

Q₁.

(a) $\because Q$ Orthogonal \therefore By def: $Q^T Q = Q Q^T = I$ $Q \in \mathbb{R}^{n \times n}$ $n \in \mathbb{R}$

$$(i) (Q^T)^T Q^T$$

$$= Q Q^T$$

$$= I$$

$\therefore Q^T$ is orthogonal

$$(Q^{-1})^T (Q^{-1})$$

$$= (Q^T)^{-1} Q^{-1}$$

From $Q^T Q = I$, We also get $Q^T = Q^{-1}$

$$\therefore (Q^T)^{-1} Q^{-1} = (Q^{-1})^{-1} Q^{-1} = Q Q^{-1} = I$$

$\therefore Q^{-1}$ is orthogonal

(ii)

Let $Q v_i = \lambda_i v_i$ $\langle \lambda_i, v_i \rangle$ is a pair of eigenvalue - eigenvector

$$(Q v_i)^T = v_i^T Q^T = (\lambda_i v_i)^T = \lambda_i v_i^T \quad (\lambda_i \in \mathbb{R})$$

$$\therefore (Q v_i)^T Q v_i = \lambda_i v_i^T \cdot \lambda_i v_i$$

$$v_i^T Q^T Q v_i = \lambda_i^2 v_i^T v_i$$

$$\because Q^T Q = I$$

$$\therefore v_i^T I v_i = \lambda_i^2 v_i^T v_i$$

$$\therefore v_i^T I = v_i$$

$$\therefore v_i^T v_i = \lambda_i^2 v_i^T v_i$$

$$\therefore \|v_i\|_2^2 = \lambda_i^2 \|v_i\|_2^2$$

$$\therefore \lambda_i^2 = 1$$

$$\therefore \lambda_i = 1 \text{ OR } \lambda_i = -1$$

(iii) $\det(Q) = \prod_{i=1}^n \lambda_i$

$\therefore \lambda_i = 1 \text{ OR } -1$

$\therefore \det(Q) = 1 \text{ OR } -1$

(iv) $Qv_i = \lambda_i v_i$

$$\|Qv_i\|_2^2 = v_i^T Q^T Q v_i = v_i^T v_i = \|v_i\|_2^2$$

\therefore Length of v_i is preserved

(b) Consider Matrix A diagonalizable

Assume A not Full Rank

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

- D Diagonal

$$A = U \Sigma V^T = U_i S V_i^T ;$$

- S Consists of non-zero Singular values of A
- Columns of U_i, V_i^T as orthonormal basis of $\text{Col}(A)$

- U, V Orthogonal ; Σ diagonal

(ii) According to how SVD is calculated

σ_i , (Singular values of A)

$$= \sqrt{\lambda_i} \quad (\lambda_i \text{ is the eigenvalue of } A^T A)$$

$$= \sqrt{\lambda_i'} \quad (\lambda_i' \text{ is the eigenvalue of } A A^T)$$

(i) Do SVD on A

$$A = U \Sigma V^T = U_i S V_i^T$$

$$A A^T = U_i S V_i^T \underbrace{V_i S^T V_i^T}_{I} U_i^T = U_i S^2 U_i^T$$

Let's do Diagonalization on $A A^T$

$$\therefore (A A^T)^T = A A^T, \quad A A^T \text{ is Symmetric}$$

$$\therefore A A^T = T D T^T$$

$$\therefore U_i S^2 U_i^T = T D T^T$$

$$\therefore \text{Again: } S^2 = D \rightarrow \sigma_i^2 = \lambda_i$$

$$U_i = T \rightarrow \text{Singular vectors } (U_i) \text{ of } A = \text{Eigenvectors of } A A^T$$

Similarly

$$\begin{aligned} A^T A &= V_1 S^T U_1^T U_1 S V_1^T \\ &= V_1 S^2 V_1^T \\ &= T D T^T \end{aligned}$$

$$\therefore S^2 = D$$

$$V_1 = T$$

$$\rightarrow \sigma_i^2 = \lambda_i$$

\rightarrow singular vector v_i = eigenvector of $A^T A$

(c)

i. False

Counter E.g. $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

$$\lambda_1 = \lambda_2 = 1$$

Should be at most n distinct eigenvalues

ii. False

$$Av_1 = \lambda_1 v_1$$

$$Av_2 = \lambda_2 v_2$$

$$A(v_1 + v_2) = \lambda_1 v_1 + \lambda_2 v_2, \text{ Generally } \neq \lambda'(v_1 + v_2)$$

iii. True

A : P.S.D. Matrix

$$\text{Let } Av_i = \lambda_i v_i$$

$$\begin{aligned} v_i^T A v_i &= v_i^T \lambda_i v_i \\ &= \lambda_i v_i^T v_i \quad (\lambda_i \in \mathbb{R}) \\ &= \lambda_i \|v_i\|_2^2 \end{aligned}$$

$$\therefore \|v_i\|_2^2 \geq 0$$

$$A \text{ P.S.D.} \therefore v_i^T A v_i \geq 0$$

$$\therefore \lambda_i \text{ must be } \geq 0$$

iv. True

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

$$\text{Rank}(A) = 2$$

$$\lambda_1 = \lambda_2 = 1$$

$$\# \text{ of distinct } \lambda_i = 1 < \text{Rank}(A)$$

v. True

$$Av_1 = \lambda v_1$$

$$Av_2 = \lambda v_2$$

$$\begin{aligned} A(v_1 + v_2) &= \lambda v_1 + \lambda v_2 \\ &= \lambda(v_1 + v_2) \end{aligned}$$

$$\text{Let } v_1 + v_2 = v' \Rightarrow Av' = \lambda v'$$

Q2.

