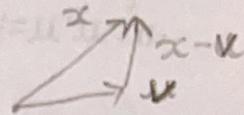


## Assignment - 04

05/03/2025

1. Let ~~data~~  $\{ \underbrace{x^{(1)}, x^{(2)}, \dots, x^{(n)} }_{\text{zero mean and unit variance}} \}$



(a) in each coordinate.

define :  $f_u(x) = \arg \min_{u \in \mathcal{Y}} \|x - v\|^2$        $\mathcal{Y} = \{\alpha u : \alpha \in \mathbb{R}\}$

$$= (x^T u) u$$

Problem statement : project our data ( $d$ -dimensional) to  $k$  dimension ( $k=1$ , here), we can formulate it as a minimization problem, (the direction ~~in~~ which gives best representation of our data)  $\min_{u: u^T u = 1} \sum_{i=1}^n \|x^{(i)} - f_u(x^{(i)})\|_2^2$

take,

$$\begin{aligned} \sum_{i=1}^n \|x^{(i)} - f_u(x^{(i)})\|_2^2 &= \sum_{i=1}^n \|x^{(i)} - ((x^{(i)})^T u) u\|_2^2 \\ &= \sum_{i=1}^n (x^{(i)} - (x^{(i)\top} u) u)^T (x^{(i)} - (x^{(i)\top} u) u) \\ &= \sum_{i=1}^n \{(x^{(i)\top} x^{(i)} - ((x^{(i)})^T u)^2 - ((x^{(i)})^T u) + (x^{(i)\top} u)^2 u^T u\} \\ &= \sum_{i=1}^n \{(x^{(i)\top} x^{(i)} - 2((x^{(i)})^T u)^2 + (x^{(i)\top} u)^2 u^T u\} \\ &= \sum_{i=1}^n (x^{(i)\top} x^{(i)} - (2 - u^T u)((x^{(i)})^T u)^2) \end{aligned}$$

---


$$\arg \min_{u: u^T u = 1} \sum_{i=1}^n \|x^{(i)} - f_u(x^{(i)})\|_2^2$$

simplified expression.

$$\begin{aligned} &\Leftrightarrow \arg \min_{u: u^T u = 1} \sum_{i=1}^n (x^{(i)\top} x^{(i)} - (2 - u^T u)((x^{(i)})^T u)^2) \\ &\Leftrightarrow \arg \max_{u: u^T u = 1} \sum_{i=1}^n ((x^{(i)\top} u))^2 \quad (\because 2 - u^T u \geq 0 \text{ and } \|u\|^2 = 1) \end{aligned}$$

$$\underset{\|u\|: \|u^T u\|=1}{\arg \max} \sum_{i=1}^n ((x^i)^T u)^2 \Leftrightarrow \underset{\|u\|: \|u^T u\|=1}{\arg \max} \sum_{i=1}^n \{u^T x^i (x^i)^T u\}$$

$$\Leftrightarrow \underset{\|u\|: \|u^T u\|=1}{\arg \max} u^T \left( \sum_{i=1}^n (x^i)(x^i)^T \right) u$$

but we know that,

$$\begin{bmatrix} \uparrow & \uparrow & \dots & \uparrow \\ x^1 & x^2 & \dots & x^n \\ \downarrow & \downarrow & & \downarrow \\ \underbrace{x^T} & & & x \end{bmatrix} \begin{bmatrix} \leftarrow x^{1T} \\ \leftarrow x^{2T} \\ \vdots \\ \leftarrow x^{nT} \end{bmatrix} = \sum_{i=1}^n (x^i)(x^i)^T$$

$$\therefore \Leftrightarrow \underset{\|u\|: \|u^T u\|=1}{\arg \max} u^T (x^T x) u$$

~~Q.E.D.~~ Let  $A = x^T x$  (symmetric and real)

$$\therefore \exists Q \in \mathbb{R}^{d \times d}, \Lambda \in \mathbb{R}^{d \times d} \text{ s.t.}$$

$$A = Q \Lambda Q^T \quad \text{s.t. } Q = \text{orthogonal matrix}$$

$\Lambda = \text{diagonal matrix}$

$$\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$$

$\lambda_1 \geq \lambda_2 \geq \lambda_3 \dots \geq \lambda_d \rightarrow \text{eigen vectors}$

$\therefore \text{we get}$

$$\Leftrightarrow \underset{\|u\|: \|u^T u\|=1}{\arg \max} u^T (Q \Lambda Q^T) u$$

$$u: u^T u=1$$

column vectors of matrix  $Q$ .

