) = (1, u, u^2, \dots, u^m)$

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

$$A = \begin{bmatrix} \phi(u_1)^\top \\ \phi(u_2)^\top \\ \vdots \\ \phi(u_n)^\top \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$(A^\top A)\theta = A^\top y \rightarrow \text{using the transformed features}$$

Regularised version of linear regression is a.k.a. **Ridge regression**

instead of solving $\min_{\theta} \frac{1}{2} \sum_{i=1}^n (u_i^\top \theta - y_i)^2$,

we solve → $\min_{\theta} \bar{L}(\theta) = \frac{1}{2} \sum_{i=1}^n (u_i^\top \theta - y_i)^2 + \underbrace{\lambda \|\theta\|^2}_{\text{regularisation term}}$

Eigenvalues and Eigenvectors

↳ Motivation: ODEs

$$\frac{dv}{dt} = 4v - 5w, v=8 \text{ at } t=0$$

$$\frac{dw}{dt} = 2v - 3w, w=5 \text{ at } t=0$$

$$u(t) = \begin{bmatrix} v(t) \\ w(t) \end{bmatrix}; u(0) = \begin{bmatrix} 8 \\ 5 \end{bmatrix}$$

$$\frac{du}{dt} = Au, \text{ where } A = \begin{bmatrix} 4 & -5 \\ 2 & -3 \end{bmatrix} \text{ and } u=u(0) \text{ at } t=0$$

→ if want solutions of this form:

$$v(t) = e^{kt}y, w(t) = e^{kt}z$$

or $u(t) = e^{kt}n$, where $n = \begin{bmatrix} y \\ z \end{bmatrix}$

$$\text{Substitute: } \lambda e^{kt}y = 4e^{kt}y - 5e^{kt}z$$

$$\lambda e^{kt}z = 2e^{kt}y - 3e^{kt}z$$

↓

$$\lambda y - 5z = \lambda y$$

$$2y - 3z = \lambda z$$

eigenvalue equation
 $\Leftrightarrow An = \lambda n$
 eigenvalue → eigen vector

→ can solve $\frac{du}{dt} = Au$ using solutions of the form $u(t) = e^{kt}n$

If $An = \lambda n$ can be solved.

↳ For a matrix A , λ is eigenvalue and $n \neq 0$ is eigenvector

If $An = \lambda n$



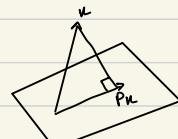
If n is eigenvector then An either stretches or shrinks n , but in the same direction.

→ If $\lambda = 0$, then $n \in N(A)$

↳ examples:

① Projection matrix P

↳ project onto a plane



(i) $Pn = n$ for any n in the plane

↳ $\lambda = 1$ is an eigenvalue and any n in the plane is an eigenvector

(ii) if an n is \perp to the plane. $Pn = 0$

so, $\lambda = 0$ is an eigenvalue and any $n \perp$ to the plane is an eigenvector.

② Permutation matrix

↪ identity matrix with the rows shuffled.

$$\text{e.g. } B = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$B\mathbf{u} = \mathbf{u} \text{ for } \mathbf{u} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad B\mathbf{u} = -\mathbf{u} \text{ for } \mathbf{u} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

↪ Finding eigenvalues:

↪ $A\mathbf{u} = \lambda\mathbf{u} \Leftrightarrow (A - \lambda I)\mathbf{u} = 0 \Leftrightarrow (A - \lambda I)$ is singular

$$\Leftrightarrow \underbrace{\det(A - \lambda I) = 0}_{\text{characteristic polynomial of matrix } A}$$

→ If A is $n \times n$, then characteristic polynomial is of degree n

$$\Leftrightarrow (a_{11} - \lambda) \dots (a_{nn} - \lambda) = 0$$

• "n" roots of characteristic polynomial = eigenvalues

$$\begin{aligned} \rightarrow A\mathbf{u} &= \lambda\mathbf{u} \\ \Rightarrow A\mathbf{u} - \lambda\mathbf{u} &= 0 \\ \Rightarrow A\mathbf{u} - \lambda I\mathbf{u} &= 0 \\ \Rightarrow (A - \lambda I)\mathbf{u} &= 0 \\ \Rightarrow A - \lambda I &= 0 \\ \det(A - \lambda I) &= 0 \end{aligned}$$

For eigenvectors: Given λ , want $(A - \lambda I)\mathbf{u} = 0$

$$\text{i.e., } \mathbf{u} \in \underbrace{N(A - \lambda I)}_{\text{find using Gaussian elimination}}$$

↪ example

$$\textcircled{1} \quad A = \begin{bmatrix} 3 & 1 \\ 1 & 3 \end{bmatrix} \quad \lambda I = \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} \quad A - \lambda I = \begin{bmatrix} 3-\lambda & 1 \\ 1 & 3-\lambda \end{bmatrix}$$

$$\det(A - \lambda I) = (3-\lambda)^2 - 1 = 9 - 6\lambda + \lambda^2 - 1 = \lambda^2 - 6\lambda + 8$$

$$\lambda^2 - 6\lambda + 8 = 0 \rightarrow \lambda^2 - 4\lambda - 2\lambda + 8 = 0 \rightarrow \lambda(\lambda - 4) - 2(\lambda - 4) = 0$$

$$\leftarrow (\lambda - 2)(\lambda - 4) = 0$$

$$\boxed{\text{eigenvalues: } \lambda = 2, \lambda = 4}$$

$$\lambda_1 + \lambda_2 = \text{trace}(A)$$

$$\lambda_1 \cdot \lambda_2 = \det(A)$$

$$\text{eigenvectors: } A - 4I = \begin{bmatrix} -1 & 1 \\ 1 & -1 \end{bmatrix}; (A - 4I)\mathbf{u} = 0 \text{ for } \boxed{\mathbf{u} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}}$$

$$A - 2I = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}; (A - 2I)\mathbf{u} = 0 \text{ for } \boxed{\mathbf{u} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}}$$

③ variation of ①

$$B = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad B + 3I = A$$

$$A\mathbf{u} = (B + 3I)\mathbf{u} = \lambda\mathbf{u} + 3\mathbf{u} = (\lambda + 3)\mathbf{u}$$

eigenvectors of B = eigenvectors of A

↪ Remarks

① Suppose $A\mathbf{u} = \lambda_1\mathbf{u}$ and $B\mathbf{u} = \lambda_2\mathbf{u}$

Then, $(A + B)\mathbf{u} \neq (\lambda_1 + \lambda_2)\mathbf{u}$ because \mathbf{u} can be diff for A and B .

② Symmetric matrix has "real" eigenvalues.

for matrix $A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$, eigenvalues: $\lambda = i, \lambda_2 = -i$

③ A matrix with linearly dependent eigenvectors:

$$A = \begin{bmatrix} 3 & 1 \\ 0 & 3 \end{bmatrix} \quad \lambda = \lambda_2 = 3 \quad \boxed{\mathbf{u} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}}$$

→ Similarity and Diagonalisation

↳ Matrix A is diagonalisable if \exists an invertible matrix S such that $S^{-1}AS = \Lambda$ \rightarrow diagonal matrix

↳ Matrix is diagonalisable if it has enough independent eigenvectors

Suppose A is $n \times n$ matrix with n lin. independent eigenvectors $= \{u_1, u_2, \dots, u_n\}$

let $S = \begin{bmatrix} 1 & 1 & 1 \\ u_1 & u_2 & \dots & u_n \\ 1 & 1 & 1 \end{bmatrix}$ S is invertible because $\text{rank}(S) = n$

$$AS = A \begin{bmatrix} 1 & 1 & 1 \\ u_1 & u_2 & \dots & u_n \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 \\ Au_1 & Au_2 & \dots & Au_n \\ 1 & 1 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 1 & 1 \\ \lambda_1 u_1 & \lambda_2 u_2 & \dots & \lambda_n u_n \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 \\ u_1 & u_2 & \dots & u_n \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_n \end{bmatrix} = S\Lambda, \Lambda = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_n \end{bmatrix}$$

$$\text{So, } AS = S\Lambda \Leftrightarrow S^{-1}AS = \Lambda$$

↳ eigenvectors corresponding to different eigenvalues are linearly independent.

↳ Remarks:

① $S^{-1}AS = \Lambda$ or $A = S\Lambda S^{-1}$

S is not unique

② $A = S\Lambda S^{-1}$

Suppose col. of S is y . Col. of $S\Lambda = \lambda_1 y$. Col. of $AS = Ay$

$Ay = \lambda_1 y \Rightarrow \lambda_1$ is an eigenvalue and y is an eigenvector

③ Powers of A

Suppose λ is an eigenvalue and u is an eigenvector

$$A^2u = A(Au) = \lambda Au = \lambda^2 u$$

Then, $S^{-1}A^2S = \Lambda^2$

④ Not all matrices are diagonalisable

↳ Fibonacci sequence

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

system of equations : $F_{k+2} = F_{k+1} + F_k$

$$u_k = \begin{bmatrix} F_{k+1} \\ F_k \end{bmatrix}, \quad u_{k+1} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} u_k \quad \Leftrightarrow \quad \begin{bmatrix} F_{k+2} \\ F_{k+1} \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} F_{k+1} \\ F_k \end{bmatrix}$$

$$\text{let } A = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \quad u_0 \xrightarrow{A} u_1 \xrightarrow{A} u_2 \xrightarrow{A} \dots$$

$u_k = A^k u_0$ is a solution to $u_{k+1} = Au_k$

If A is diagonalisable, then $u_0 = c_1 u_1 + c_2 u_2 + \dots + c_n u_n$

$$u_k = A^k u_0 = c_1 \lambda_1^k u_1 + c_2 \lambda_2^k u_2 + \dots + c_n \lambda_n^k u_n$$

$$A = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}, \quad \lambda_1 = \frac{1+\sqrt{5}}{2}, \quad \lambda_2 = \frac{1-\sqrt{5}}{2}$$

with $u_1 = \begin{bmatrix} \lambda_1 \\ 1 \end{bmatrix}$ and $u_2 = \begin{bmatrix} \lambda_2 \\ 1 \end{bmatrix}$

$$(A - \lambda_1 I) \begin{bmatrix} \lambda_1 \\ 1 \end{bmatrix} = \begin{bmatrix} \lambda_1^2 - \lambda_1 - 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$(A - \lambda_2 I) \begin{bmatrix} \lambda_2 \\ 1 \end{bmatrix} = \begin{bmatrix} \lambda_2^2 - \lambda_2 - 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

So, u_1 and u_2 are the eigenvectors

→ u_0 as a linear combination of u_1, u_2 :

$$u_0 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} = c_1 \begin{bmatrix} (1+\sqrt{5})/2 \\ 1 \end{bmatrix} + c_2 \begin{bmatrix} (1-\sqrt{5})/2 \\ 1 \end{bmatrix}$$

$$c_1 = \frac{1}{\sqrt{5}}, \quad c_2 = -\frac{1}{\sqrt{5}}$$

$$u_k = c_1 \lambda_1^k u_1 + c_2 \lambda_2^k u_2$$

$$\begin{bmatrix} F_{k+1} \\ F_k \end{bmatrix} = \frac{1}{\sqrt{5}} \left(\frac{1+\sqrt{5}}{2} \right)^k \begin{bmatrix} (1+\sqrt{5})/2 \\ 1 \end{bmatrix} + \left(-\frac{1}{\sqrt{5}} \right) \left(\frac{1-\sqrt{5}}{2} \right)^k \begin{bmatrix} (1-\sqrt{5})/2 \\ 1 \end{bmatrix}$$

$$F_k = \underbrace{\frac{1}{\sqrt{5}} \left(\frac{1+\sqrt{5}}{2} \right)^k}_{\text{becomes negligible as } k \text{ increases}} - \underbrace{\frac{1}{\sqrt{5}} \left(\frac{1-\sqrt{5}}{2} \right)^k}_{\text{becomes negligible as } k \text{ increases}}$$

$$F_k \approx \frac{1}{5} \left(\frac{1+\sqrt{5}}{2} \right)^k \quad F_{100} \approx \frac{1}{\sqrt{5}} \left(\frac{1+\sqrt{5}}{2} \right)^{100} \approx 3.54 \times 10^20$$

↳ Orthogonally diagonalisable

→ If A is a real symmetric matrix, then:

① Eigenvalues of A are real.