(a)

$$(i) \quad P(A \text{ not hit}) \\ = \sum_{n=1}^{\infty} P(A \text{ not hit, } n^{\text{th}} \text{ round})$$

Assume duel ends in the n^{th} round and A not hit

$$P(A^c, n) \Leftrightarrow \frac{\text{Both Miss}}{1} \cdot \frac{\text{Both Miss}}{2} \cdot \frac{\text{Both Miss}}{3} \cdots \frac{\text{Both Miss}}{n-1} \cdot \frac{B \text{ Hit}}{n}$$

$$\begin{aligned} \therefore P(A \text{ not hit, } n^{\text{th}} \text{ round}) \\ &= [(1-P_A)(1-P_B)]^{n-1} \cdot (1-P_A) \cdot P_B \\ &= (1-P_A)^n \cdot (1-P_B)^{n-1} \cdot P_B \end{aligned}$$

$$\therefore P(A \text{ not hit}) = \sum_{n=1}^{\infty} P_B \cdot (1-P_B)^{n-1} \cdot (1-P_A)^n$$

It's a sum of ∞ Geometric Series

$$\begin{aligned} \therefore P(A \text{ not hit}) &= \frac{a_1}{1-r} \\ &= \frac{P_B (1-P_A)}{1 - (1-P_B)(1-P_A)} \end{aligned}$$

(ii) P (Both A, B are hit)

Similar to (i)

$$\begin{aligned}
 P(\text{Both A \& B hit}) &= \sum_{n=1}^{\infty} (\text{Both A \& B hit, } n^{\text{th}} \text{ Round}) \\
 &= \sum_{n=1}^{\infty} [(1-P_A)(1-P_B)]^{n-1} P_A P_B \\
 &= \frac{P_A P_B}{1 - (1-P_A)(1-P_B)}
 \end{aligned}$$

(iii) Let $S \Leftrightarrow$ duel Ending in the n^{th} Round

A	B	Both	$n^{\text{th}} S$
hit	hit	hit	Miss

$$\begin{aligned}
 \therefore P(S) &= P(A \text{ hit, } n^{\text{th}} \text{ Round}) + P(B \text{ hit, } n^{\text{th}} \text{ Round}) \\
 &\quad + P(A, B \text{ hit, } n^{\text{th}} \text{ Round}) + \underbrace{P(A, B \text{ Miss, } n^{\text{th}} \text{ Round})}_{=0 \quad \because A, B \text{ Miss} \\
 &\quad \text{Duel won't end}}
 \end{aligned}$$

$$\begin{aligned}
 &= [(1-P_A)(1-P_B)]^{n-1} P_A (1-P_B) + [(1-P_A)(1-P_B)]^{n-1} P_B (1-P_A) \\
 &\quad + [(1-P_A)(1-P_B)]^{n-1} P_A P_B \\
 &= [(1-P_A)(1-P_B)]^{n-1} [P_A (1-P_B) + P_B (1-P_A) + P_A P_B] \\
 &= [(1-P_A)(1-P_B)]^{n-1} [P_A + P_B - P_A P_B] \\
 &= [(1-P_A)(1-P_B)]^{n-1} [1 - (1-P_A)(1-P_B)]
 \end{aligned}$$

$$\begin{aligned}
 \text{(iv)} \quad & P(A^c, \text{ duel ends @ } n^{\text{th}} \text{ round}) \\
 &= P(B \text{ hits}, n^{\text{th}} \text{ round}) \\
 &= [(1-p_A)(1-p_B)]^{n-1} p_B
 \end{aligned}$$

$$\begin{aligned}
 & P(A^c \mid n^{\text{th}} \text{ round}) \\
 &= P(A^c, n^{\text{th}} \text{ round}) / P(n^{\text{th}} \text{ round}) \\
 &= \frac{[(1-p_A)(1-p_B)]^{n-1} p_B}{[(1-p_A)(1-p_B)]^{n-1} (p_A + p_B - p_A p_B)} \\
 &= \frac{p_B}{p_A + p_B - p_A p_B}
 \end{aligned}$$

$$\begin{aligned}
 \text{(v)} \quad & P(A, B \mid n^{\text{th}} \text{ round}) \\
 &= \frac{P(A, B \text{ hit}, n^{\text{th}} \text{ round})}{P(n^{\text{th}} \text{ round})} \\
 &= \frac{[(1-p_A)(1-p_B)]^{n-1} p_A p_B}{[(1-p_A)(1-p_B)]^{n-1} (p_A + p_B - p_A p_B)} \\
 &= \frac{p_A p_B}{p_A + p_B - p_A p_B}
 \end{aligned}$$

(h) Let R.V. $X \sim \#$ of Faculties Isolated

$$\text{Let } x_i \sim \begin{cases} 1 & i^{\text{th}} \text{ faculty ISO} \\ 0 & i^{\text{th}} \text{ faculty not ISO} \end{cases} \quad i \in [1, 18]$$

$$\therefore X = \sum_{i=1}^{18} x_i$$

$$\therefore E[X] = E\left[\sum_{i=1}^{18} x_i\right] = \sum_{i=1}^{18} E[x_i] = 18 \cdot E[x_i]$$

$\therefore x_i \sim \text{Bernoulli R.V.}$

$$\therefore E[x_i] = P[x_i=1] = P[i^{\text{th}} \text{ faculty ISO}]$$

Assume if $x_i \in \text{ECE}$

"ISO":

Left	ECE	Right
\uparrow		\uparrow
$12/17$		$11/16$

$$\therefore P[x_i=1] = P[x_i=1 | x_i \in \text{ECE}] P[x_i \in \text{ECE}]$$

$$+ P[x_i=1 | x_i \in \text{CS}] P[x_i \in \text{CS}]$$

$$+ P[x_i=1 | x_i \in \text{MTH}] P[x_i \in \text{MTH}]$$

$$= \left(\frac{12}{17} \cdot \frac{11}{16} \cdot \frac{6}{18} \right) \times 3$$

$$= \frac{12}{17} \cdot \frac{11}{16} = \frac{33}{68}$$

$$\therefore E[x_i] = \frac{33}{68}$$

$$\therefore E[X] = 18 \cdot \frac{33}{68}$$

$$= \frac{297}{34}$$

(ii) Similarly $E[X] = 18 E[X_i]$

Semi
-Happy

$$\begin{array}{ccc} \overline{\frac{5}{17}} & \xrightarrow{\text{ECE}} & \overline{\frac{12}{16}} \\ & \text{OR} & \\ \frac{12}{17} & & \frac{5}{16} \end{array}$$

$$\begin{aligned} \therefore E[X_i] &= \left(\frac{5}{17} \cdot \frac{12}{16} + \frac{12}{17} \cdot \frac{5}{16} \right) \cdot \frac{1}{3} \cdot 3 \\ &= \frac{15}{34} \end{aligned}$$

$$\begin{aligned} \therefore E[X] &= 18 E[X_i] \\ &= \frac{135}{17} \end{aligned}$$