For simplification,

$$u^T (Q \Lambda Q^T) u, \quad \text{let } u = \sum_{j=1}^d \alpha_j q_j \quad (\text{$Q$ is invertible so form basis for } \mathbb{R}^d)$$

$$= u^T (Q \Lambda Q^T) \sum_{j=1}^d \alpha_j q_j = u^T Q \Lambda \sum_{j=1}^d \alpha_j e_j$$

$$= u^T Q \sum_{j=1}^d \alpha_j \lambda_j e_j = u^T \sum_{j=1}^d \alpha_j \lambda_j q_j$$

$$= \sum_{j=1}^d \{ \alpha_j \lambda_j (q_j \sum_{i=1}^d \alpha_i q_i) \} = \sum_{j=1}^d \alpha_j \lambda_j q_j^T q_j$$

$$= \sum_{j=1}^d (\alpha_j)^2 \lambda_j$$

Now minimisation problem becomes  
(initial)

$$\arg \max \sum_{j=1}^d (\alpha_j)^2 \lambda_j \Leftrightarrow \arg \max_{\substack{\sum \beta_i = 1 \\ \beta_i \geq 0}} \sum_{j=1}^d (\beta_j) \lambda_j$$

$\|\alpha\|_2^2 = 1$

$\beta = (\beta_1, \dots, \beta_d)$  maximising happens at one of  
the corner points of the region given by

$$\beta_{\text{region}} = \left\{ \beta : \sum_{i=1}^d \beta_i = 1, \beta_i \geq 0 \quad i=1 \dots d \right\}$$

∴ the soln of the optimisation problem is

$$\beta = (1, 0, 0, \dots, 0) \text{ i.e. } \beta_1 = 1, \beta_2 = \beta_3 = \dots = \beta_d = 0 \quad \therefore \lambda_1 \geq \lambda_k \quad (k=2 \dots d)$$

$$\Leftrightarrow \alpha = (1, 0, \dots, 0)$$

$$\Leftrightarrow u = v_1 \text{ (eigen vector corresponding to } \lambda_1 \text{)}$$

1 (b) (i)  $M \in \mathbb{R}^{p \times q}$

then  $M^T M \in \mathbb{R}^{q \times q}$   $\therefore$  Real matrix  
and square too

NOW,

$$(M^T M)^T = M^T (M^T)^T = M^T M \quad \therefore M^T M \text{ is symmetric too}$$

$$\text{QED } M^T M = Q \Lambda Q^T \quad \Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_q)$$

s.t.  $\lambda_1 \geq \lambda_2 \dots \geq \lambda_q \in \mathbb{R}$

$Q \in \mathbb{R}^{q \times q}$  and orthogonal matrix

(ii)  $M \in \mathbb{R}^{p \times q}$  (given)

$M = U \Sigma V^T$  where  $U \in \mathbb{R}^{p \times p}$ ,  $V \in \mathbb{R}^{q \times q}$ ,

where  $U^T U = I$ ,  $V^T V = I$

$\Sigma \in \mathbb{R}^{p \times q}$  and diagonal matrix.

$$M^T M = (U \Sigma V^T)^T (U \Sigma V^T)$$

$$= (V^T \Sigma^T U^T)(U \Sigma V^T)$$

$$= V^T \Sigma^T (U^T U) \Sigma V^T$$

$$= V^T \Sigma^T \Sigma V^T = V^T \Sigma^2 V^T \quad (\because \Sigma \text{ is diagonal matrix})$$

(iii)

q2 b)

$$\text{original size} = (512 \times 512 \times 3) \times (8)$$

$\downarrow$   
#pixels       $\downarrow$  size of each pixel (0-255)  
(in bytes)      (in bytes)

$$\text{compressed size} = (512 \times 512 \times 1) \times 4 + (16 \times 3) \times 8$$

$\downarrow$  size of each pixel  
(in bytes)

each pixel can be assigned a number ~~between~~ between 0 to 15 depending on the centroid it is assigned to

each pixel's value range from (0-255)