② Eigenvectors corresponding to diff. eigenvalues are linearly independent.

③ A is orthogonally diagonalisable.

→ Matrix A is orthogonally diagonalisable if

$\exists Q$ satisfying $Q^T Q = I$ s.t. $A = Q \Lambda Q^T$

$$A = \begin{bmatrix} 1 & & & \\ u_1 & \dots & u_n \\ 1 & & 1 \end{bmatrix} \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & 0 \\ 0 & & \lambda_n \end{bmatrix} \begin{bmatrix} u_1^T \\ \vdots \\ u_n^T \end{bmatrix}$$

→ Remarks:

① A is diagonalisable + orthogonal matrix for diagonalisation \Rightarrow orthogonally diagonalisable

② Real symmetric matrix is orthogonally diagonalisable

$$\rightarrow \text{e.g. } A = \begin{bmatrix} 1 & -2 \\ -2 & -2 \end{bmatrix} \quad \det \begin{pmatrix} 1-\lambda & -2 \\ -2 & -2-\lambda \end{pmatrix} = 0 \quad ; \quad \lambda_1 = 2, \lambda_2 = -3$$

$$(A - 2\mathbb{I})\mathbf{u} = 0 \quad \begin{bmatrix} -1 & -2 \\ -2 & -4 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad u_1 = \begin{bmatrix} -2 \\ 1 \end{bmatrix} \quad u_2 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$(A + 3\mathbb{I})\mathbf{u} = 0 \quad \begin{bmatrix} 4 & -2 \\ -2 & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$q_1 = \frac{u_1}{\|u_1\|} = \frac{1}{\sqrt{5}} \begin{bmatrix} -2 \\ 1 \end{bmatrix}, \quad q_2 = \frac{u_2}{\|u_2\|} = \frac{1}{\sqrt{5}} \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$Q = \frac{1}{\sqrt{5}} \begin{bmatrix} -2 & 1 \\ 1 & 2 \end{bmatrix}$$

— X — X — X — X —

WEEK 5

→ Complex Matrices

→ \mathbb{C}^n : complex counterpart of \mathbb{R}^n
 $(u_1, u_2, \dots, u_n) \in \mathbb{C}^n$, then $u_i \rightarrow$ complex numbers

→ complex conjugate of $(a+ib)$ is $(a-ib)$

→ Linear combinations: $c_1 u_1 + c_2 u_2 + \dots + c_n u_n \xrightarrow{\text{complex numbers}} 0$

→ Inner product and length

$$\text{In } \mathbb{R}^n, \quad \underbrace{\|u\|^2}_{\text{length}} = \underbrace{u^\top u}_{\text{inner product}}$$

$$\text{In } \mathbb{C}^n, \quad u \cdot y = \bar{u}^\top y = \bar{u}_1 y_1 + \bar{u}_2 y_2 + \dots + \bar{u}_n y_n \quad \xrightarrow{\text{conjugate}} \\ \bar{u}^\top y \neq \bar{y}^\top u$$

length of a complex vector: For $u \in \mathbb{C}^n$, define $\|u\|^2 = \bar{u}^\top u$

$$\rightarrow \textcircled{1} \quad u \cdot y = \bar{y} \cdot u$$

$$\textcircled{2} \quad u \cdot (cy) = c(u \cdot y)$$

$$\textcircled{3} \quad (cu) \cdot y = \bar{c}(u \cdot y)$$

→ Conjugate transpose

$A^* = \text{conjugate transpose of } A$
 $A^* = \bar{A}^T = \overline{A^T}$

→ Remarks

→ Real matrix A , $A^* = A^T$

→ $(AB)^* = B^* A^*$

$$\rightarrow u \cdot y = [\bar{u}_1, \bar{u}_2, \dots, \bar{u}_n] \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = u^* y$$

→ Hermitian Matrix

↪ A is hermitian matrix if

$$A^* = A$$

↪ complex equivalent of symmetric matrix in \mathbb{R}^n

→ Hermitian Matrices

→ Matrix A is Hermitian if $A^* = A$

$$\text{e.g. } A = \begin{bmatrix} 2 & 3-3i \\ 3+3i & 5 \end{bmatrix} \quad A^* = \bar{A}^T = \begin{bmatrix} 2 & 3+3i \\ 3-3i & 5 \end{bmatrix}^T = \begin{bmatrix} 2 & 3-3i \\ 3+3i & 5 \end{bmatrix}$$

→ Diagonal entries in Hermitian matrices have to be real

→ Properties:

① If A is Hermitian, then all eigenvalues are real.

$$\text{e.g. } A = \begin{bmatrix} 2 & 3-3i \\ 3+3i & 5 \end{bmatrix} \quad (A - \lambda I) = \begin{bmatrix} 2-\lambda & 3-3i \\ 3+3i & 5-\lambda \end{bmatrix}$$

$$\begin{aligned} \det(A - \lambda I) &= (2-\lambda)(5-\lambda) - (3+3i)(3-3i) \\ &= (2-\lambda)(5-\lambda) - (9 - 9i^2) = (2-\lambda)(5-\lambda) - 18 \\ &= 10 - 2\lambda - 5\lambda + \lambda^2 - 18 = (\lambda-8)(\lambda+1) \end{aligned}$$

Proof: $Au = \lambda u ; (Au)^* = (\lambda u)^* = \bar{\lambda} u^*$

$$\Leftrightarrow u^* A^* = \bar{\lambda} u^* \Leftrightarrow u^* A^* u = \bar{\lambda} u^* u \Leftrightarrow u^* A u = \bar{\lambda} u^* u$$

$$\Leftrightarrow u^* \lambda u = \bar{\lambda} u^* u \Leftrightarrow \lambda u^* u = \bar{\lambda} u^* u \Leftrightarrow \lambda = \bar{\lambda}$$

② If A is Hermitian, then eigenvectors corresponding to diff. eigenvalues are orthogonal i.e., If $Au = \lambda_1 u$ and $Ay = \lambda_2 y$ where $\lambda_1 \neq \lambda_2$ then $u \cdot y = \bar{u}^T y = 0$

Proof: $Au = \lambda_1 u$ and $Ay = \lambda_2 y$, $\lambda_1 \neq \lambda_2$

Show: $u \cdot y = 0$

$$\rightarrow u \cdot Ay = u \cdot \lambda_2 y = \lambda_2 (u \cdot y)$$

$$\rightarrow u \cdot Ay = u^* Ay = u^* A^* y = (Au)^* y = (\lambda_1 u)^* y = \bar{\lambda}_1 (u \cdot y) = \lambda_1 (u \cdot y)$$

$$\rightarrow u \cdot Ay = \lambda_2 (u \cdot y) = \lambda_1 (u \cdot y)$$

Since $\lambda_1 \neq \lambda_2$, we have $u \cdot y = u^* y = 0$

→ Remarks

① equivalent of Hermitian matrices in $\mathbb{R}^n \rightarrow$ symmetric matrix
 ↳ All real symmetric matrices are Hermitian.

② If no eigenvalues are repeated, then A is diagonalisable.

→ Unitary Matrices

→ A matrix is unitary if it is a square matrix and has orthonormal columns.

→ In \mathbb{R}^n , $Q^T Q = I$, then Q is orthogonal matrix

$$\hookrightarrow \begin{bmatrix} 1 & 1 & 1 \\ v_1 & v_2 & \dots & v_n \\ 1 & 1 & 1 \end{bmatrix} \text{ then } v_i^T v_j \quad \forall i \neq j \quad \& \quad \|v_i\| = 1 \quad \forall i = 1, 2, \dots, n$$

→ In \mathbb{C}^n , $U^* U = I$, then U is unitary and $U^{-1} = U^*$

→ e.g.

$$U = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{i}{\sqrt{2}} \\ \frac{i}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}, \quad U = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}$$

→ Properties:

① length unchanged

$$\hookrightarrow \|Uu\| = \|u\|$$

Proof:

$$Uu \cdot Uy = (Uu)^* Uy = u^* U^* Uy = u^* y = u \cdot y$$

② Eigenvalues of unitary matrix have absolute value 1

Proof:

$$Uu = \lambda u ; \|Uu\| = \|u\| \Rightarrow \|\lambda u\| = \|u\| \Rightarrow |\lambda| = 1$$

③ Eigenvectors corresponding to diff. eigenvalues are orthogonal.

Proof:

$$Uu = \lambda_1 u, \quad Uy = \lambda_2 y, \quad \lambda_1 \neq \lambda_2$$

$$u \cdot y = Uu \cdot Uy = (\lambda_1 u) \cdot (\lambda_2 y) = \bar{\lambda}_1 \lambda_2 (u \cdot y)$$

$$\Rightarrow (\bar{\lambda}_1 \lambda_2 - 1)(u \cdot y) = 0$$

⇒ For $(u \cdot y) = 0$, λ_1 needs to be equal to λ_2 .

→ Diagonalization of Hermitian Matrices

$$\rightarrow A^* = \bar{A}^t$$

Matrix A is Hermitian if $A^* = A$

Matrix U is Unitary if $U^*U = I$ and U is square matrix.