(iii)

$$\begin{array}{ccc} \overline{\frac{5}{17}} & \xrightarrow{\text{ECE}} & \overline{\frac{4}{16}} \end{array}$$

$$\therefore E[X_i] = \frac{5}{17} \cdot \frac{4}{16} = \frac{5}{68}$$

$$\therefore E[X] = 18 E[X_i] = \frac{45}{34}$$

(d)

$$\begin{aligned} E[Ax+b] &= E[Ax] + E[b] \\ &= A E[x] + b \\ &= A \cdot \begin{bmatrix} E[x_1] \\ E[x_2] \\ \vdots \\ E[x_n] \end{bmatrix} + b \end{aligned}$$

(e)

$$\begin{aligned} \text{Cov}[x] &= E[(x - E[x])(x - E[x])^T] \\ \text{Cov}[Ax+b] &= E[(Ax+b - A \cdot E[x] - b)(Ax+b - A \cdot E[x] - b)^T] \\ &= E[A(x - E[x]) \cdot (A(x - E[x]))^T] \\ &= E[A \cdot [(x - E[x])(x - E[x])^T] \cdot A^T] \\ &= E[A] \cdot E[(x - E[x])(x - E[x])^T] \cdot E[A^T] \\ &= A \cdot \text{Cov}[x] \cdot A^T \end{aligned}$$

(c) $D \sim \text{Cancer Pangoous}$

$+|D \sim \text{True Positive Test}$

$+|D^c \sim \text{False Positive Test}$

$$P[D] = 0.0005 \quad P[+|D] = 0.9 \quad P[+|D^c] = 0.01$$

$$(i) \quad P[D|+] = \frac{P[D, +]}{P[+]}$$

$$= \frac{P[+|D] P[D]}{P[+|D] P[D] + P[+|D^c] P[D^c]}$$

$$= \frac{0.9 \times 0.0005}{0.9 \times 0.0005 + 0.01 \times (1 - 0.0005)}$$

$$= \frac{90}{2089} \approx 0.043$$

$$(ii) \quad P[D|-] = \frac{P[D, -]}{P[-]} = \frac{P[-|D] P[D]}{P[-|D] P[D] + P[-|D^c] P[D^c]}$$
$$= \frac{(1 - 0.9) \times 0.0005}{(1 - 0.9) \times 0.0005 + (1 - 0.01) \times (1 - 0.0005)}$$
$$\approx 5.0528 \times 10^{-5}$$

Q3.

$$(a) \nabla_x x^T A y = \frac{\partial (x^T A y)}{\partial x} = \frac{\partial (x^T (A y))}{\partial x}$$

$$\therefore A \in \mathbb{R}^{n \times m}$$

$$y \in \mathbb{R}^{m \times 1}$$

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

$$A y \in \mathbb{R}^n$$

$$\therefore \nabla_x x^T A y = A y \quad (\text{property 69})$$

check: ① $A y \in \mathbb{R}^n \quad x \in \mathbb{R}^n$

② if all scalar

$$\frac{\partial (x^T A y)}{\partial x} = \frac{\partial (x A y)}{\partial x} = A y$$

$$(b) \nabla_y x^T A y = \frac{\partial ((x^T A) y)}{\partial y} = \frac{\partial ((A^T x)^T y)}{\partial y}$$

$$A^T \in \mathbb{R}^{m \times n}$$

$$\therefore A^T x \in \mathbb{R}^m$$

$$y \in \mathbb{R}^m$$

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

$$\therefore \nabla_y x^T A y = \nabla_y (A^T x)^T y = A^T x \quad (\text{property 69})$$

check: ① $A^T x \in \mathbb{R}^m \quad y \in \mathbb{R}^m$

② if all scalar

$$\frac{\partial x^T A y}{\partial y} = x^T A = A^T x$$

$$(c) \nabla_A x^T A y = x y^T \quad (\text{property 70})$$

check: ① $x y^T \in \mathbb{R}^{n \times m} \quad A \in \mathbb{R}^{n \times m}$

② if all scalar

$$\frac{\partial x^T A y}{\partial A} = x^T y = x y^T$$

$$(d) f = x^T A x + b^T x$$

$$\nabla_x f = \frac{\partial (x^T A x + b^T x)}{\partial x} = \frac{\partial (x^T A x)}{\partial x} + \frac{\partial (b^T x)}{\partial x}$$

$$= (A + A^T) x + b$$

(property 97)

check: ①

$$\begin{matrix} (A + A^T) x + b \\ \uparrow \quad \uparrow \quad \uparrow \quad \uparrow \\ n \times n \quad n \times n \quad n \times 1 \quad n \times 1 \end{matrix} \in \mathbb{R}^{n \times 1}$$

$$x \in \mathbb{R}^{n \times 1}$$

② if all scalar

$$\frac{\partial f}{\partial x} = \frac{\partial (A x^T + b x)}{\partial x}$$

$$= 2 A x + b = (A + A^T) x + b$$

(e) $f = \text{Tr}(AB)$

$$\nabla_A f = \frac{\partial (\text{Tr}(AB))}{\partial A} = B^T \quad (\text{property 100})$$

$$A \in \mathbb{R}^{n \times n}$$

$$B \in \mathbb{R}^{n \times n}$$

Check:

① $B^T \in \mathbb{R}^{n \times n}$

$$A \in \mathbb{R}^{n \times n}$$

② if all scalar

$$\frac{\partial (\text{Tr}(AB))}{\partial A} = \frac{\partial (AB)}{\partial A} = B = B^T$$