16 centroids having 3 color channel

$$\text{compressed factor} = \frac{\text{original size (bytes)}}{\text{compressed size (bytes)}}$$

$$= \frac{512 \times 512 \times 3 \times 8}{512 \times 512 \times 4 + 16 \times 3 \times 8}$$

$$\approx 6$$

Q3 @ T.P.  $D_{KL}(P||Q) \geq 0$  &  $P, Q$   
and  $D_{KL}(P||Q) = 0$  iff  $P = Q$ .

by Jensen's inequality,

$$E(f(x)) \geq f(E(x))$$

consider the function  $f(x) = -\log x$ , which is strictly convex. On using Jensen's inequality:

$$\sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)} = E_P \left[ \log \frac{P(x)}{Q(x)} \right] \geq \log E_P \left[ \frac{P(x)}{Q(x)} \right]$$

$$D_{KL}(P||Q) = \sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)}$$

$$E_P \left( -\log \left( \frac{Q(x)}{P(x)} \right) \right) \geq -\log E_P \left( \frac{Q(x)}{P(x)} \right)$$

expanding we get

$$-\sum_{x \in X} P(x) \log \left( \frac{Q(x)}{P(x)} \right) \geq -\log \sum_{x \in X} \left( P(x) \frac{Q(x)}{P(x)} \right)$$

$$D_{KL}(P||Q) \geq -\log 1 = 0 \quad \therefore \sum_{x \in X} Q(x) = 1$$

$$D_{KL}(P||Q) = \sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)} \geq 0$$

equality holds iff argument of expectation fn  
satisfies  $\log \frac{P(x)}{Q(x)} = 0 \quad \forall x \quad \} \text{ for Jensen equality}$   
to hold.

$$\Leftrightarrow \frac{P(x)}{Q(x)} = 1 \quad \forall x \in X$$

$$\Leftrightarrow P(x) = Q(x) \quad \forall x \in X$$

Q3(b) Expanding joint KL divergence,

$$D_{KL}(P(x,y) \parallel Q(x,y)) = \sum_{x,y} P(x,y) \log \left( \frac{P(x,y)}{Q(x,y)} \right)$$

$$\text{but } P(x,y) = P(x) P(y|x), \quad Q(x,y) = Q(x) Q(y|x)$$

$$= \sum_x \sum_y P(x) P(y|x) \left[ \log \left( \frac{P(x)}{Q(x)} \right) + \log \left( \frac{P(y|x)}{Q(y|x)} \right) \right]$$

~~split into two terms~~

$$= \sum_x \left[ P(x) \log \frac{P(x)}{Q(x)} \left( \sum_y P(y|x) \right) + P(x) \sum_y P(y|x) \log \frac{P(y|x)}{Q(y|x)} \right]$$

$$= \sum_x \left[ P(x) \log \frac{P(x)}{Q(x)} \sum_y P(y|x) \right] + \sum_x \left[ P(x) \sum_y P(y|x) \log \frac{P(y|x)}{Q(y|x)} \right]$$

$$= D_{KL}(P(x) \parallel Q(x)) + \sum_{P(x)} [D_{KL} P(y|x) \parallel Q(y|x)]$$

Q3(c)  $D_{KL}(\hat{P} \parallel P_\theta) = \sum_x \hat{P}(x) \log \frac{\hat{P}(x)}{P_\theta(x)}$  const.

$$= -\cancel{P_\theta(x)} - \sum_x \hat{P}(x) \log P_\theta(x) + K$$

$$\therefore - \sum_x \left\{ \cancel{\sum_{i=1}^n 1\{x_i=x\}} \right\} \log P_\theta(x) + K$$