→ A is unitary diagonalizable \exists matrix U s.t.

$$A = U \Lambda U^*$$
, where Λ is a diagonal matrix

→ Any $n \times n$ matrix is similar to an upper triangular matrix

$$\hookrightarrow A = \underset{\text{unitary}}{U} \underset{\text{upper triangular}}{T} \underset{\text{unitary}}{U}^*$$

→ Shur's Theorem

$$\hookrightarrow A = U T U^*$$

↪ Proof for $n=3$:

let $p(\lambda)$ be characteristic polynomial of A.

let λ_1 be a root of $p(\lambda)$. let z_1 be corresponding eigenvector

Extend $\{z_1\}$ to a basis, and make it orthonormal (Gramm-Schmidt)

let $\{z_1, u, v\}$ be the orthonormal basis.

↪ i.e., $\|z_1\| = \|u\| = \|v\| = 1$

$$z_1^* u = z_1^* v = u^* v = 0$$

$$\text{let } U_1 = \begin{bmatrix} 1 & 1 & 1 \\ z_1 & u & v \\ 1 & 1 & 1 \end{bmatrix} \quad AU_1 = A \begin{bmatrix} 1 & 1 & 1 \\ z_1 & u & v \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & & 1 \\ \lambda_1 z_1 & Au & Av \\ 1 & & 1 \end{bmatrix}$$

\downarrow
 $Az_1 = \lambda_1 z_1$

$$U_1^* A U_1 = \begin{bmatrix} -\bar{z}_1^* & & \\ -\bar{u}^* & & \\ -\bar{v}^* & & \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ \lambda_1 z_1 & Au & Av \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} \lambda_1 & * & * \\ 0 & B & \\ 0 & & \end{bmatrix}$$

some non-zero entries
(2x2) matrix

$$\bar{z}_1^* \lambda_1 z_1 = \lambda_1, \bar{z}_1^* z_1 = 1, \bar{u}^* z_1 = 0 \text{ and } \bar{v}^* z_1 = 0$$

$$U_1^* A U_1 = \begin{bmatrix} \lambda_1 & * & * \\ 0 & B & \\ 0 & & \end{bmatrix}$$

→ Repeat the procedure with B to get eigenvalue λ_2 of B and a unitary matrix P s.t. $P^* B P = \begin{bmatrix} \lambda_2 & * \\ 0 & \lambda_3 \end{bmatrix}$

$$\text{let } U_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & P & \\ 0 & & \end{bmatrix} \text{ . Then, } U_2^* = \begin{bmatrix} 1 & 0 & 0 \\ 0 & P^* & \\ 0 & & \end{bmatrix}$$

$$U_2^* U_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \text{ since } P^* P = I$$

$$\begin{aligned} U_2^* (U_1^* A U_1) U_2 &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & P^* & \\ 0 & & \end{bmatrix} \begin{bmatrix} \lambda_1 & * & * \\ 0 & B & \\ 0 & & \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & P & \\ 0 & & \end{bmatrix} = \begin{bmatrix} \lambda_1 & * & * \\ 0 & \lambda_2 & * \\ 0 & 0 & \lambda_3 \end{bmatrix} \\ &= \begin{bmatrix} \lambda_1 & * & * \\ 0 & P^* B P & \\ 0 & & \end{bmatrix} \end{aligned}$$

$$\rightarrow \text{e.g. } A = \begin{bmatrix} 5 & 8 & 16 \\ 5 & 0 & 9 \\ -3 & -5 & -10 \end{bmatrix} \quad p(\lambda) = -(\lambda-1)(\lambda+3)^2 \quad \lambda_1 = 1 ; \lambda_2 = -3$$

Step ① eigenvector z_1 corresponding to $\lambda_1 = \begin{bmatrix} -2 \\ -1 \\ 1 \end{bmatrix}$

Step ② $\{z_1, e_1, e_2\} = \left\{ \begin{bmatrix} -2 \\ -1 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \right\}$ ← basis for \mathbb{R}^3

Graham-Schmidt procedure:

$$w_1 = \frac{z_1}{\|z_1\|} = \frac{1}{\sqrt{6}} \begin{bmatrix} -2 \\ -1 \\ 1 \end{bmatrix}$$

$$v_2 = z_2 - \langle z_2, w_1 \rangle w_1 ; \quad w_2 = \frac{v_2}{\|v_2\|}$$

$$w_2 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} - \left(-\frac{1}{3}\right) \begin{bmatrix} -2 \\ -1 \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{4}{3} \\ \frac{2}{3} \\ \frac{1}{3} \end{bmatrix}$$

$$w_3 = \begin{bmatrix} 0 \\ \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$$

Step ③ $U_1 = \begin{bmatrix} -\frac{2}{\sqrt{6}} & \frac{1}{\sqrt{3}} & 0 \\ -\frac{1}{\sqrt{6}} & -\frac{1}{\sqrt{3}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{2}} \end{bmatrix} \quad U_1^* A U_1 = U_1^* A U_1 = \begin{bmatrix} 1 & -8\sqrt{2} & -12\sqrt{3} \\ 0 & 1 & 0 \\ 0 & \sqrt{6} & -3 \end{bmatrix} \underbrace{\qquad\qquad\qquad}_{B}$

Step ④

Find an eigenvalue of B

$$\lambda_2 = -3, \quad e_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad \text{eigenvector corresponding to } \lambda_2$$

$$\text{extend } \{e_2\} = \{e_2, e_1\} = \left\{ \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \right\}$$

$$\text{unitary matrix } P = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$P^* B P = P^* B P = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} -3 & 0 \\ \sqrt{6} & -3 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} -3 & \sqrt{6} \\ 0 & -3 \end{bmatrix}$$

$$U_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & P \\ 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

$$U_2^* U_1^* A U_1 U_2 = \begin{bmatrix} 1 & -8\sqrt{2} & -12\sqrt{3} \\ 0 & -3 & \sqrt{6} \\ 0 & 0 & -3 \end{bmatrix}$$

→ Spectral Theorem

→ Hermitian matrix A is unitarily diagonalisable

i.e., \exists unitary U s.t. $U^* A U = D$ → diagonal matrix with real numbers

→ Proof:

Schur's theorem: $U^* A U = T$

$$T^* = U^* A^* U \xrightarrow{U^* A U = T} \text{because } A^* = A$$

$T \rightarrow$ upper Δ matrix and $T^* \rightarrow$ lower triangular matrix

$T^* = T \Rightarrow T$ is a diagonal matrix

$\hookrightarrow T_{ii} = T_{ii}^* = \bar{T}_{ii} \Rightarrow T_{ii}$ is a real number

→ example:

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

$$p(\lambda) = (\lambda - 1)(\lambda - 4) \quad \text{eigenvalues: } \lambda_1 = 1, \lambda_2 = 4$$

$$\text{eigen vectors: } z_1 = \begin{bmatrix} -1+i \\ 1 \end{bmatrix}; \quad z_2 = \begin{bmatrix} 1-i \\ 2 \end{bmatrix}$$

$$\frac{z_1}{\|z_1\|} = u_1 = \frac{1}{\sqrt{3}} \begin{bmatrix} -1+i \\ 1 \end{bmatrix}; \quad u_2 = \frac{1}{\sqrt{6}} \begin{bmatrix} 1-i \\ 2 \end{bmatrix}$$

$$U = \begin{bmatrix} \frac{-1+i}{\sqrt{3}} & \frac{1-i}{\sqrt{6}} \\ \frac{1}{\sqrt{3}} & \frac{2}{\sqrt{6}} \end{bmatrix}; \quad U^* A U = \begin{bmatrix} 1 & 0 \\ 0 & 4 \end{bmatrix}$$

→ Remark: Hermitian ⇒ unitarily diagonalisable
but that does not mean only unitarily diagonalisable matrices are Hermitian.



WEEK 6

→ Singular Value Decomposition

→ Spectral theorem

↳ If A is a real symmetric matrix, then:

- ① All eigenvalues of A are real
- ② A is orthogonally diagonalisable

→ Every matrix cannot be diagonalised but any "real" $m \times n$ matrix A can be decomposed in SVD form i.e.,

A can be written as $A = Q_1 \Sigma Q_2^T$, where Q_1 and Q_2 are orthogonal

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

→ Proof:

If A is $m \times n$, A^T is $n \times m$, $A^T A$ is symmetric and real

⇒ there exists a basis of orthonormal eigenvectors $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ corresponding to eigenvalues $\{\lambda_1, \dots, \lambda_n\}$

$$A^T A \mathbf{u}_i = \lambda_i \mathbf{u}_i, \forall i=1, \dots, n$$

orthonormal eigenvector

$$(A^T A \mathbf{u}_i) \cdot \mathbf{u}_i = (\lambda_i \mathbf{u}_i) \cdot \mathbf{u}_i = \lambda_i \text{ since } \|\mathbf{u}_i\|^2 = 1$$

Order the eigenvalues of $A^T A$: $\lambda_1, \lambda_2, \dots, \lambda_r, \lambda_{r+1}, \dots, \lambda_n$

$$\lambda_1 > 0, \dots, \lambda_r > 0; \lambda_{r+1} = \dots = \lambda_n = 0$$

σ_i 's are singular values and $\forall i=1, \dots, r, \sigma_i = \sqrt{\lambda_i}$

Let $y_i = \frac{1}{\sigma_i} A \mathbf{u}_i$. For $i=1, \dots, r$, $y_i \in \mathbb{R}^m$

$$\|y_i\| = \frac{1}{\sigma_i} \|A \mathbf{u}_i\| = \frac{\sqrt{\lambda_i}}{\sigma_i} = 1$$

$\underbrace{\quad}_{\text{because } \|A \mathbf{u}_i\|^2 = \lambda_i}$

$$y_i \cdot y_j = \frac{1}{\sigma_i \sigma_j} (A \mathbf{u}_i) \cdot (A \mathbf{u}_j) = \frac{1}{\sigma_i \sigma_j} \mathbf{u}_i^T A^T A \mathbf{u}_j = \frac{1}{\sigma_i \sigma_j} \mathbf{u}_i^T \lambda_j \mathbf{u}_j$$

$$= \frac{\lambda_j}{\sigma_i \sigma_j} \mathbf{u}_i^T \mathbf{u}_j = 0 \text{ since } \{\mathbf{u}_1, \dots, \mathbf{u}_n\} \text{ are orthonormal}$$

→ So we have a set $\{y_1, \dots, y_r\}$ of orthonormal vectors.