(f) $f = \text{Tr}(BA + A^T B + A^2 B)$

$$\therefore \text{Tr}(A^T) = \text{Tr}(A)$$

$$\text{Tr}(AB) = \text{Tr}(BA)$$

$$\text{Tr}(A+B) = \text{Tr}(A) + \text{Tr}(B)$$

$$\begin{aligned} \therefore f &= \text{Tr}(BA) + \text{Tr}(A^T B) + \text{Tr}(A^2 B) \\ &= \text{Tr}(AB) + \text{Tr}(A^T B) + \text{Tr}(A^2 B) \end{aligned}$$

$$\therefore \nabla_A f = \frac{\partial (\text{Tr}(AB))}{\partial A} + \frac{\partial (\text{Tr}(A^T B))}{\partial A} + \frac{\partial (\text{Tr}(A^2 B))}{\partial A}$$

$$= \underset{\substack{\uparrow \\ \text{property} \\ 100}}{B^T} + \underset{\substack{\uparrow \\ \text{property} \\ 103}}{B} + \underset{\substack{\uparrow \\ \text{property} \\ 107}}{(AB+BA)^T}$$

Check:

① $\underset{\substack{\uparrow \\ n \times n}}{B^T} + \underset{\substack{\uparrow \\ n \times n}}{B} + (\underset{\substack{\uparrow \\ n \times n}}{AB} + \underset{\substack{\uparrow \\ n \times n}}{BA})^T \in \mathbb{R}^{n \times n}$

$$A \in \mathbb{R}^{n \times n}$$

② if all scalar

$$\begin{aligned} f &= \text{Tr}(BA + A^T B + A^2 B) \\ &= 2AB + A^2 B \end{aligned}$$

$$\begin{aligned} \frac{\partial f}{\partial A} &= 2B + 2AB \\ &= B^T + B + (AB + BA)^T \end{aligned}$$

$$(9) \|A + \lambda B\|_F^2$$

$$= \text{Tr}((A + \lambda B)^T (A + \lambda B))$$

$$= \text{Tr}((A^T + \lambda B^T)(A + \lambda B))$$

$$= \text{Tr}(A^T A + A^T \lambda B + \lambda B^T A + \lambda^2 B^T B)$$

$$= \text{Tr}(A^T A) + \lambda \text{Tr}(A^T B) + \lambda \text{Tr}(B^T A) + \lambda^2 \text{Tr}(B^T B)$$

$$= \text{Tr}(A^T A) + \lambda \cdot \text{Tr}(A^T B) + \lambda \text{Tr}((B^T A)^T) + \lambda^2 \text{Tr}(B^T B)$$

$$= \text{Tr}(A^T A) + 2\lambda \cdot \text{Tr}(A^T B) + \lambda^2 \text{Tr}(B^T B)$$

Dropping term with no λ dependence, $f = \text{Tr}(A^T A) + 2\lambda \cdot \text{Tr}(A^T B)$

$$\therefore \nabla_A f = 2A + 2\lambda \cdot B$$

\uparrow
property
115

\uparrow
property
103

Check: ① $A \in \mathbb{R}^{n \times n}$

$$\underset{n \times n}{2A} + 2\lambda \cdot \underset{n \times n}{B} \in \mathbb{R}^{n \times n}$$

② if all scalar

$$\begin{aligned} \frac{\partial f}{\partial A} &= \frac{\partial (A + \lambda B)^2}{\partial A} \\ &= 2(A + \lambda B) \cdot 1 \\ &= 2A + 2\lambda B \end{aligned}$$

Q4.

$$X \in \mathbb{R}^{m \times n}, X = \begin{bmatrix} x_1 & x_2 & \dots & x_n \\ 1 & 1 & & 1 \end{bmatrix} \quad \|x_i\|_2^2 = \sum_{j=1}^m (x_{ij})^2 \Leftarrow \text{Sum all entries in Column } i$$

$$\|X\|_F^2 = \sum_{i=1}^n \|x_i\|_2^2 \Leftarrow \text{Sum Column norms}$$

$$= \sum_{i=1}^n \sum_{j=1}^m (x_{ij})^2$$

$$\hat{y} = Wx \quad \therefore W \in \mathbb{R}^{n \times n}$$

\uparrow \mathbb{R}^n \uparrow \mathbb{R}^n

$$\mathcal{L}(w) = \frac{1}{2} \sum_{i=1}^n \|y^{(i)} - Wx^{(i)}\|^2$$

$$\therefore \text{Find: } \frac{\partial \mathcal{L}(w)}{\partial w} / \nabla_w (\mathcal{L}(w))$$

$$\text{Let } Y = \begin{bmatrix} y_1^{(1)} & y_1^{(2)} & \dots & y_1^{(n)} \\ 1 & 1 & & 1 \end{bmatrix}$$

$$X = \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(n)} \\ 1 & 1 & & 1 \end{bmatrix}$$

$$\therefore \text{According to } \|X\|_F^2 = \sum_{i=1}^n \|x_i\|_2^2$$

$$\therefore \mathcal{L}(w) = \frac{1}{2} \sum_{i=1}^n \|y^{(i)} - Wx^{(i)}\|^2$$

$$= \frac{1}{2} \cdot \|Y - WX\|_F^2$$

$$= \frac{1}{2} \cdot \text{Tr}((Y - WX)^T \cdot (Y - WX))$$

$$= \frac{1}{2} \cdot \text{Tr}((Y^T - X^T W^T) \cdot (Y - WX))$$

$$= \frac{1}{2} \cdot \text{Tr}(Y^T Y - Y^T W X - X^T W^T Y + X^T W^T W X)$$

$$= \frac{1}{2} \cdot [\text{Tr}(Y^T Y) - \text{Tr}(Y^T W X) - \text{Tr}((Y^T W X)^T) + \text{Tr}(X^T W^T W X)]$$

$$= \frac{1}{2} \cdot [\cancel{\text{Tr}(Y^T Y)} - 2 \cdot \text{Tr}(Y^T W X) + \text{Tr}(X^T W^T W X)]$$

Drop terms with no dependence to W

$$\therefore \mathcal{L}(W) = -\text{Tr}(Y^T W X) + \frac{1}{2} \cdot \text{Tr}(X^T W^T W X)$$

$$\therefore \text{Tr}(ABC) = \text{Tr}(CAB) = \text{Tr}(BCA)$$

$$\therefore \text{Tr}(\underline{X^T W^T W X}) = \text{Tr}(W^T W X X^T)$$

$$\begin{aligned} \therefore \frac{\partial \mathcal{L}(W)}{\partial W} &= -(Y^T)^T X^T + \frac{1}{2} \cdot (W \cdot (X X^T)^T + W \cdot (X X^T)) \\ &= -Y X^T + \frac{1}{2} \cdot 2 \cdot W X X^T \\ &= -Y X^T + W X X^T \quad (\text{property 101, 113}) \end{aligned}$$

Set Left Hand to 0

$$\therefore Y X^T = W X X^T$$

$$\therefore W = Y X^T (X X^T)^{-1}$$

Q5.