$$= - \sum_{i=1}^n \log P_\theta(x^i) + K$$

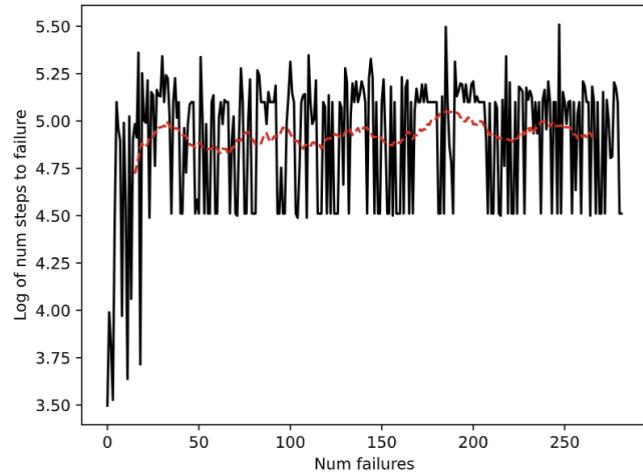
$$\therefore \arg \min_{\theta} D_{KL}(\hat{P} \parallel P_\theta) \Leftrightarrow \arg \max_{\theta} \sum_{i=1}^n \log P_\theta(x^i)$$

Q2 a)



Q2 b)

**Q4 a)**



**Q4 b i)**

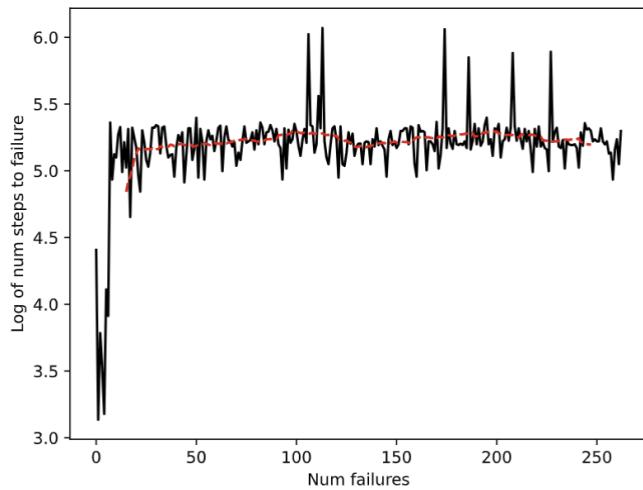
[INFO] Failure number 282

It took 282 trials for the algorithm to converge.

**Q4 b ii)**

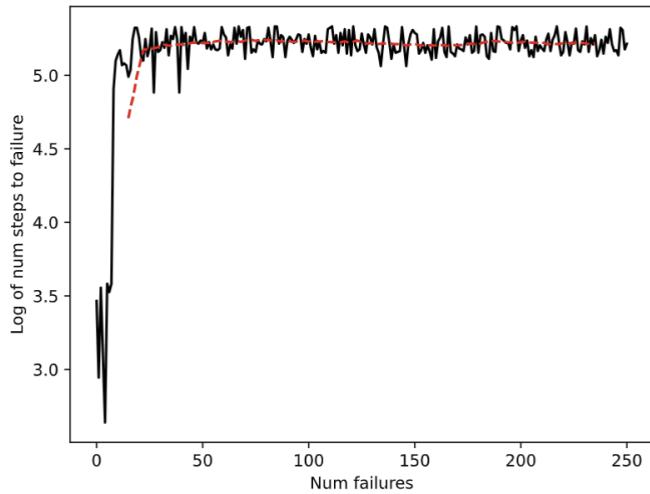
SEED = 1

[INFO] Failure number 263 ( #trials before convergence )



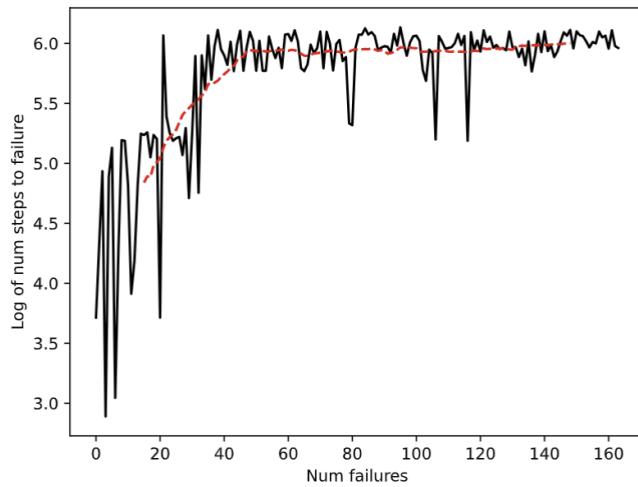
SEED = 2

[INFO] Failure number 251



SEED = 3

[INFO] Failure number 164



Different seeds cause different initialization and random transitions, therefore the algorithm's convergence changes