→ Since $r \leq m$, set $\{y_1, \dots, y_r\}$ is not a basis. So extend the set $\{y_1, \dots, y_r\}$ to form orthonormal basis of \mathbb{R}^m giving the set $\{y_1, \dots, y_m\}$

$$Q_1 = \begin{bmatrix} 1 & & & \\ y_1 & \dots & y_m \\ | & & | \end{bmatrix} \text{ and } Q_2 = \begin{bmatrix} 1 & & & \\ \mathbf{u}_1 & \dots & \mathbf{u}_n \\ | & & | \end{bmatrix}$$

$$\begin{aligned}\Sigma &= Q_1^T A Q_2 \\ &= \begin{bmatrix} -y_1^T & \dots & -y_m^T \end{bmatrix} \begin{bmatrix} 1 & & \\ & A_{K_1} & \dots & A_{K_m} \\ 1 & & \end{bmatrix}\end{aligned}$$

$\Sigma_{ij} = y_i^T (A x_j)$, where $\sigma_i = \sqrt{\lambda_i}$ and x_i = a basis vector, λ_i = eigenvalue

$$y_i^\top A x_j = \begin{cases} \sigma_j & \text{If } i=j \\ 0 & \text{otherwise} \end{cases}$$

For $j \leq r$ 
 For $j > r$ 

$\|A x_j\|^2 = \lambda_j = 0$
 $\Rightarrow y_i^\top A x_j = 0$

$$\sum = \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \\ \hline 0 & 0 \end{bmatrix}$$

→ Remarks

$$\textcircled{1} \quad A = Q, \Sigma Q^T$$

$$AA^T = Q, \Sigma \Sigma^T Q^T \quad \left. \begin{array}{l} \{ \\ \rightarrow \end{array} \right. \begin{array}{l} \text{eigen decomposition of } AA^T \\ \text{Q}_1 \text{ contains eigenvectors of } AA^T \end{array}$$

$$\textcircled{2} \quad A^T A = Q_2 \Sigma^T \Sigma Q_2^T$$

→ Q_2 contains eigenvectors of $A A^T$

→ Procedure of SVD:

- ① Find eigenvalues and eigenvectors of $A^T A$. Using eigenvectors form an orthonormal basis.

- ② Use $\sigma_i y_i = A u_i$, $\sigma_2 y_2 = A u_2$, where $\{u_1, u_2\}$ is the basis from step ①, and $\sigma_i = \sqrt{\lambda_i}$. Find y_1 and y_2 .

- $$③ A = Q_1 \sum Q_2^T \text{ s.t.}$$

$$Q_1 = \begin{bmatrix} 1 & 1 \\ y_1 & \dots & y_n \\ 1 & 1 \end{bmatrix}, \quad Q_2 = \begin{bmatrix} 1 & 1 \\ x_1 & \dots & x_n \\ 1 & 1 \end{bmatrix} \quad \text{and} \quad \Sigma = \begin{bmatrix} \sqrt{\lambda_1} & 0 \\ 0 & \sqrt{\lambda_n} \end{bmatrix}$$

→ example:

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

Is A diagonalisable? No, because $\lambda_1 = \lambda_2 = \sqrt{2}$

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

① $A^T A = \begin{bmatrix} 2 & \sqrt{2} \\ \sqrt{2} & 3 \end{bmatrix} \quad \lambda_1 = 4, \lambda_2 = 1$
 $\sigma_1 = 2, \sigma_2 = 1$

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

eigenvectors of $A^T A$: $A^T A - 4I = \begin{bmatrix} -2 & \sqrt{2} \\ \sqrt{2} & -1 \end{bmatrix} \quad z_1 = \begin{bmatrix} 1 \\ \sqrt{2} \end{bmatrix}$

$$A^T A - I = \begin{bmatrix} 1 & \sqrt{2} \\ \sqrt{2} & 2 \end{bmatrix} \quad z_2 = \begin{bmatrix} \sqrt{2} \\ -1 \end{bmatrix}$$

normalise eigenvectors: $v_1 = \frac{1}{\sqrt{3}} \begin{bmatrix} 1 \\ \sqrt{2} \end{bmatrix}$ and $v_2 = \frac{1}{\sqrt{3}} \begin{bmatrix} \sqrt{2} \\ -1 \end{bmatrix}$

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

② Use $\sigma_i y_i = A v_i$ and $\sigma_2 y_2 = A v_2$

to obtain $y_1 = \frac{1}{\sqrt{3}} \begin{bmatrix} \sqrt{2} \\ 1 \end{bmatrix} \quad y_2 = \frac{1}{\sqrt{3}} \begin{bmatrix} 1 \\ -\sqrt{2} \end{bmatrix}$

$$Q_1 = \frac{1}{\sqrt{3}} \begin{bmatrix} \sqrt{2} & 1 \\ 1 & -\sqrt{2} \end{bmatrix}$$

Positive Definiteness

→ A function that is that has first derivative 0 at (0,0) and is strictly positive at other points is called positive definite.

→ example:

$$f(x,y) = 2x^2 + 4xy + y^2$$

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

(0,0) is a stationary point.

$$\frac{\partial^2 f}{\partial x^2} = 4 \quad \frac{\partial^2 f}{\partial y^2} = 2 \quad \frac{\partial^2 f}{\partial x \partial y} = 4 \quad \Rightarrow f \text{ has a minima at (0,0)}$$

→ conditions for a function to be positive definite.

let $f(x,y)$ be of the form $ax^2 + 2bxy + cy^2$

① If $f > 0$, then $a > 0$

② If $f > 0$, then $c > 0$

③ If $f > 0$, then $ac > b^2$

→ Remarks:

④ If $ac = b^2$, then $f(x,y) = ax^2 + 2bxy + cy^2$ is $\begin{cases} \text{positive semi-definite if } a > 0 \\ \text{negative semi-definite if } a < 0 \end{cases}$

⑤ We have saddle point at $(0,0)$ if $ac < b^2$

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

$$\text{let } v = \begin{bmatrix} x \\ y \end{bmatrix} \text{ and } A = \begin{bmatrix} a & b \\ b & c \end{bmatrix}$$

$$\text{Then, } ax^2 + 2bxy + cy^2 = v^T A v$$

In general,

$$[x_1 \dots x_n] \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j$$

At $v = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$, $f(v) = 0$ i.e., $\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$ is a stationary point of f

→ Positive Definite Matrices

→ $A = \begin{bmatrix} a & b \\ b & c \end{bmatrix}$ is positive definite if $a > 0$ and $ac - b^2 > 0$

If $a > 0$ and $ac - b^2 > 0$, then both eigenvalues of A are > 0 .

$ac - b^2 = \det(A) = \lambda_1 \lambda_2 > 0$, $\text{Trace}(A) = \lambda_1 + \lambda_2 = a + c > 0$

→ A real symmetric $n \times n$ matrix is positive definite if:

$v^T A v > 0 \quad \forall v \in \mathbb{R}^n, v \neq 0$



WEEK 7

→ Principal Component Analysis

→ Feature Selection: start with as many features as you can and then find the good subset of features.

→ Idea behind PCA → project data onto a lower dimensional subspace s.t.

① Reconstruction error is minimised.

② Variance of projected data is maximised.

→ Problem:

Given: Dataset $D: \{x_1, \dots, x_m\}$, $x_i \in \mathbb{R}^d$

Goal: Project D onto a m -dimensional subspace.

↳ m : input parameter

PCA algorithm

① Let $B = \{u_1, \dots, u_m\}$ orthonormal basis for m -dimensional subspace

↳ fix a subspace, find the best projection

↳ how to choose subspace optimally? later...

② Extend B to an orthonormal basis for \mathbb{R}^d .

Extended basis $B' = \{u_1, \dots, u_m, u_{m+1}, \dots, u_d\}$

③ Any vector $x \in \mathbb{R}^d$ can be written using B' as:

$$x = \alpha_1 u_1 + \dots + \alpha_d u_d, \text{ where } \alpha_j = x^\top u_j \quad \forall j=1, \dots, d$$

Expressing each datapoint in $D = \{x_1, \dots, x_m\}$ using B' :

$$x_i = \sum_{j=1}^d (x_i^\top u_j) u_j \quad \leftarrow \text{original datapoint}$$

Approximate x_i by \tilde{x}_i as:

$$\tilde{x}_i = \sum_{j=1}^m z_{ij} u_j + \sum_{j=m+1}^d \beta_j u_j \quad \leftarrow \begin{array}{l} \text{projected datapoint} \\ (\text{belongs to } m\text{-dim subspace}) \end{array}$$

④ Find optimal z_{ij} and β_j to minimize square error:

$$J = \frac{1}{n} \sum_{i=1}^n \|x_i - \tilde{x}_i\|^2$$

$$J = \frac{1}{n} \sum_{i=1}^n \left\| \sum_{j=1}^m (x_i^\top u_j - z_{ij}) u_j + \sum_{j=m+1}^d (x_i^\top u_j - \beta_j) u_j \right\|^2$$

$$J = \frac{1}{n} \sum_{i=1}^n \left[\sum_{j=1}^m (x_i^\top u_j - z_{ij})^2 + \sum_{j=m+1}^d (x_i^\top u_j - \beta_j)^2 \right]$$

To minimise J , take partial derivative w.r.t. z_{ij} and β_j

$$\frac{\partial J}{\partial z_{ij}} = 0 \Rightarrow 2(x_i^\top u_j - z_{ij}) = 0 \Rightarrow z_{ij} = x_i^\top u_j$$

$$\frac{\partial J}{\partial \beta_j} = 0 \Rightarrow \frac{1}{n} \sum_{i=1}^n (x_i^\top u_j - \beta_j) = 0 \Rightarrow \beta_j = \underbrace{(\frac{1}{n} \sum_{i=1}^n x_i)^\top u_j}_{\bar{x}^\top} \rightarrow \text{mean of data}$$

→ Conclusion:

$$\tilde{x}_i = \sum_{j=1}^m (x_i^\top u_j) u_j + \sum_{j=m+1}^d (\bar{x}^\top u_j) u_j$$

$$\rightarrow u_i - \bar{u}_i = \sum_{j=m+1}^d (u_i^\top u_j - \bar{u}^\top u_i) u_j \quad \sum_{j=1}^m (u_i^\top u_j) u_j \text{ vanishes}$$

$$\|u_i - \bar{u}_i\|^2 = \sum_{j=m+1}^d (u_i^\top u_j - \bar{u}^\top u_i)^2 = \sum_{j=m+1}^d ((u_i - \bar{u})^\top u_j)^2, \text{ where } \bar{u}^\top \text{ is the mean of data}$$

With optimised z_{ij}, β_{ij} , the square error becomes:

$$\begin{aligned} J^* &= \frac{1}{n} \sum_{i=1}^n \sum_{j=m+1}^d ((u_i - \bar{u})^\top u_j)^2 \\ &= \frac{1}{n} \sum_{j=m+1}^d \sum_{i=1}^n ((u_i - \bar{u})^\top u_j) ((u_i - \bar{u})^\top u_j) \\ &= \frac{1}{n} \sum_{j=m+1}^d \sum_{i=1}^n u_j^\top (u_i - \bar{u})(u_i - \bar{u})^\top u_j \\ &= \sum_{j=m+1}^d u_j^\top \underbrace{\left[\frac{1}{n} \sum_{i=1}^n (u_i - \bar{u})(u_i - \bar{u})^\top \right] u_j}_{C u_j} \\ J^* &= \sum_{j=m+1}^d u_j^\top C u_j, \text{ where } C = \frac{1}{n} \sum_{i=1}^n (u_i - \bar{u})(u_i - \bar{u})^\top \end{aligned}$$

→ example:

Goal: want to choose u_{m+1}, \dots, u_d s.t. J^* is minimised

Constrained optimisation: $\min_u u^\top C u$ s.t. $u^\top u = 1$

Lagrangian $L(u, \lambda) = u^\top C u + \lambda(1 - u^\top u)$

↑ primal variable ↑ Lagrange multiplier

$$\nabla_u L(u, \lambda) = 0 \Rightarrow C u = \lambda u \Leftrightarrow u^\top C u = \lambda$$

$\min_u u^\top C u$ is the smallest eigenvalue λ of C .

$$C = \frac{1}{n} \sum_{i=1}^n (u_i - \bar{u})(u_i - \bar{u})^\top$$

real symmetric matrix and
all eigenvalues are real and
there exists an orthonormal basis of eigenvectors

- let the basis of eigenvectors of C be $\{u_1, \dots, u_m, u_{m+1}, \dots, u_d\}$ corresponding to eigenvalues $\{\lambda_1, \dots, \lambda_m, \lambda_{m+1}, \dots, \lambda_d\}$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$

- To minimise $J^* = \sum_{j=m+1}^d u_j^\top C u_j$:

choose $\{u_{m+1}, \dots, u_d\}$ to be $(d-m)$ eigenvectors corresponding to the $(d-m)$ least eigenvalues $\{\lambda_{m+1}, \dots, \lambda_d\}$

Remaining $\{u_1, \dots, u_m\}$ are chosen to be top- m eigenvectors of C .

→ Procedure of PCA:

① Data: $\{x_1, \dots, x_n\} \quad x_i \in \mathbb{R}^d \quad \forall i$

pick the largest m eigenvalues and their eigenvectors as a proxy for "max explainability"

essentially, the variables with the least eigenvalues in the covariance matrix are being omitted.

② Let $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ and $C = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^\top$

③ Find eigenvalues $\{\lambda_1, \dots, \lambda_d\}$ $\lambda_1 \geq \dots \geq \lambda_d$ with corresponding eigenvectors $\{u_1, \dots, u_d\}$

④ Projected data: $\tilde{x}_i = \sum_{j=1}^m (\bar{x}^\top u_j) u_j + \sum_{j=m+1}^d (\bar{x}^\top u_j) u_j$

→ example:

$$\text{dataset } D = \left\{ \begin{pmatrix} x_1 \\ -1 \end{pmatrix}, \begin{pmatrix} x_2 \\ 0 \end{pmatrix}, \begin{pmatrix} x_3 \\ 1 \end{pmatrix} \right\}$$

project D onto " $m=1$ "-dimensional subspace

$$\bar{x} = \frac{1}{3} \sum_{i=1}^3 x_i = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$C = \frac{1}{3} \sum_{i=1}^3 (x_i - \bar{x})(x_i - \bar{x})^\top = \frac{1}{3} \sum_{i=1}^3 x_i x_i^\top$$

$$= \frac{1}{3} \left[\begin{bmatrix} -1 \\ -1 \end{bmatrix} \begin{bmatrix} -1 & -1 \end{bmatrix} + 0 + \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix} \right] = \frac{2}{3} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

eigenvalues of C : $\lambda_1 = 4/3, \lambda_2 = 0$

eigenvectors of C : $u_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}, u_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix}$

Best 1-dimensional space is spanned by $u_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$$\tilde{x}_i = (x_i^\top u_1) u_1 + (\bar{x}^\top u_2) u_2 = (x_i^\top u_1) u_1$$

$$\tilde{x}_1 = \frac{1}{2} [-1 -1] \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$$

$$\tilde{x}_2 = \frac{1}{2} [0 0] \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 0$$

$$\tilde{x}_3 = \frac{1}{2} [1 1] \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Reconstruction error: $J^* = \frac{1}{3} \sum_{i=1}^3 \|x_i - \tilde{x}_i\|^2 = 0$

→ PCA as maximising variance

→ Projection of a datapoint x_i onto line along u is $(x_i^T u)u$

Dataset $D = \{x_1, \dots, x_n\}$ $x_i \in \mathbb{R}^d$, mean $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

- x_i 's projection = $(x_i^T u)u$, mean projected value = $(\bar{x}^T u)u$

- variance = $(x_i^T u - \bar{x}^T u)^2$

- Projected variance over all points: $\frac{1}{n} \sum_{i=1}^n (x_i^T u - \bar{x}^T u)^2$

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

- Goal: $\max_u u^T C u$ s.t. $u^T u = 1$

→ maximiser of $u^T C u$ is an eigenvector of C corresponding to the largest eigenvalue of C

$$\rightarrow \max_u u^T C u \text{ s.t. } u^T u = 1 \Leftrightarrow \max_u \frac{u^T C u}{u^T u}$$

let $u = (u^{(1)}, \dots, u^{(n)})$

$$\frac{\partial u^T u}{\partial u^{(i)}} = \frac{\partial}{\partial u^{(i)}} (u^{(1)^2} + \dots + u^{(n)^2}) = 2u^{(i)}$$

$$\frac{\partial u^T C u}{\partial u^{(i)}} = \frac{\partial}{\partial u^{(i)}} \left(\sum_{i=1}^n \sum_{j=1}^n C_{ij} u^{(i)} u^{(j)} \right) = 2 \sum_j C_{ij} u^{(j)} = 2(Cu)^{(i)}$$

(i,j) th entry of matrix i-th coordinate of u

$$\frac{\partial}{\partial u^{(i)}} \left(\frac{u^T C u}{u^T u} \right) = u^T u 2(Cu)^{(i)} - (u^T C u) 2u^{(i)} = 0 \quad \text{---} \textcircled{1}$$

$$\textcircled{1} \text{ in vector form: } u^T u C u = (u^T C u) u \Leftrightarrow C u = \left(\frac{u^T C u}{u^T u} \right) u$$

$$\Rightarrow C u = \lambda u$$

So, the maximiser of $\frac{u^T C u}{u^T u}$ is an eigenvector of C

and $\max_u \frac{u^T C u}{u^T u} = \lambda_1$, where λ_1 is the largest eigenvalue

→ Can be extended to cases where $m > 1$.

→ Pick the eigenvectors corresponding to the largest eigenvalues to maximise the projected variance.

→ From a "maximising variance" viewpoint, PCA picks largest- m eigenvalues of C with corresponding $\{u_1, \dots, u_m\}$ $u_i \rightarrow$ principal directions; projected values \rightarrow principal components

→ Revisiting the example:

$$D = \left\{ \begin{pmatrix} -1 \\ -1 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right\}; m=1; u_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\text{Projections} = \left\{ \begin{pmatrix} -1 \\ -1 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right\}$$

$$x_1^T u_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} -1 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = -\sqrt{2}$$

$$x_2^T u_1 = 0$$

$$x_3^T u_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \sqrt{2}$$

$$\text{Projected variance} = \frac{1}{3} ((-\sqrt{2})^2 + (\sqrt{2})^2) = \frac{4}{3}$$

is equal to the largest eigenvalue of C .

→ PCA in higher dimensions

→ Dataset $D = \{x_1, \dots, x_n\}$ $x_i \in \mathbb{R}^d$ $\forall i$

Feature dimension $d \gg$ number of datapoints n

↳ easier to handle $n \times n$ matrix than $d \times d$

PCA requires finding eigenvectors of $C = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$, $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
 ↳ $d \times d$ matrix

Goal: Reformulate the problem as finding eigenvectors of $n \times n$ matrix

Rank (C) $\leq n$

⇒ $(d-n)$ eigenvalues are zero; it is not necessary to find $(d-n)$ eigenvectors

$$\text{Let } A = \begin{bmatrix} (x_1 - \bar{x})^T \\ \vdots \\ (x_n - \bar{x})^T \end{bmatrix} \quad \text{Then, } C = \frac{1}{n} A^T A$$

Let u_i be an eigenvector of C corresponding to eigenvalue $\lambda_i > 0$

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

$$\text{Proof: } \lambda_i(A u_i) = A(\lambda_i u_i)$$

$$= A\left(\frac{1}{n} A^T A u_i\right) \rightarrow \text{since } \lambda_i \text{ is an eigenvalue of } \frac{1}{n} A^T A$$

$$\lambda_i(A u_i) = \frac{1}{n} A A^T (A u_i)$$

$$\text{i.e., } \left(\frac{1}{n} A A^T\right)(A u_i) = \lambda_i(A u_i) \Rightarrow \lambda_i \text{ is eigenvalue of } \frac{1}{n} A A^T$$

→ It is enough to find eigenvectors of $\frac{1}{n} A A^T$, where $A = \begin{bmatrix} (x_1 - \bar{x})^T \\ \vdots \\ (x_n - \bar{x})^T \end{bmatrix}$



WEEK 8

→ General form of an optimisation problem

$$\rightarrow \min_{x \in \mathbb{R}^k} f(x) ; g_i(x) \leq 0 \quad \forall i=1, \dots, k ; h_j(x) = 0 \quad \forall j=1, \dots, k$$

variable / parameter
 objective function
 inequality constraints
 equality constraints

→ Unconstrained Optimisation

$$\rightarrow \text{example} \quad \min_x (x-5)^2$$

$$f(x) = (x-5)^2 ; f'(x) = 2(x-5)$$

using: $x_{t+1} = x_t + d$, where $d = -f'(x)$ the solution will not converge to the answer. The direction is correct, magnitude is not.

update: $x_{t+1} = x_t + n_t d$, where $d = -f'(x)$

$$\text{good step size sequence} \left\{ n_t = \frac{1}{t+1} \right. \quad \left. n_t = \frac{1}{2^t} \right\} \text{bad step size sequence}$$

→ Gradient Descent Algorithm

↳ Initialise at $x_0 \in \mathbb{R}$

for $t=1, 2, \dots$

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

end

↳ Properties:

→ If $n_t = 1/(t+1)$, the algorithm converges

→ Gradient descent converges to a "local minimum"

→ Taylor Series

$$f(x+n\Delta) = f(x) + n\Delta f'(x) + \frac{n^2 \Delta^2}{2} f''(x) + \dots$$

↳ Local information gives global information

↳ for small enough n :

$$f(x+n\Delta) \approx f(x) + n\Delta f'(x)$$

$$\underbrace{f(x+n\Delta)}_{\text{function at new point}} - \underbrace{f(x)}_{\text{function at current point}} \approx n\Delta f'(x)$$

↳ minimize $f(x)$, choose direction Δ s.t. $f(x+n\Delta) - f(x) < 0$

$$\Rightarrow \text{s.t. } n\Delta f'(x) < 0 \Rightarrow \begin{matrix} \text{want} \\ \Delta \text{ s.t. } \end{matrix} \Delta f'(x) < 0$$

small
 +ve constant

→ Gradient Descent for multivariate functions

$$\hookrightarrow f(u_1, u_2) = u_1^2 + 4u_2 + 8u_2^2$$

$$\nabla f \left(\begin{bmatrix} 1 \\ 3 \end{bmatrix} \right) = \begin{bmatrix} 2u_1 \\ 4+16u_2 \end{bmatrix}_{u_1=1, u_2=3} = \begin{bmatrix} 2 \\ 52 \end{bmatrix}$$

$$\hookrightarrow \vec{u}_{t+1} = \vec{u}_t + n(-\nabla f(\vec{u}_t))$$

→ Taylor Series for multivariable Functions

$$\hookrightarrow f(u + nd) = f(u) + nd^\top \nabla f(u) + \dots$$

$$\hookrightarrow \text{want } d \text{ s.t. } f(u + nd) - f(u) < 0$$

$$\Rightarrow \text{want } d \text{ st. } d^\top \nabla f(u) < 0$$

if $d = -\nabla f(u)$, then $\underbrace{\|d\|}_{d^\top \nabla f(u) = -\|\nabla f(u)\|^2}$ norm is always positive



WEEK 9

→ Constrained Optimisation

$$\rightarrow \min_u f(u) \text{ s.t. } g(u) \leq 0 \quad \text{--- ①}$$

Given u^* , it is claimed to solve ① (i.e., u^* is optimal). How to check the claim?

→ check constraint $\rightarrow g(u^*) \leq 0$

→ No "descent direction" should be a "feasible direction"

any direction
that reduces $f(u)$

satisfies
 $g(u) \leq 0$ any direction that
takes to a point that
is feasible

→ If there is a direction from u^* which is a descent direction for both $f(u)$ and $g(u)$, i.e. $d^\top \nabla f(u^*) < 0$ and $d^\top \nabla g(u^*) < 0$, then u^* is not optimal because there is a direction that decreases $f(u)$ and also satisfies $g(u) \leq 0$.

→ Necessary condition for u^* to be optimal: $\nabla f(u^*)$ and $\nabla g(u^*)$ should be antiparallel

$$\hookrightarrow \boxed{\nabla f(u^*) = -\lambda \nabla g(u^*)} \quad \text{where } \lambda \text{ is a positive scalar}$$

↑
Lagrange multiplier

→ Method of Lagrange Multiplier

$$\rightarrow \min_{\mathbf{x}} f(\mathbf{x}) \quad g(\mathbf{x}) = 0$$

Method:

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

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

→ example:

$$f(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^2 + 2\mathbf{x}_2 + 4\mathbf{x}_2^2$$

$$g(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^2 + \mathbf{x}_2^2 - 1$$

$$\nabla f(\mathbf{x}_1, \mathbf{x}_2) = (2\mathbf{x}_1, 2+8\mathbf{x}_2)$$

$$\nabla g(\mathbf{x}_1, \mathbf{x}_2) = (2\mathbf{x}_1, 2\mathbf{x}_2)$$

Using method of lagrangian Multiplier:

$$\begin{bmatrix} 2\mathbf{x}_1 \\ 2+8\mathbf{x}_2 \end{bmatrix} = -\lambda \begin{bmatrix} 2\mathbf{x}_1 \\ 2\mathbf{x}_2 \end{bmatrix} \Rightarrow \begin{aligned} 2\mathbf{x}_1 &= -\lambda 2\mathbf{x}_1 \\ 2+8\mathbf{x}_2 &= -\lambda 2\mathbf{x}_2 \end{aligned} \quad \textcircled{1} \quad \textcircled{2}$$

$$\textcircled{1} \Rightarrow 2\mathbf{x}_1(1+\lambda) = 0 \rightarrow \text{either } \lambda = -1 \text{ or } \mathbf{x}_1 = 0$$

If $\lambda = -1$, using $\textcircled{2}$, $\mathbf{x}_2 = -\frac{1}{3}$

$$\mathbf{x}_1 = \sqrt{1-\mathbf{x}_2^2} = \pm\sqrt{\frac{8}{3}} \text{ or } -\sqrt{\frac{8}{3}}$$

$$\left\{ \begin{bmatrix} -\sqrt{\frac{8}{3}} \\ -\frac{1}{3} \end{bmatrix}, \begin{bmatrix} \sqrt{\frac{8}{3}} \\ -\frac{1}{3} \end{bmatrix} \right\}$$

If $\mathbf{x}_1 = 0$, $\mathbf{x}_2 = 1$ or -1

$$\left\{ \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \end{bmatrix} \right\}$$

$$f(0, 1) = 6, f(0, -1) = 2, f\left(\sqrt{\frac{8}{3}}, -\frac{1}{3}\right) = \frac{2}{3}, f\left(-\sqrt{\frac{8}{3}}, -\frac{1}{3}\right) = \frac{2}{3}$$

maximizer

minimizer

→ Projected Gradient Descent

→ Gradient descent algorithm worked for unconstrained optimisation, but in the case of constrained optimisation, the algorithm might go outside the scope of constraints.

→ Solution: after every step, take a projection onto feasible set.

→ Convex Sets

→ A set $S \subseteq \mathbb{R}^d$ is a convex set if $\forall \mathbf{x}_1, \mathbf{x}_2 \in S$.

$$\text{then } (\lambda \mathbf{x}_1 + (1-\lambda) \mathbf{x}_2) \in S \quad \forall \lambda \in [0, 1]$$

→ Properties:

① Intersection of convex sets is convex.

↪ System of linear equations is a convex set

② Convex combinations: let $S = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \subseteq \mathbb{R}^d$. $\mathbf{z} \in \mathbb{R}^d$ is a convex combination of points in S if $\exists \lambda_1, \dots, \lambda_n$ s.t. $\lambda_i \geq 0, \sum \lambda_i = 1$

$$\mathbf{z} = \lambda_1 \mathbf{x}_1 + \lambda_2 \mathbf{x}_2 + \dots + \lambda_n \mathbf{x}_n$$

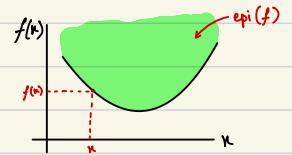
$$\text{Convex-Hull}(S) = \left\{ \mathbf{z} : \mathbf{z} = \sum_{i=1}^n \lambda_i \mathbf{x}_i \text{ for some } \lambda_1, \dots, \lambda_n \geq 0 \text{ s.t. } \sum \lambda_i = 1 \right\}$$

Convex Functions

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

\hookrightarrow any convex set

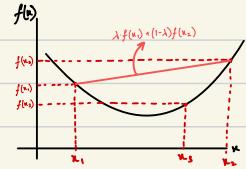
$$\overbrace{\text{epi}(f)}^{\subseteq \mathbb{R}^{d+1}} = \left\{ \begin{bmatrix} x \\ z \end{bmatrix} \in \mathbb{R}^{d+1} : z \geq f(x) \right\}$$



① **Definition:** A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function
if $\text{epi}(f)$ is a convex set

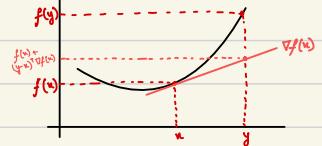
\rightarrow For any point x_3 , $\lambda f(x_1) + (1-\lambda)f(x_2) > f(x_3)$

$$\text{where } x_3 = \lambda x_1 + (1-\lambda)x_2$$



② **Definition:** A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function
iff $\forall x_1, x_2 \in \mathbb{R}^d$ and all $\lambda \in [0, 1]$,
 $f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2)$

\rightarrow Assume $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is differentiable
 f is always greater than the linear approximation.



③ **Definition:** A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function
iff $f(y) \geq f(x) + (y-x)^T \nabla f(x)$

\rightarrow Assume $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is twice differentiable
 f always has a non-negative curvature

④ **Definition:** A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function

if matrix H is a positive semi-definite matrix
where $H \in \mathbb{R}^{d \times d}$, $H_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$. \rightarrow all eigenvalues(H) ≥ 0



WEEK 10

Properties of Convex Functions

→ If f is a convex function, then all local minima of f are also global minima.

→ Necessary conditions for optimality of convex functions
 $\hookrightarrow \min_u f(u)$

$x^* \in \mathbb{R}^n$ is a global minimum of f iff $\nabla f(x^*) = 0$

Proof: if $\nabla f(x^*) \neq 0$, then f can be lowered by moving in direction $-\nabla f(x^*)$

→ If $\exists u^* \text{ s.t. } \nabla f(u^*) = 0 \Rightarrow u^*$ is a global minimum

Proof:

By definition of convexity:

$$\begin{aligned} f(y) &\geq f(u) + \nabla f(u)^T (y - u) \quad \forall u, y \\ \Rightarrow f(y) &\geq f(u^*) + \underbrace{\nabla f(u^*)^T (y - u^*)}_{=0} \quad \forall y \end{aligned}$$

$\Rightarrow f(y) \geq f(u^*) \quad \forall y \Rightarrow u^*$ is a global min

Additional properties of Convex Functions

① If $f: \mathbb{R}^n \rightarrow \mathbb{R}$ and $g: \mathbb{R}^n \rightarrow \mathbb{R}$ are both convex functions, then $h(u) := f(u) + g(u)$ is also a convex function.

→ Sum of convex functions is a convex function

Proof:

$$h(\lambda u + (1-\lambda)y) = f(\lambda u + (1-\lambda)y) + g(\lambda u + (1-\lambda)y) \quad \lambda \in [0, 1]$$

$$h(\lambda u + (1-\lambda)y) \leq (\lambda f(u) + (1-\lambda)f(y)) + (\lambda g(u) + (1-\lambda)g(y))$$

$$h(\lambda u + (1-\lambda)y) \leq \lambda(f(u) + g(u)) + (1-\lambda)(f(u) + g(y))$$

$h(\lambda u + (1-\lambda)y) \leq \lambda h(u) + (1-\lambda)h(y) \longrightarrow$ proves convexity of $h(u)$

② Composition of functions

Let $f: \mathbb{R} \rightarrow \mathbb{R}$ be a convex non-decreasing function i.e., $f(u) \geq f(y) \quad \forall u \geq y$

Let $g: \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function

Define $h := f \circ g \quad h(u) = f(g(u)) \quad h: \mathbb{R}^n \rightarrow \mathbb{R}$

→ Composition of convex with convex + non-decreasing is a convex function

③ Let $f: \mathbb{R} \rightarrow \mathbb{R}$ be convex

Let $g: \mathbb{R}^n \rightarrow \mathbb{R}$ be linear

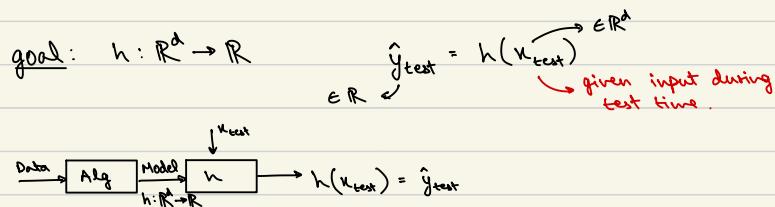
Define $h := f \circ g$

→ Composition of linear with convex is a convex function

→ Application of optimization in ML

→ Linear Regression:

$$\text{Dataset: } \begin{cases} x_1, x_2, \dots, x_n \\ y_1, y_2, \dots, y_n \end{cases} \quad \begin{array}{l} x_i \in \mathbb{R}^d \\ y_i \in \mathbb{R} \end{array} \quad \text{Training data}$$



- In this case, h is a linear function

$$h(x) = w^T x \quad \text{for some } w \in \mathbb{R}^d$$

- The "best" $h(x)$ from training data is characterised by

good the w is Performance measure: sum of Square Error

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

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

$$f(w) = \sum_{i=1}^n \underbrace{(w^T x_i - y_i)^2}_{h_i(w)} \quad \text{Is } f(w) \text{ convex?}$$

$\sum_{i=1}^n h_i(w)$ sum of convex functions is convex

If $h_i(w)$ is convex $\forall i$, then $f(w)$ is convex

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

$$= f(g(w)) \quad \text{where } g(w) = w^T x_i - y_i \quad g: \mathbb{R}^d \rightarrow \mathbb{R} \quad \leftarrow \text{linear}$$

$$\Rightarrow \text{convexity} \quad f(z) = z^2 \quad f: \mathbb{R} \rightarrow \mathbb{R} \quad \leftarrow \text{convex}$$

Conclusion: SSE is a convex function

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

$$= \frac{1}{2} \left\| \begin{bmatrix} (w^T x_1 - y_1) \\ \vdots \\ (w^T x_n - y_n) \end{bmatrix} \right\|_2^2 \quad \text{norm}$$

$$f(w) = \frac{1}{2} \|xw - y\|_2^2 \quad \text{where} \quad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \quad w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

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

$$= \frac{1}{2} (w^T x^T x - 2w^T x^T y + y^T y)$$

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

The global minimum at w^* : $(x^T x)w^* = (x^T y)$

$$w^* = (x^T x)^{-1}(x^T y)$$

Issue with $(x^T x)^{-1}(x^T y)$: needs an inverse computation $\sim O(d^3)$

Gradient descent can be done without inverse; iteratively

$$\hookrightarrow w_{t+1} = w_t - \eta_t \nabla f(w_t) \quad \xrightarrow{(x^T x)w - x^T y}$$

↳ can be made faster by using
Stochastic Gradient Descent

↳ samples a small set of data points uniformly at random

↳ pretend sampled points form the new dataset
and compute gradient w.r.t. w

→ Constrained Optimisation

$$\rightarrow \min_{\mathbf{x}} f(\mathbf{x}) \text{ s.t. } h(\mathbf{x}) \leq 0 \quad f: \mathbb{R}^n \rightarrow \mathbb{R}$$

→ Lagrangian Function

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

vector scalar

$$\text{Fix } \mathbf{x} \in \mathbb{R}^n \quad \max_{\lambda \geq 0} L(\mathbf{x}, \lambda) = \max_{\lambda \geq 0} f(\mathbf{x}) + \lambda h(\mathbf{x})$$

↪ if $h(\mathbf{x}) \leq 0$, $L(\mathbf{x}, \lambda)$ at max will be $f(\mathbf{x})$

if $h(\mathbf{x}) > 0$, $L(\mathbf{x}, \lambda)$ at max will be ∞

→ By maximising $L(\mathbf{x}, \lambda)$, the solution is only possible if the condition of $h(\mathbf{x}) \leq 0$ is met.

Therefore,

$$\boxed{\begin{aligned} \min_{\mathbf{x}} f(\mathbf{x}) \\ \text{s.t. } h(\mathbf{x}) \leq 0 \end{aligned}} \quad \equiv \quad \min_{\mathbf{x}} \left[\max_{\lambda \geq 0} L(\mathbf{x}, \lambda) \right] \quad \text{where } L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda h(\mathbf{x})$$

$$\rightarrow \min_{\mathbf{x}} \left[\max_{\lambda \geq 0} L(\mathbf{x}, \lambda) \right] \quad \text{Primal problem} \quad \mathbf{x}^* \rightarrow \text{Primal solution}$$

$$\max_{\lambda \geq 0} \left[\min_{\mathbf{x}} f(\mathbf{x}) + \lambda h(\mathbf{x}) \right] \quad \text{Dual Problem} \quad = \quad \max_{\lambda \geq 0} g(\lambda) \quad \text{, where } g(\lambda) = \min_{\mathbf{x}} L(\mathbf{x}, \lambda)$$

Unconstrained Optimisation

Concave Function

→ Relation b/w Primal and Dual problem

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

$$\text{Define } J(\mathbf{x}) = \begin{cases} f(\mathbf{x}) & \text{if } h(\mathbf{x}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$$

→ For any $\lambda \geq 0$, $L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda h(\mathbf{x})$

↪ $L(\mathbf{x}, \lambda) \leq f(\mathbf{x})$ if $h(\mathbf{x}) \leq 0$ and $\leq \infty$ otherwise

→ Fix $\lambda \geq 0$,

$$L(\mathbf{x}, \lambda) \leq J(\mathbf{x}) \quad \forall \mathbf{x}$$

$$\min_{\mathbf{x}} L(\mathbf{x}, \lambda) \leq \min_{\mathbf{x}} J(\mathbf{x}) = f(\mathbf{x}^*)$$

Solution of the
Primal Problem

$$\max_{\lambda \geq 0} \min_{\mathbf{x}} L(\mathbf{x}, \lambda) \leq f(\mathbf{x}^*)$$

→ The value at Dual optimum \leq value at primal optimum ← **Weak-Duality**

→ If f, h are convex, then **Strong-Duality holds**

objective constraint

$$\rightarrow \min_{\mathbf{x}} \left[\max_{\lambda \geq 0} L(\mathbf{x}, \lambda) \right] \quad (\text{or}) \quad \max_{\lambda \geq 0} \left[\underbrace{\min_{\mathbf{x}} L(\mathbf{x}, \lambda)}_{\text{easier to solve}} \right]$$

→ KKT conditions

→ Assume f, h are convex \Rightarrow Strong Duality holds
let x^*, λ^* are the primal and dual optimal solutions

$$x^* = \underset{u}{\operatorname{argmin}} \left[\max_{\lambda \geq 0} f(u) + \lambda h(u) \right]$$

$$\lambda^* = \underset{\lambda \geq 0}{\operatorname{argmax}} \left[\min_u f(u) + \lambda h(u) \right]$$

By strong duality,

$$f(x^*) = g(\lambda^*) \Rightarrow \nabla f(x^*) + \lambda^* \nabla h(x^*) = 0 \\ = \min_u f(u) + \lambda^* h(u) \leftarrow \text{this is minimum over all possible } u$$

$$f(x^*) \leq f(u^*) + \underbrace{\lambda^* h(u^*)}_{\leq 0} \leq f(u^*)$$

$$\Rightarrow \lambda^* h(u^*) = 0$$

→ If f, h are convex \Rightarrow strong duality
 $\Rightarrow x^*, \lambda^*$ must satisfy:

$$\textcircled{1} \quad \nabla f(x^*) + \lambda^* \nabla h(x^*) = 0 \quad \text{stationarity condition}$$

$$\textcircled{2} \quad \lambda^* h(x^*) = 0 \quad \text{complementary slackness condition}$$

$$\textcircled{3} \quad h(x^*) \leq 0 \quad \text{primal feasibility}$$

$$\textcircled{4} \quad \lambda^* \geq 0 \quad \text{dual feasibility}$$

→ In general, if (x^*, λ^*) satisfies the above conditions $\Rightarrow x^*$ is a local optima

→ Optimisation problem with multiple constraints:

$$\min_u f(u) \text{ s.t. } h_i(u) \leq 0 \quad \forall i=1, \dots, m, \quad l_j(u) = 0 \quad \forall j=1, \dots, n$$

$$L(u, u_v, v) = f(u) + \sum_{i=1}^m u_i h_i(u) + \sum_{j=1}^n v_j l_j(u) \quad \text{where } u = \begin{bmatrix} u_1 \\ \vdots \\ u_m \end{bmatrix}, \quad v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$$

KKT conditions:

$$\textcircled{1} \quad \nabla f(x^*) + \sum_{i=1}^m u_i^* \nabla h_i(x^*) + \sum_{j=1}^n v_j^* \nabla l_j(x^*) = 0$$

$$\textcircled{2} \quad u_i^* h_i(x^*) = 0 \quad \forall i$$

$$\textcircled{3} \quad h_i(x^*) \leq 0 \quad \forall i; \quad l_j(x^*) = 0 \quad \forall j$$

$$\textcircled{4} \quad u_i^* \geq 0 \quad \forall i$$

→ Example ML algorithm: Support Vector Machine (SVM)

→ Optimisation problem

$$\min_w \frac{1}{2} \|w\|^2 \quad \text{objective}$$

Dataset: $\{(u_i, y_i), \dots, (u_n, y_n)\}$

\curvearrowright Quadratic \Rightarrow convex. $\|w\|^2 = \sum_i w_i^2$

$$\text{s.t. } w^T u_i y_i \geq 1 \quad \forall i \quad \text{constraints}$$

\curvearrowright Linear constraints \Rightarrow convex

→ Tutorial on Linear Programming

→ Example:

$$\begin{array}{ccc} \text{Primal Problem} & \longleftrightarrow & \text{Dual Problem} \\ \min 5u_1 + 8u_2 & & \max 5y_1 + 10y_2 + 8y_3 \\ \text{s.t. } \begin{cases} 3u_1 + 0u_2 \geq 6 \\ 2u_1 + 4u_2 \geq 10 \\ 2u_1 + 5u_2 \geq 8 \\ u_1, u_2 \geq 0 \end{cases} & & \begin{cases} 3y_1 + 2y_2 + 2y_3 \leq 50 \\ 0y_1 + 4y_2 + 5y_3 \leq 80 \\ y_1, y_2, y_3 \geq 0 \end{cases} \end{array}$$

WEEK 11

→ Continuous Random Variable

→ $X: \Omega \rightarrow \mathbb{R}$ probability function X maps the sample space Ω to \mathbb{R}

↳ Domain and Range are uncountable

→ PDF: $f_x(x) = \frac{P(X \in [x, x+dx])}{dx}; P(X=x)=0$

→ CDF: $F_x(x) = P(X \leq x)$

→ Properties:

$$① f_x(x) \geq 0$$

$$② \int_{-\infty}^{\infty} f_x(x) dx = 1$$

$$③ F_x(-\infty) = 0$$

$$④ F_x(\infty) = 1$$

→ Conditional PDF

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

→ Function of random variables

example:

$$f_x(u) = y_2 \text{ if } u \in [-1, 1]$$

$$y = x/2 \quad f_y(y) = \frac{P(Y \in [y, y+dy])}{dy} = \frac{P(X \in [2y, 2y+2dy])}{dy}$$

= 0 if $2y \notin [-1, 1]$; $2 \cdot \frac{1}{2}$ if $y \in [-\frac{1}{2}, \frac{1}{2}]$

$$X \sim \text{Uniform } [-1, 1] \quad y = x/2 \in [-\frac{1}{2}, \frac{1}{2}]$$

$$F_y(y) = P(Y \leq y) = P(X \leq 2y)$$

$$P(X \leq 2y) = \int_{-1}^{2y} \frac{1}{2} du = \left. u/2 \right|_{-1}^{2y} = y + \frac{1}{2}$$

$$F_y(y) = y + \frac{1}{2} \quad f_y(y) = 1$$

→ Expectation

$$X: \Omega \rightarrow \mathbb{R}$$

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

Properties of expectation:

$$\textcircled{1} \quad E[X+Y] = E[X] + E[Y]$$

$$\textcircled{2} \quad Y = g(X); \quad E[Y] = \int_{-\infty}^{\infty} g(x) f_X(x) dx$$

→ Variance

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

→ Conditional Expectation

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

→ Multiple Random Variables

→ Joint Functions

$$f_{XY}(x, y) = \frac{P(X \in [x, x+dx], Y \in [y, y+dy])}{dx dy}$$

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

→ Marginal

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

→ Conditional

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



WEEK 12

→ Bivariate and multivariate normal

→ Standard Normal Vector

$$Z_1 \sim N(0, 1), \dots, Z_d \sim N(0, 1) \quad Z = \begin{bmatrix} Z_1 \\ \vdots \\ Z_d \end{bmatrix}$$

$$f_Z(z) = \prod_{i=1}^d \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2} z_i^2) = \frac{1}{(2\pi)^{d/2}} \exp(-\frac{1}{2} \|z\|^2)$$

→ Super Linear Transform of 2D-Normal

$$X_1 = Z_1; \quad X_2 = \rho Z_1 + \sqrt{1-\rho^2} Z_2, \text{ where } \rho \in [-1, 1]$$

$$X = \begin{bmatrix} 1 & 0 \\ \rho & \sqrt{1-\rho^2} \end{bmatrix} Z, \text{ where } X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \text{ and } Z = \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix}$$

$$\rightarrow A \quad A^{-1} = \begin{bmatrix} 1 & 0 \\ -\frac{\rho}{\sqrt{1-\rho^2}} & \frac{1}{\sqrt{1-\rho^2}} \end{bmatrix} \quad X = AZ \quad Z = A^{-1}X$$

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

$$E[X_1] = E[X_2] = 0$$

$$\begin{aligned} \text{Cov}[X_1, X_2] &= E[X_1 X_2] - E[X_1] \cdot E[X_2] \\ &= E[X_1 X_2] = E[\rho Z_1^2 + \sqrt{1-\rho^2} Z_1 Z_2] \\ &= \rho \end{aligned}$$

$$\text{Var}[X_1] = 1; \quad \text{Var}[X_2] = 1$$

$$\text{Cov}[X] = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} = AA^\top$$

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

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

$$\begin{aligned} f_X(u) &= f_Z(A^{-1}u) \cdot |\det(A^{-1})| \\ &= \frac{1}{\sqrt{1-\rho^2}} \cdot \frac{1}{2\pi} \exp(-\frac{1}{2} \|A^{-1}u\|^2) \end{aligned}$$

$$\left. \begin{aligned} &= \frac{1}{2\pi \sqrt{1-\rho^2}} \exp(-\frac{1}{2} u^\top A^{-1} A^{-1} u) \quad A^{-1} A^{-1} = \Sigma^{-1} \\ &\text{simplifies to} \\ &= \frac{1}{2\pi \sqrt{1-\rho^2}} \exp\left(-\frac{1}{2} \cdot \frac{1}{1-\rho^2} [u_1 u_2] \begin{bmatrix} 1 & -\rho \\ -\rho & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}\right) \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} u_1^2\right) \cdot \frac{1}{\sqrt{2\pi \sqrt{1-\rho^2}}} \exp\left(-\frac{1}{2(1-\rho^2)} (u_2 - \rho u_1)^2\right) \end{aligned} \right.$$

$$f_X(u) = f_{X_1}(u_1) \cdot f_{X_2|X_1}(u_2 | u_1)$$

→ Bivariate Normal

$$X = AZ$$

$$f_X(u) = f_Z(A^{-1}u) \cdot |\det(A^{-1})|$$

$$\Sigma = AA^T = E[X X^T] = \text{Cov}[X]$$

$$\det(A^{-1}) = \frac{1}{\det(A)} = \frac{1}{\sqrt{\det(\Sigma)}}$$

$$f_X(u) = \frac{1}{2\pi\sqrt{\det(\Sigma)}} \exp\left(-\frac{1}{2} u^T \Sigma^{-1} u\right)$$

Only positive definite matrices can be written as Σ

$$\Sigma = \begin{bmatrix} a^2 & pab \\ pab & b^2 \end{bmatrix} \quad p \in [-1, 1]$$

→ Multivariate Normal

$$X = AZ + \mu$$

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

$$E[X] = \mu \quad \text{Cov}[X] = \Sigma$$

→ Properties of multivariate normal

$$X \sim N(\mu, \Sigma)$$

$$\textcircled{1} \quad Y = a^T X \Rightarrow Y \sim N(a^T \mu, a^T \Sigma a)$$

$$\textcircled{2} \quad Y = B \cdot X \Rightarrow Y \sim N(B \cdot \mu, B \Sigma B^T)$$

$$\textcircled{3} \quad X_i, X_j \text{ are ind} \Leftrightarrow \text{Cov}[X_i, X_j] = 0$$

only if they happen to be coordinates of multivariate random vector

→ Estimation of parameters using ML

→ $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ → family of distributions parameterised by θ , which takes values from Θ space.

$\underbrace{X_1, X_2, \dots, X_n}$ drawn iid from P_θ for $\theta \in \Theta$
use this data to estimate θ .

$$\rightarrow L(\theta) = P(X_1 = x_1, \dots, X_n = x_n | \theta)$$

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

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

$$R(\theta) = -L(\theta)$$

$$\log(L(\theta)) = \sum_{i=1}^n \log(P_\theta(x_i)) \leftarrow \text{maximising log likelihood} \Leftrightarrow \text{maximising likelihood}$$

→ example: $\mathcal{P} = \{\text{Bernoulli}(\theta) : \theta \in [0, 1]\}$

$$X_1, X_2, \dots, X_n \in \{0, 1\}$$

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

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

$$= -\sum_{i=1}^n u_i \log(\theta) + (1-u_i) \log(1-\theta) = a \log(\frac{1}{\theta}) + (n-a) \log(\frac{1}{1-\theta})$$

$$\hat{\theta}_{ML} = \frac{a}{n}$$

$$a = \sum_{i=1}^n u_i$$

→ Gaussian Mixture Models and Expectation Maximisation

$$\rightarrow f_x(u) = \sum_{k=1}^K \pi_k N(u_n | \mu_k, \Sigma_k), \text{ where } K = \# \text{ of components}$$

π_k = what % of data is from k-th component

→ ML Estimation

Assume K is known.

Data = $\{x_1, x_2, \dots, x_N\}$, $x_i \rightarrow d$ -dimensional vector

$$P(\text{Data} | \pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K) = \prod_{n=1}^N \left(\sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k) \right)$$

Goal: Find $\max_{\pi_1, \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K}$

$$R(\theta) = \sum_{n=1}^N -\log \left(\sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k) \right)$$

not easy to solve because of summation
inside the log function

→ If the 'clustering' of data points were given, estimating parameters for each 'cluster' is straightforward. If the parameters of each cluster were given, assigning datapoints to clusters is also possible

→ Cluster Responsibilities

→ Start by fixing the parameters $\{\pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K\}$ to some random numbers. This gives the joint density of the entire space.

→ Calculate $P(z=k|x_n)$, i.e., for the n-th datapoint what is the probability that it belongs to the k-th component/cluster.

$$P(z=k|x_n) = \frac{P(x_n|z=k) P(z=k)}{P(x_n)} = \frac{P(x_n|z=k) P(z=k)}{\sum_{j=1}^K P(x_n|z=j) P(z=j)}$$

$$\gamma(z_{nk}) = \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_n | \mu_j, \Sigma_j)}$$

Probability that n-th datapoint belongs to k-th cluster

$$\mu_k = \frac{\sum_{n=1}^N \gamma(z_{nk}) x_n}{N_k}, \text{ where } \gamma(z_{nk}) = \begin{cases} 1 & \text{if } n\text{-th datapoint is assigned to } k\text{-th cluster} \\ 0 & \text{otherwise} \end{cases}$$

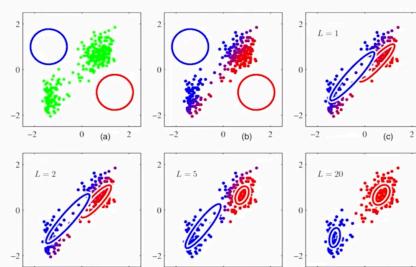
N_k # of points assigned to k-th cluster

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk})(x_n - \mu_k)(x_n - \mu_k)^T$$

These parameters are 'tamer' than the random parameters. With this, you can iteratively come closer to an increasingly less random parameters.

Illustration

$K=2$



→ Markov Inequality

→ X is a positive random variable.

$$P(X \geq t) \leq \frac{E[X]}{t}$$

→ Chebyshen Inequality

→ $E[X] = \mu$; $\text{var}[X] = \sigma^2$

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

→ Hoeffding Inequality

→ X_1, X_2, \dots, X_n are i.i.d.

$$E[X_i] = \mu \quad a \leq X_i \leq b$$

$$\bar{X}_n = \frac{1}{n} \sum_i X_i \quad E[\bar{X}] = \mu \quad \text{Var}[\bar{X}_n] = \frac{\sigma^2}{n}$$

$$P(|\bar{X}_n - \mu| \geq \delta) \leq \frac{\sigma^2}{n\delta^2}$$

$$P(|\bar{X}_n - \mu| \geq \delta) \leq \frac{\sigma^2}{n\delta^2} \longrightarrow \text{Chebyshen's Inequality}$$

$$P(|\bar{X}_n - \mu| \geq \delta) \leq 2 \exp\left(\frac{-2nt^2}{(b-a)^2}\right)$$