We've known from the Lecture that:

$$\begin{aligned} \mathcal{L}(\theta) &= \frac{1}{2} \sum_{i=1}^N (y^{(i)} - \theta^T \hat{x}^{(i)})^2 \\ &= \frac{1}{2} (Y^T Y - 2Y^T X \theta + \theta^T X^T X \theta) \end{aligned}$$

$$\begin{aligned} \text{Now: } \mathcal{L}(\theta) &= \frac{1}{2} \sum_{i=1}^N (y^{(i)} - \theta^T \hat{x}^{(i)})^2 + \frac{\lambda}{2} \|\theta\|_2^2 \\ &= \frac{1}{2} (\cancel{Y^T Y} - 2Y^T X \theta + \theta^T X^T X \theta) + \frac{\lambda}{2} \theta^T \theta \end{aligned}$$

Find $\frac{\partial \mathcal{L}(\theta)}{\partial \theta}$: : Get rid of $Y^T Y$ since it's not θ related

$$\mathcal{L}(\theta) = -Y^T X \theta + \frac{1}{2} \theta^T X^T X \theta + \frac{\lambda}{2} \theta^T \theta$$

$$\begin{aligned} \nabla_{\theta} (-Y^T X \theta) &= -(Y^T X)^T \quad (\text{property 6?}) \\ &= -X^T Y \end{aligned}$$

$$\begin{aligned} &\frac{1}{2} \theta^T X^T X \theta + \frac{\lambda}{2} \theta^T \theta \\ &= \frac{1}{2} \theta^T (X^T X \theta + \lambda I \cdot \theta) \\ &= \frac{1}{2} \theta^T \cdot (X^T X + \lambda I) \cdot \theta \end{aligned}$$

$$\begin{aligned} \therefore \nabla_{\theta} \left(\frac{1}{2} \theta^T X^T X \theta + \frac{\lambda}{2} \theta^T \theta \right) &= \frac{1}{2} \cdot (X^T X + \lambda I + (X^T X + \lambda I)^T) \theta \\ &= \frac{1}{2} \cdot (X^T X + \lambda I + X^T X + \lambda I) \theta \\ &= (X^T X + \lambda I) \theta \end{aligned}$$

$$\therefore \nabla_{\theta} (\mathcal{L}(\theta)) = -X^T Y + (X^T X + \lambda I) \theta$$

$$\text{Let } \nabla_{\theta} (\mathcal{L}(\theta)) = 0$$

$$\therefore X^T Y = (X^T X + \lambda I) \theta^* \quad \therefore \theta^* = (X^T X + \lambda I)^{-1} X^T Y$$

linear_regression

January 22, 2024

0.1 Linear regression workbook

This workbook will walk you through a linear regression example. It will provide familiarity with Jupyter Notebook and Python. Please print (to pdf) a completed version of this workbook for submission with HW #1.

ECE C147/C247, Winter Quarter 2024, Prof. J.C. Kao, TAs: T.Monsoor, Y. Liu, S. Rajesh, L. Julakanti, K. Pang

```
[19]: import numpy as np
import matplotlib.pyplot as plt

#allows matlab plots to be generated in line
%matplotlib inline
```

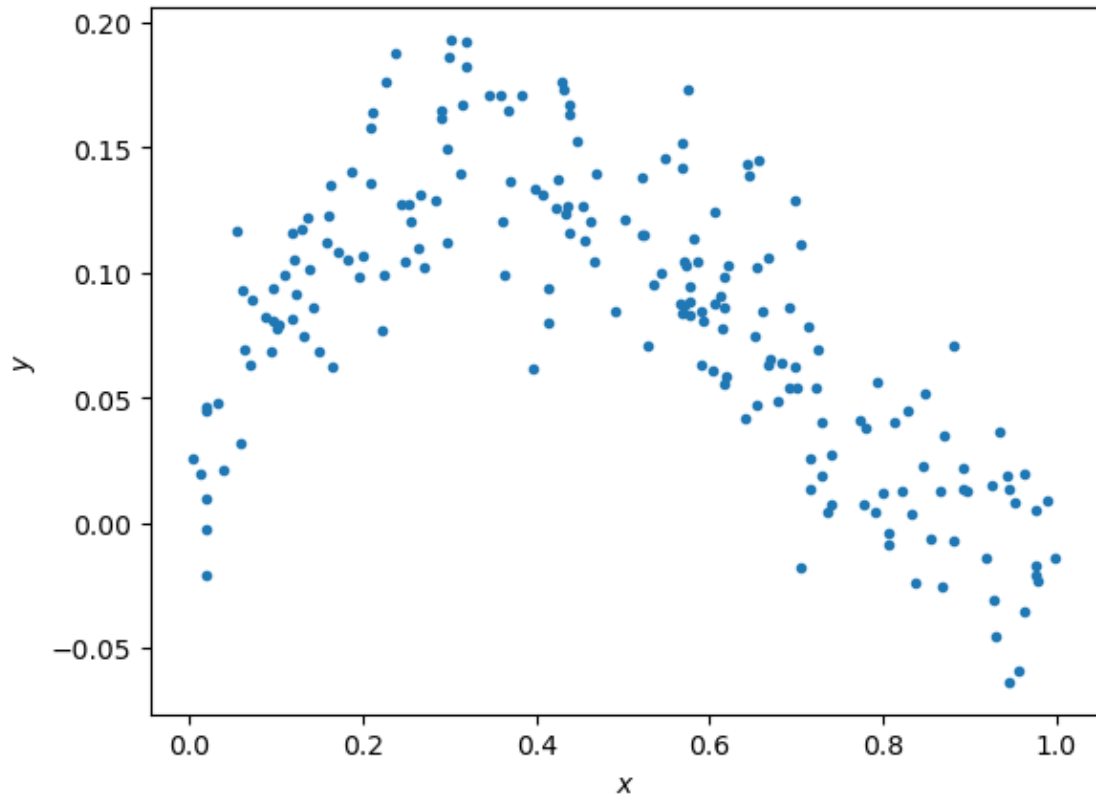
0.1.1 Data generation

For any example, we first have to generate some appropriate data to use. The following cell generates data according to the model: $y = x - 2x^2 + x^3 + \epsilon$

```
[20]: np.random.seed(0) # Sets the random seed.
num_train = 200 # Number of training data points

# Generate the training data
x = np.random.uniform(low=0, high=1, size=(num_train,))
y = x - 2*x**2 + x**3 + np.random.normal(loc=0, scale=0.03, size=(num_train,))
f = plt.figure()
ax = f.gca()
ax.plot(x, y, '.')
ax.set_xlabel('$x$')
ax.set_ylabel('$y$')
```

```
[20]: Text(0, 0.5, '$y$')
```



0.1.2 QUESTIONS:

Write your answers in the markdown cell below this one:

- (1) What is the generating distribution of x ?
- (2) What is the distribution of the additive noise ϵ ?

0.1.3 ANSWERS:

- (1) x is a uniform distribution
- (2) ϵ is a normal distribution with mean 0 and STD 0.03

0.1.4 Fitting data to the model (5 points)

Here, we'll do linear regression to fit the parameters of a model $y = ax + b$.

```
[21]: # xhat = (x, 1)
xhat = np.vstack((x, np.ones_like(x)))

# ===== #
# START YOUR CODE HERE #
# ===== #
```

```

# GOAL: create a variable theta; theta is a numpy array whose elements are [a,
↪ b]

theta = (np.linalg.inv((xhat)@(xhat.T))) @ (xhat@y) # please modify this line

# ===== #
# END YOUR CODE HERE #
# ===== #

```

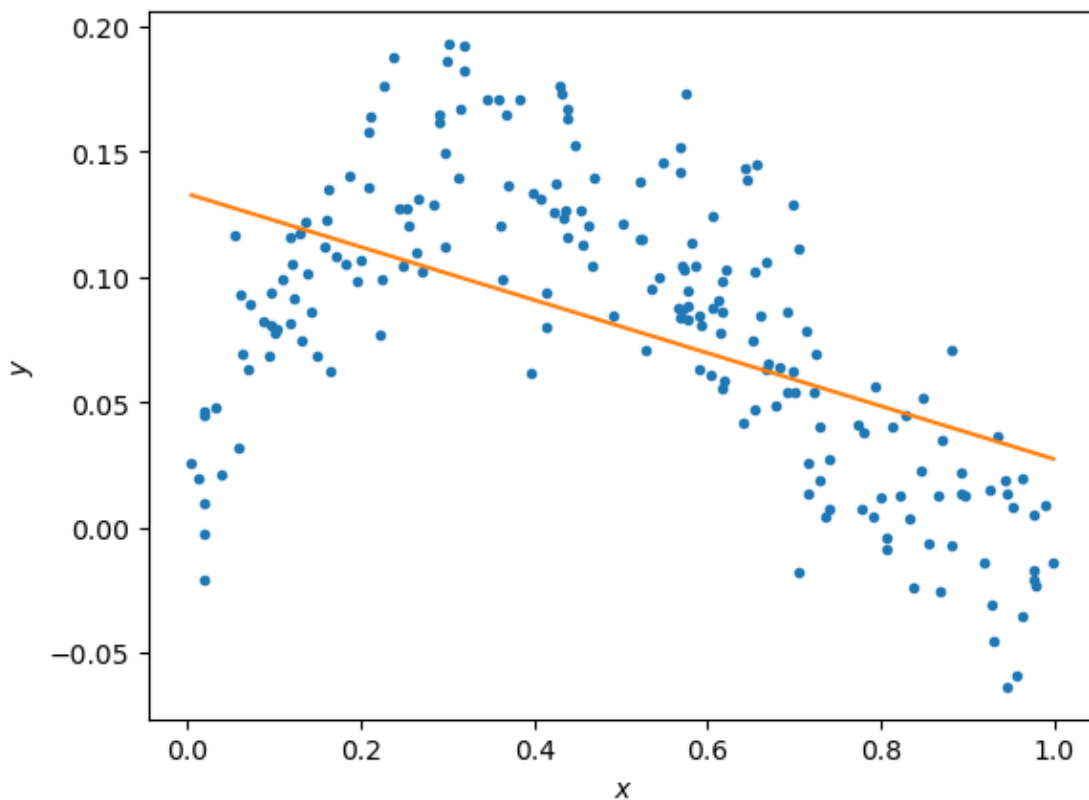
```

[22]: # Plot the data and your model fit.
f = plt.figure()
ax = f.gca()
ax.plot(x, y, '.')
ax.set_xlabel('$x$')
ax.set_ylabel('$y$')

# Plot the regression line
xs = np.linspace(min(x), max(x), 50)
xs = np.vstack((xs, np.ones_like(xs)))
plt.plot(xs[0,:], theta.dot(xs))

```

[22]: [[matplotlib.lines.Line2D](#) at 0x1e7a296c790>]



0.1.5 QUESTIONS

- (1) Does the linear model under- or overfit the data?
- (2) How to change the model to improve the fitting?

0.1.6 ANSWERS

- (1) It underfits the data.
- (2) We can instead use a higher order polynomial.

0.1.7 Fitting data to the model (5 points)

Here, we'll now do regression to polynomial models of orders 1 to 5. Note, the order 1 model is the linear model you prior fit.

```
[23]: N = 5
xhats = []
thetas = []

# ===== #
# START YOUR CODE HERE #
# ===== #

# GOAL: create a variable thetas.
# thetas is a list, where theta[i] are the model parameters for the polynomial
#   ↪ fit of order i+1.
#   i.e., thetas[0] is equivalent to theta above.
#   i.e., thetas[1] should be a length 3 np.array with the coefficients of the
#   ↪  $x^2$ ,  $x$ , and 1 respectively.
#   ... etc.

xhats.append(xhat)
thetas.append(theta)
xhat_i = xhat
for i in range(2,6):
    x_power = np.array(x**i)
    xhat_i = np.vstack((x_power, xhat_i))
    theta_i = (np.linalg.inv((xhat_i)@(xhat_i.T))) @ (xhat_i@y)
    xhats.append(xhat_i)
    thetas.append(theta_i)

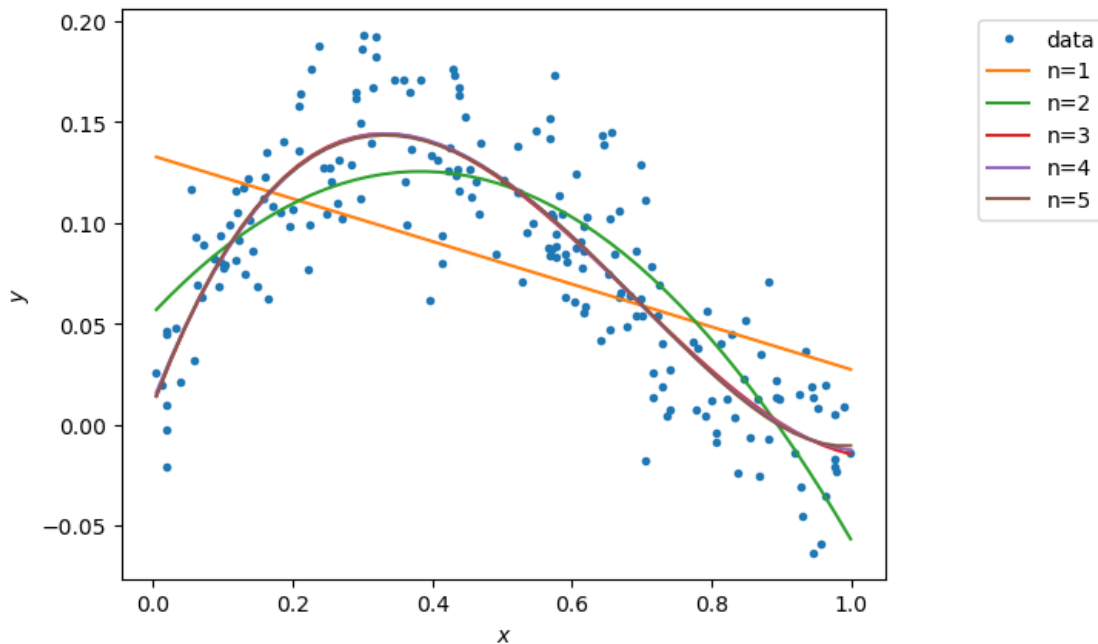
# ===== #
# END YOUR CODE HERE #
# ===== #
```

```
[24]: # Plot the data
f = plt.figure()
ax = f.gca()
ax.plot(x, y, '.')
ax.set_xlabel('$x$')
ax.set_ylabel('$y$')

# Plot the regression lines
plot_xs = []
for i in np.arange(N):
    if i == 0:
        plot_x = np.vstack((np.linspace(min(x), max(x), 50), np.ones(50)))
    else:
        plot_x = np.vstack((plot_x[-2]**(i+1), plot_x))
    plot_xs.append(plot_x)

for i in np.arange(N):
    ax.plot(plot_xs[i][-2:], thetas[i].dot(plot_xs[i]))

labels = ['data']
[labels.append('n={}'.format(i+1)) for i in np.arange(N)]
bbox_to_anchor=(1.3, 1)
lgd = ax.legend(labels, bbox_to_anchor=bbox_to_anchor)
```



0.1.8 Calculating the training error (5 points)

Here, we'll now calculate the training error of polynomial models of orders 1 to 5.

```
[25]: training_errors = []

# ===== #
# START YOUR CODE HERE #
# ===== #

# GOAL: create a variable training_errors, a list of 5 elements,
# where training_errors[i] are the training loss for the polynomial fit of
# order i+1.
for i in range(5):
    y_predict = thetas[i]@xhats[i]

    my_errors = 0

    for j in range(len(y)):
        error_ = ((y[j] - y_predict[j])**2)
        my_errors = my_errors + error_

    training_errors.append(my_errors)

# ===== #
# END YOUR CODE HERE #
# ===== #

print ('Training errors are: \n', training_errors)
```

Training errors are:

```
[0.4759922176725402, 0.21849844418537054, 0.16339207602210745,
0.16330707470593958, 0.16322958391050585]
```

0.1.9 QUESTIONS

- (1) What polynomial has the best training error?
- (2) Why is this expected?

0.1.10 ANSWERS

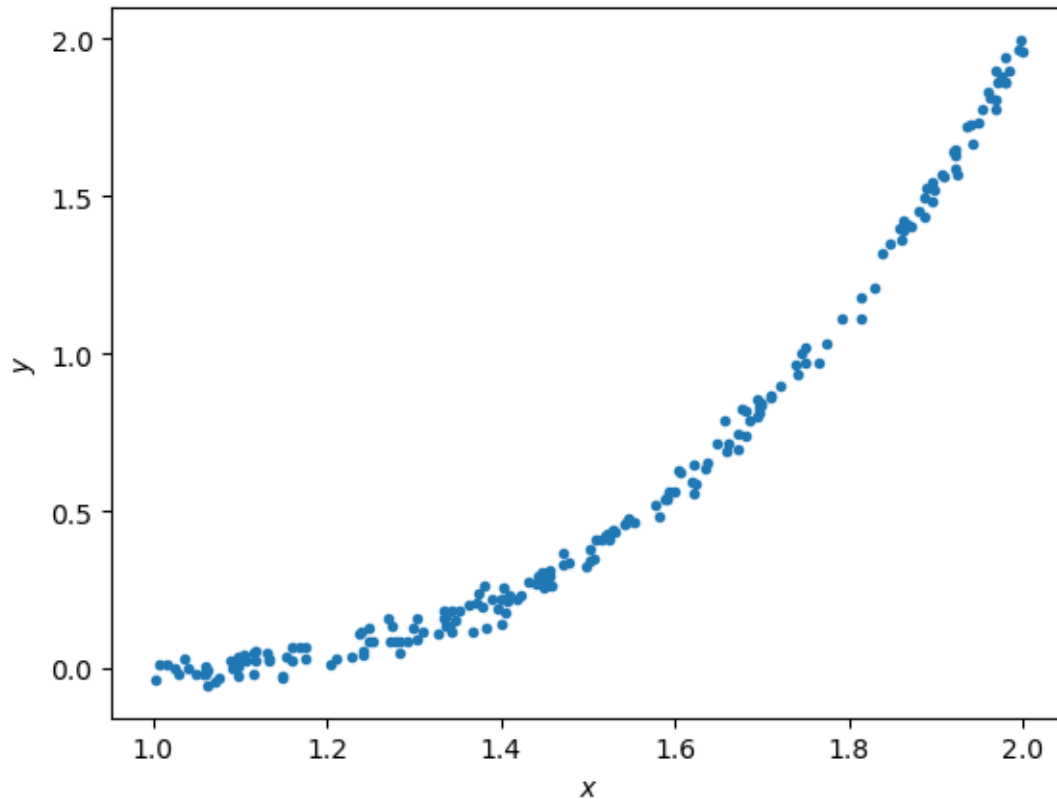
- (1) 5th order polynomial has the best training error.
- (2) This is expected because we can fit more data training data points with higher order polynomial models.

0.1.11 Generating new samples and testing error (5 points)

Here, we'll now generate new samples and calculate testing error of polynomial models of orders 1 to 5.


```
[26]: x = np.random.uniform(low=1, high=2, size=(num_train,))
y = x - 2*x**2 + x**3 + np.random.normal(loc=0, scale=0.03, size=(num_train,))
f = plt.figure()
ax = f.gca()
ax.plot(x, y, '.')
ax.set_xlabel('$x$')
ax.set_ylabel('$y$')
```

```
[26]: Text(0, 0.5, '$y$')
```



```
[27]: xhats = []
for i in np.arange(N):
    if i == 0:
        xhat = np.vstack((x, np.ones_like(x)))
        plot_x = np.vstack((np.linspace(min(x), max(x), 50), np.ones(50)))
    else:
        xhat = np.vstack((x**(i+1), xhat))
        plot_x = np.vstack((plot_x[-2]**(i+1), plot_x))

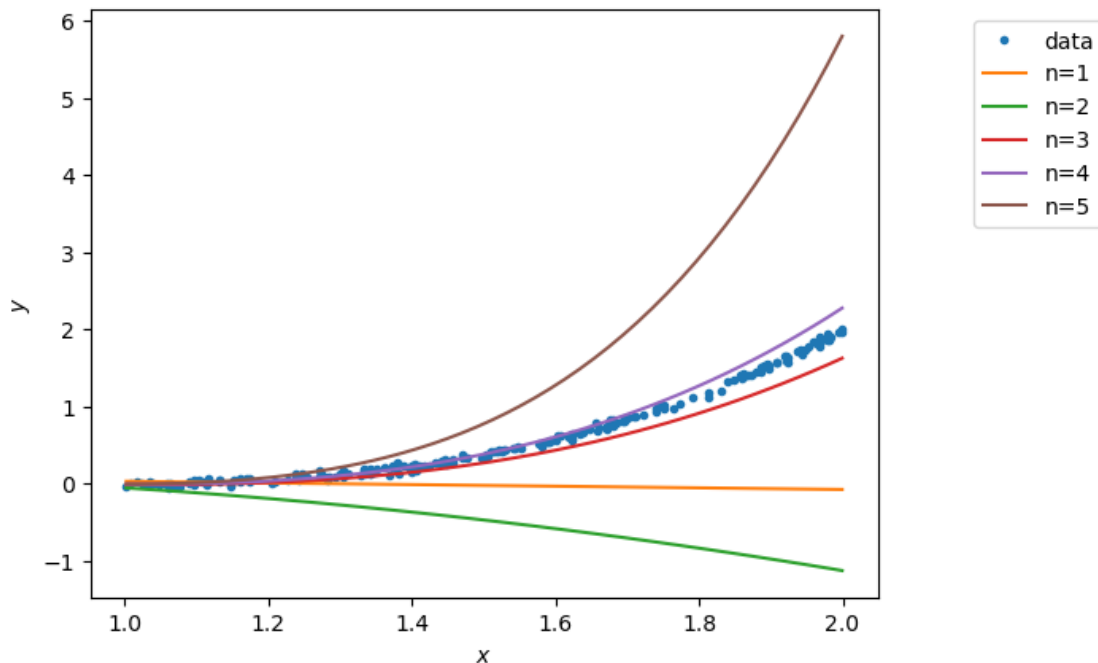
    xhats.append(xhat)
```

```
[28]: # Plot the data
f = plt.figure()
ax = f.gca()
ax.plot(x, y, '.')
ax.set_xlabel('$x$')
ax.set_ylabel('$y$')

# Plot the regression lines
plot_xs = []
for i in np.arange(N):
    if i == 0:
        plot_x = np.vstack((np.linspace(min(x), max(x), 50), np.ones(50)))
    else:
        plot_x = np.vstack((plot_x[-2]**(i+1), plot_x))
    plot_xs.append(plot_x)

for i in np.arange(N):
    ax.plot(plot_xs[i][-2:], thetas[i].dot(plot_xs[i]))

labels = ['data']
[labels.append('n={}'.format(i+1)) for i in np.arange(N)]
bbox_to_anchor=(1.3, 1)
lgd = ax.legend(labels, bbox_to_anchor=bbox_to_anchor)
```



```
[29]: testing_errors = []

# ===== #
# START YOUR CODE HERE #
# ===== #

# GOAL: create a variable testing_errors, a list of 5 elements,
# where testing_errors[i] are the testing loss for the polynomial fit of order  $i+1$ .
for i in range(5):
    y_predict = thetas[i]@xhats[i]
    my_errors = 0
    for j in range(len(y)):
        error_ = ((y[j] - y_predict[j])**2)
        my_errors = my_errors + error_
    testing_errors.append(my_errors)

# ===== #
# END YOUR CODE HERE #
# ===== #

print ('Testing errors are: \n', testing_errors)
```

Testing errors are:

```
[161.72330369101164, 426.38384890115776, 6.251394216552787, 2.3741530378949403,
429.82043635305246]
```

0.1.12 QUESTIONS

- (1) What polynomial has the best testing error?
- (2) Why polynomial models of orders 5 does not generalize well?

0.1.13 ANSWERS

- (1) 4th order polynomial has the best testing error.
- (2) Because it has the issue of overfitting. It fits the training data well but instead loses generality since it can't predict unseen data well.

[]: