

Due Monday, 20 Jan 2025, by 11:59pm to Gradescope.  
Covers material up to Introduction to machine learning refresher 1.  
100 points total.

**1. (25 points) Linear algebra refresher.**

- (a) (12 points) Let  $\mathbf{Q}$  be a real orthogonal matrix.
  - i. (3 points) Show that  $\mathbf{Q}^T$  and  $\mathbf{Q}^{-1}$  are also orthogonal.
  - ii. (3 points) Show that  $\mathbf{Q}$  has eigenvalues with norm 1.
  - iii. (3 points) Show that the determinant of  $\mathbf{Q}$  is either +1 or -1.
  - iv. (3 points) Show that  $\mathbf{Q}$  defines a length preserving transformation.
- (b) (8 points) Let  $\mathbf{A}$  be a matrix.
  - i. (4 points) What is the relationship between the singular vectors of  $\mathbf{A}$  and the eigenvectors of  $\mathbf{A}\mathbf{A}^T$ ? What about  $\mathbf{A}^T\mathbf{A}$ ?
  - ii. (4 points) What is the relationship between the singular values of  $\mathbf{A}$  and the eigenvalues of  $\mathbf{A}\mathbf{A}^T$ ? What about  $\mathbf{A}^T\mathbf{A}$ ?
- (c) (5 points) True or False. Partial credit on an incorrect solution may be awarded if you justify your answer.
  - i. Every linear operator in an  $n$ -dimensional vector space has  $n$  distinct eigenvalues.
  - ii. A non-zero sum of two eigenvectors of a matrix  $\mathbf{A}$  is an eigenvector.
  - iii. If a matrix  $\mathbf{A}$  has the positive semidefinite property, i.e.,  $\mathbf{x}^T\mathbf{A}\mathbf{x} \geq 0$  for all  $\mathbf{x}$ , then its eigenvalues must be non-negative.
  - iv. The rank of a matrix can exceed the number of distinct non-zero eigenvalues.
  - v. A non-zero sum of two eigenvectors of a matrix  $\mathbf{A}$  corresponding to the same eigenvalue  $\lambda$  is always an eigenvector.

**2. (25 points) Probability refresher.**

- (a) (10 points) A and B are involved in a duel. The rules of the duel are that they are to pick up their guns and shoot at each other simultaneously. If one or both are hit, then the duel is over. If both shots miss, then they repeat the process. Suppose that the results of the shots are independent and that each shot of A will hit B with probability  $p_A$  and each shot of B will hit A with probability  $p_B$ . What is:
  - i. (2 points) the probability that A is not hit?
  - ii. (2 points) the probability that both duelists are hit?
  - iii. (2 points) the probability that the duel ends after the  $n^{th}$  round of shots?

- iv. (2 points) the conditional probability that the duel ends after the  $n^{th}$  round of shots given that A is not hit?
- v. (2 points) the conditional probability that the duel ends after the  $n^{th}$  round of shots given that both duelists are hit?
- (b) (5 points) Let  $X$  be a binary signal, such that  $P(X = +1) = P(X = -1) = 0.5$ . Suppose  $X$  is sent across a noisy channel, with noise  $N$  modeled by a zero-mean Gaussian distribution with variance  $\sigma^2$ , where the noise is independent of the signal that was sent. The received signal is  $Y = X + N$ .
- (2 points) Find the conditional PDFs of  $Y$  given both  $\{X = +1\}$  and  $\{X = -1\}$ .
  - (2 points) Suppose a detector compares the received signal  $Y$  to a fixed threshold  $\gamma$  to decide which signal was sent. Specifically, the detector decides that  $X = +1$  was sent if  $Y \geq \gamma$  and that  $X = -1$  was sent otherwise. For a given threshold  $\gamma$ , express the probability of error in terms of the  $\Phi$  function,  $\gamma$ , and  $\sigma$ . Recall, if  $Z \sim \mathcal{N}(0, 1)$ , then

$$P(Z \leq z) = \Phi(z)$$

- (1 point) What is the optimal value of  $\gamma$  to minimize the probability of error? For that optimal value of  $\gamma$ , what is the probability of error? Leave your answer for the optimal value of  $\gamma$  as a real number and the probability of error in terms of the  $\Phi$  function and  $\sigma$ .
- (c) (5 points) There is a screening test for lung cancer that looks at the level of LSA (lung specific antigen) in the blood. There are a number of reasons besides lung cancer that a man can have elevated LSA levels. In addition, many types of lung cancer develop so slowly that they are never a problem. Unfortunately, there is currently no test to distinguish the different types and using the test is controversial because it's hard to quantify the accuracy rates and the harm done by false positives. For this problem, we will call a positive test a true positive if it catches a dangerous type of lung cancer. Also, we will assume the following numbers:

- Rate of dangerous type of lung cancer among men over 30 = 0.0005
- True positive rate for the test = 0.9
- False positive rate for the test = 0.01

Suppose you randomly select a man over 30 and perform a screening test.

- (3 points) What is the probability that the man has a dangerous type of the disease given that he had a positive test?
  - (2 points) What is the probability that the man has a dangerous type of the disease given that he had a negative test?
- (d) (5 points) A family with three daughters and three sons needs to go to the grocery store. Besides the father, who is driving the car, exactly three of the children can come along to the grocery store with him. Suppose that the three children to join the father are chosen randomly, and all such choices are equally likely. Let  $X$  denote the number of daughters who accompany the father to the grocery. Let  $X_i = 1$  if the  $i^{th}$  child who joins the father is a girl, and  $X_i = 0$  otherwise.

- i. (1 point) Find  $Cov(X_i, X_i) = E(X_i X_i) - E(X_i)E(X_i)$ , for  $1 \leq i \leq 3$ . Leave your answer as a fraction.
- ii. (2 points) Find  $Cov(X_i, X_j) = E(X_i X_j) - E(X_i)E(X_j)$ , for  $1 \leq i, j \leq 3$  and  $i \neq j$ . Leave your answer as a fraction.
- iii. (2 points) Find  $Var(X)$ . Leave your answer as a fraction.

3. (10 points) **Multivariate derivatives.**

- (a) (1 points) Let  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{y} \in \mathbb{R}^m$ , and  $\mathbf{A} \in \mathbb{R}^{n \times m}$ . What is  $\nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{y}$ ?
- (b) (1 points) Let  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{y} \in \mathbb{R}^m$ , and  $\mathbf{A} \in \mathbb{R}^{n \times m}$ . What is  $\nabla_{\mathbf{y}} \mathbf{x}^T \mathbf{A} \mathbf{y}$ ?
- (c) (1 points) Let  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{y} \in \mathbb{R}^m$ , and  $\mathbf{A} \in \mathbb{R}^{n \times m}$ . What is  $\nabla_{\mathbf{A}} \mathbf{x}^T \mathbf{A} \mathbf{y}$ ?
- (d) (1 points) Let  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , and let  $f = \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x}$ . What is  $\nabla_{\mathbf{x}} f$ ?
- (e) (1 points) Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{n \times n}$  and  $f = \text{tr}(\mathbf{AB})$ . What is  $\nabla_{\mathbf{A}} f$ ?
- (f) (2 points) Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{n \times n}$  and  $f = \text{tr}(\mathbf{BA} + \mathbf{A}^T \mathbf{B} + \mathbf{A}^2 \mathbf{B})$ . What is  $\nabla_{\mathbf{A}} f$ ?
- (g) (3 points) Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{n \times n}$  and  $f = \|\mathbf{A} + \lambda \mathbf{B}\|_F^2$ . What is  $\nabla_{\mathbf{A}} f$ ?

4. (10 points) **Deriving least-squares with matrix derivatives.**

In least-squares, we seek to estimate some multivariate output  $\mathbf{y}$  via the model

$$\hat{\mathbf{y}} = \mathbf{W} \mathbf{x}$$

In the training set we're given paired data examples  $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$  from  $i = 1, \dots, n$ . Least-squares is the following quadratic optimization problem:

$$\min_{\mathbf{W}} \quad \frac{1}{2} \sum_{i=1}^n \left\| \mathbf{y}^{(i)} - \mathbf{W} \mathbf{x}^{(i)} \right\|^2$$

Derive the optimal  $\mathbf{W}$ .

Where  $\mathbf{W}$  is a matrix, and for each example in the training set, both  $\mathbf{x}^{(i)}$  and  $\mathbf{y}^{(i)}$   $\forall i = 1, \dots, n$  are vectors.

Hint: you may find the following derivatives useful:

$$\begin{aligned} \frac{\partial \text{tr}(\mathbf{WA})}{\partial \mathbf{W}} &= \mathbf{A}^T \\ \frac{\partial \text{tr}(\mathbf{WAW}^T)}{\partial \mathbf{W}} &= \mathbf{WA}^T + \mathbf{WA} \end{aligned}$$

5. (10 points) **Regularized least squares**

In lecture, we worked through the following least squares problem

$$\arg \min_{\theta} \frac{1}{2} \sum_{i=1}^N (y^{(i)} - \theta^T \hat{\mathbf{x}}^{(i)})^2$$

However, the least squares has a tendency to overfit the training data. One common technique used to address the overfitting problem is regularization. In this problem, we work through one of the regularization techniques namely ridge regularization which is also known as the regularized least squares problem. In the regularized least squares we solve the following optimization problem

$$\arg \min_{\theta} \frac{1}{2} \sum_{i=1}^N (y^{(i)} - \theta^T \hat{\mathbf{x}}^{(i)})^2 + \frac{\lambda}{2} \|\theta\|_2^2$$

where  $\lambda$  is a tunable regularization parameter. From the above cost function it can be observed that we are seeking least squares solution with a smaller 2-norm. Derive the solution to the regularized least squares problem, i.e Find  $\theta^*$ .

6. (20 points) **Linear regression.**

Complete the Jupyter notebook `linear_regression.ipynb`. Print out the Jupyter notebook as a PDF and submit it to Gradescope.

1) (a) i. Square matrix is orthogonal, if & only if:

$$Q^T Q = Q Q^T = I$$

$$Q^T Q = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I$$

If orthogonal,  $Q^T = Q^{-1}$ , where  $Q Q^{-1} = I$

$$Q^{-1} = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} = \frac{1}{\det(Q)} \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} = -1 \cdot \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$$

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

ii. We know for an orthogonal matrix  $Q$ ,  $Q^T Q = I$

Let  $\lambda$  be the eigenvalue of  $Q$  &  $v$  be the corresponding eigenvector, so  $Qv = \lambda v$

$$\|Qv\| = \|\lambda v\|$$

$$\|v\| = \|\lambda v\|, \text{ if } v \neq 0, \text{ then } \|v\| \neq 0, \text{ so } \|\lambda\| = 1$$

$$\text{Eg: } \|Q\| = \sqrt{(1)^2 + (-1)^2} = 1$$

$\therefore Q$  has  $\lambda$ s with norm 1

iii.  $\det(Q^T Q) = \det(I)$ , where  $\det(I) = 1$ .

Can also be written as,  $\det(Q^T) \cdot \det(Q) = \det(I)$

Since  $\det(Q^T) = \det(Q)$ , we get  $(\det(Q))^2 = 1$

$$\therefore \det(Q) = \pm 1$$

iv. Let's transform  $Q$  by  $v$ . The norm of  $v$  is defined by  $\|v\| = \sqrt{v^T v}$ .

$$\begin{aligned}\|Qv\| &= \sqrt{(Qv)^T (Qv)} = \sqrt{(Q^T v^T) \cdot (Q \cdot v)} \\ &= \sqrt{(Q^T \cdot Q) (v^T v)} \\ &= \sqrt{I \cdot (v^T v)} \\ &= \sqrt{v^T v}\end{aligned}$$

Thus  $\|Qv\| = \|v\|$ , prove that  $Q$  preserves length.

For any orthogonal matrix, its eigenvalues are unitary (e.g. have length 1). The eigenvalues are placed on the unit circle.

b) i. Any  $n \times m$  matrix  $A$  can be written as or decompose into

$$A = U \Sigma V^T,$$

where

$$U = \text{eigenvectors of } AA^T \quad n \times n \quad (\text{left singular vectors})$$

$$\Sigma = \sqrt{\text{diag}(\text{eig}(AA^T))} \quad n \times m \quad (\text{singular values})$$

$$V = \text{eigenvectors of } A^T A \quad m \times m \quad (\text{right singular vectors})$$

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

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

$$= U (\Sigma^2) V^T$$

The eigenvectors of  $AA^T$  are column of  $U$  and  
the eigenvalues are the square of the singular values of  $A$ .

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

$$= V (\Sigma^2) V^T$$

The eigenvectors of  $A^T A$  are the column of  $V$  and  
the eigenvalues are the square of the singular values of  $A$ .

ii. The singular values of matrix A are the square roots of the non-zero eigenvalues of both  $A^T A$  &  $A A^T$ .

Both  $A^T A$  &  $A A^T$  share the same set of eigenvalues, since the singular values are shared.

If  $\sigma_1, \sigma_2, \dots$  are singular values of A, then the eigenvalues of  $A^T A$  and  $A A^T$  are  $\sigma_1^2, \sigma_2^2, \dots$ .

c) i. False, only if the vector spaces are linearly independent (e.g. I matrix has  $\lambda=1$ )

ii. False, Only true if two eigenvectors correspond to the same eigenvalue.

iii. True, A pos. Semi-dif.  $\rightarrow \text{eig} \left( \frac{A+A^H}{2} \right) \geq 0$

iv. True, The rank of a matrix is dependent on the number of linearly independent rows or columns.  
 $A = \text{diag}(1, 1, 1) \rightarrow \text{Rank } 3$ , but only 1 distinct eigenvalue.

v. True, If the vectors are linearly independent, their sum will be non-zero.

2. a) i. Since  $n=1, 2, 3, 4, \dots$ , where n denotes the ground in which the duel ends, partitions the sample space.

S					
$n=1$	$n=2$	$n=3$	$n=4$	$\dots$	$n=\infty$

Events:

$$P_1 = A \text{ hits } B \quad P(A \text{ not hit}) = \sum_{n=1}^{\infty} P(A \text{ not hit}, n)$$

$P_2 = B \text{ hits } A$  By Bernoulli trial, outcomes are defined by success or failure:

failure  $\rightarrow 1-P$

success  $\rightarrow P$

Duel modeled as a infinite sum of geometric series

$$= [(1-P_1)(1-P_2)]^{(n-1)} P_1 (1-P_2)$$

$$= P_1 (1-P_1)^{n-1} (1-P_2)^n$$

$$P(A \text{ not hit}) = P_1 \sum_{n=1}^{\infty} (1-P_1)^{(n-1)} (1-P_2)^n$$

$$= \frac{P_1(1-P_2)}{1 - (1-P_1)(1-P_2)}$$

ii. Following a similar setup

$$\begin{aligned} P(\text{both hit}) &= \sum_{n=1}^{\infty} [(1-P_1)(1-P_2)]^{n-1} P_1 P_2 \\ &= \frac{P_1 P_2}{1 - (1-P_1)(1-P_2)} \end{aligned}$$

iii. Sample space for game ending in  $n$  rounds

A hit	B hit	Both hit
----------	----------	-------------

S

There may the duel can end  
 $\rightarrow$  A hits  
 $\rightarrow$  B hits  
 $\rightarrow$  A  $\geq$  B hit

$$\begin{aligned} P(\text{game ends in } n \text{ rounds}) &= P(n, \text{A hit}) + P(n, \text{B hit}) + P(n, \text{both hit}) \\ &= [(1-P_1)(1-P_2)]^{(n-1)} (1-P_2) P_2 + \\ &\quad [(1-P_1)(1-P_2)]^{(n-1)} P_1 (1-P_2) + \\ &\quad [(1-P_1)(1-P_2)]^{(n-1)} P_1 P_2 \\ &= [(1-P_1)(1-P_2)]^{(n-1)} [1 - (1-P_1)(1-P_2)] \end{aligned}$$

iv. E: duel ends after  $n^{\text{th}}$  round

$$P(E \mid \text{A not hit}) = \frac{P(\text{A not hit} \mid E) P(E)}{P(\text{A not hit})}$$

$$\text{we know, } P(\text{A not hit}) = \frac{P_1(1-P_2)}{1 - (1-P_1)(1-P_2)},$$

$$\text{and } P(E) = [(1-P_1)(1-P_2)]^{(n-1)} [1 - (1-P_1)(1-P_2)]$$

$$P(E \text{ and A not hit}) = [(1-P_1)(1-P_2)]^{(n-1)} [1 - (1-P_1)(1-P_2)] \times \frac{P_1(1-P_2)}{1 - (1-P_1)(1-P_2)}$$

$$\text{let } q = (1-P_1)(1-P_2)$$

$$= q_1^{n-1} \left[ \frac{1-p_1}{1-q_1} \right] \times \frac{p_1(1-p_2)}{1-q_1}$$

$$= q_1^{n-1} \times p_1(1-p_2)$$

$$= [(1-p_1)(1-p_2)]^{n-1} \times p_1(1-p_2)$$

$$= p_1(1-p_1)^{n-1} (1-p_2)^n$$

$$P(E \mid A \text{ not hit}) = \frac{p_1(1-p_1)^{n-1} (1-p_2)^n}{(1-p_2)^n} = p_1(1-p_1)^{n-1}$$

V. E: dual ends after  $n^{\text{th}}$  round

$$P(E \mid \text{both hit}) = \frac{P(\text{both hit} \mid E) P(E)}{P(\text{both hit})}$$

$$\text{we know, } P(\text{both hit}) = \frac{p_1 p_2}{1 - (1-p_1)(1-p_2)},$$

$$\text{and } P(E) = [(1-p_1)(1-p_2)]^{(n-1)} [1 - (1-p_1)(1-p_2)],$$

$$P(\text{both hit and } E) = P(\text{both hit}) \cdot P(E)$$

$$= \frac{p_1 p_2}{1 - (1-p_1)(1-p_2)} \cdot [(1-p_1)(1-p_2)]^{(n-1)} \cdot [1 - (1-p_1)(1-p_2)]$$

$$= p_1 p_2 (1-p_1)^{(n-1)} (1-p_2)^{(n-1)}$$

$$P(E \mid \text{both hit}) = \frac{p_1 p_2 (1-p_1)^{(n-1)} (1-p_2)^{(n-1)}}{p_1 p_2} = (1-p_1)^{(n-1)} (1-p_2)^{(n-1)}$$

b) i.

$$P(X=+1) = P(X=-1) = 0.5 \quad X \rightarrow N \sim N(0, \sigma^2) \rightarrow Y = X + N$$

Conditional = Joint  
Marginal

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad z = \frac{x-\mu}{\sigma}$$

Conditional PDF of  $Y$  given  $\{x = +1\}$  and  $\{x = -1\}$

If  $x = +1$ ,  $y = 1 + N$ ,  $y|x=+1) \sim N(1, \sigma^2)$

$$f(y|x=+1) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(y-1)^2}{2\sigma^2}}$$

If  $x = -1$ ,  $y = -1 + N$ ,  $y|x=-1) \sim N(-1, \sigma^2)$

$$f(y|x=-1) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(y+1)^2}{2\sigma^2}}$$

ii. If  $y \geq \gamma \rightarrow x = +1$

If  $y < \gamma \rightarrow x = -1$

$y \geq \gamma$	$y < \gamma$	$\hookrightarrow$	partition the sample space
-----------------	--------------	-------------------	----------------------------

$$P(\text{Error}) = P(x = +1) \cdot P(y < \gamma | x = +1) + P(x = -1) \cdot P(y \geq \gamma | x = -1)$$

$$\text{We know } P(x = +1) = P(x = -1) = 0.5$$

$$P(\text{Error}) = 0.5 \cdot P(y < \gamma | x = +1) + 0.5 \cdot P(y \geq \gamma | x = -1)$$

Expressing in terms of  $\phi$  function:

$$\begin{aligned} P(y < \gamma | x = +1) + P(y \geq \gamma | x = -1) \\ = 1 \end{aligned}$$

$$P(y < \gamma | x = +1) = P(1 + N < \gamma) = P(N < \gamma - 1) = P\left(\frac{N}{\sigma} < \frac{\gamma - 1}{\sigma}\right) = \phi\left(\frac{\gamma - 1}{\sigma}\right)$$

$$\therefore P(\text{Error}) = 0.5 \cdot \phi\left(\frac{\gamma - 1}{\sigma}\right) + 0.5 \cdot \left[1 - \phi\left(\frac{\gamma + 1}{\sigma}\right)\right]$$

iii. To minimize the probability of error, we should choose a  $\gamma$  that equalizes the error from both cases. In our case, letting  $\gamma = 0$  will be the optimal threshold because the prior probabilities change from  $x = +1$  &  $x = -1$

$$\therefore P(\text{Error}) = 0.5 \cdot \phi\left(\frac{0 - 1}{\sigma}\right) + 0.5 \cdot \left[1 - \phi\left(\frac{0 + 1}{\sigma}\right)\right]$$

$$= 0.5 \phi\left(-\frac{1}{\sigma}\right) + 0.5 \left[ 1 - \phi\left(\frac{1}{\sigma}\right)\right]$$

$$= 1 - \phi\left(\frac{1}{\sigma}\right)$$

c) i)  $P(D) = 0.0005$

S			
$d_1$	$d_2$	...	$d_i$

TPR  $\rightarrow P(T|D) = 0.9$

FPR  $\rightarrow P(T|D^c) = 0.01$

$$P(D) = \sum_{i=1}^n d_i$$

$$P(D|T) = \frac{P(T|D) P(D)}{P(T)}$$

$$P(T) = P(T|D) P(D) + P(T|D^c) P(D^c)$$

$$P(D|T) = \frac{0.9 \times 0.0005}{(0.9 \times 0.0005) + (0.01 \times 0.9995)}$$

$$= 0.04308$$

ii)  $P(D|T^c) = \frac{P(T^c|D) P(D)}{P(T^c|D) P(D) + P(T^c|D^c) P(D^c)}$

$$= \frac{0.1 \times 0.0005}{(0.1 \times 0.0005) + (0.99)(0.0005)}$$

d)  $D = 3$   
 $S = 3$        $x_i = \begin{cases} 1 & \text{if the } i^{\text{th}} \text{ child is a girl} \\ 0 & \text{if the } i^{\text{th}} \text{ child is a boy} \end{cases}$

i.  $Cov(x_i, x_i) = E(x_i x_i) - E(x_i) E(x_i)$ , for  $1 \leq i \leq 3$

$$Cov(x_i, x_i) = Var(x_i) = E(x_i) - (E(x_i))^2 = \nu - \sigma^2$$

but  
case ,  $P(x_i = 1) = \frac{\binom{5}{2}}{\binom{6}{3}} = \frac{1}{2}$

$$Cov(x_i, x_i) = \frac{1}{2} - \left(\frac{1}{2}\right)^2 = \frac{1}{4}$$

ii.  $\text{Cov}(x_i, x_j) = E(x_i x_j) - E(x_i)E(x_j)$ , for  $1 \leq i, j \leq 3$  and  $i \neq j$

$$\text{We know, } E(x_i) = E(x_j) = \frac{1}{2}$$

$$\therefore E(x_i)E(x_j) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$$

$x_i x_j = 1$  occurs when both  $i^{\text{th}}$  &  $j^{\text{th}}$  related children are daughter

$$E(x_i x_j) = P(\text{both } i \text{ & } j \text{ are daughter})$$

$$P(x_i x_j = 1) = \frac{\text{Ways to choose 2 daughters & 1 other child}}{\text{Ways to choose 3 out of 6}}$$

$$P(x_i x_j = 1) = \frac{\binom{3}{2} \binom{3}{1}}{\binom{6}{3}} = \frac{9}{20}$$

$$\therefore \text{Cov}(x_i, x_j) = \frac{9}{20} - \frac{1}{4} = \frac{9-5}{20} = \frac{4}{20} = \frac{1}{5}$$

iii.  $\text{Var}(x)$

$x \rightarrow \text{no. of daughter that accompany the father}$

$$x = x_1 + x_2 + x_3$$

$$\text{Var}(x) = \sum_{i=1}^3 \text{Var}(x_i) + \sum_{1 \leq i < j \leq 3} \text{Cov}(x_i, x_j)$$

$$\text{we know, } \text{Var}(x_i) = \frac{1}{4}$$

$$\therefore \sum_{i=1}^3 \text{Var}(x_i) = \frac{3}{4}$$

$$\text{we know, } \text{Cov}(x_i, x_j) = \frac{1}{5}$$

$$\therefore \sum_{i=1}^3 \text{Cov}(x_i, x_j) = \frac{3}{5}$$

$$\therefore \text{Var}(x) = \frac{3}{4} + \frac{3}{5} = \frac{15+12}{20} = \frac{27}{20}$$

### 3. Multivariate derivative

a)  $x \in \mathbb{R}^n$        $\nabla_x x^T A y$   
 $y \in \mathbb{R}^m$        $(1 \times n) \times (n \times m) \times (m \times 1) \rightarrow \text{Scalar}$   
 $A \in \mathbb{R}^{n \times m}$        $\nabla_x x^T A y \in \mathbb{R}^1$

$$\boxed{\nabla_x x^T A y = A y}$$

b)  $x \in \mathbb{R}^n$   $y \in \mathbb{R}^m$   $A \in \mathbb{R}^{n \times m}$

$$\nabla_y x^T A y = A^T x$$

$$(1 \times n)(n \times m)(m \times 1)$$

$$\nabla_y x^T A y \in \mathbb{R}^1$$

c)  $x \in \mathbb{R}^n$   $y \in \mathbb{R}^m$   $A \in \mathbb{R}^{n \times m}$

$$\nabla_A (x^T A y) = x y^T$$

$$\nabla_A x^T A y = \begin{bmatrix} \frac{\partial x^T A y}{\partial a_{1,1}} & \frac{\partial x^T A y}{\partial a_{1,2}} & \dots & \frac{\partial x^T A y}{\partial a_{1,n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial x^T A y}{\partial a_{m,1}} & \dots & \dots & \frac{\partial x^T A y}{\partial a_{m,n}} \end{bmatrix}$$

d)  $x \in \mathbb{R}^n$   $A \in \mathbb{R}^{n \times n}$   $b \in \mathbb{R}^n$

$\nabla_x f$  where  $f = x^T A x + b^T x$   
 $x^T A x \rightarrow \text{quadratic}$   
differentiating  $x^T A x$  w.r.t  $x \rightarrow A x + A^T$   
assuming  $A$  is symmetric

$$\nabla_x f = A x + A^T x + b$$

e)  $A, B \in \mathbb{R}^{n \times n}$

$$\nabla_A f, \text{ where } f = \text{tr}(AB)$$

$$\text{tr}(AB) = \text{tr}(BA) = \text{tr} \left( \begin{bmatrix} -b_1^T \\ -b_2^T \\ \vdots \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & \dots \\ a_{1,1} & a_{1,2} & a_{1,3} & \dots \\ a_{2,1} & a_{2,2} & a_{2,3} & \dots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix} \right)$$

derivation of  $f$  w.r.t  $A$

$$\nabla_A f = B^T$$

f)  $A, B \in \mathbb{R}^{n \times n}$

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

$$\text{tr}(BA) : \nabla_A = B^T$$

$$\text{tr}(A^T B) = \text{tr}(BA) : \nabla_A = B^T$$

$$\text{tr}(A^2 B) : \nabla_A = AB + BA$$

$$\therefore \nabla_A f = B^T + B + (AB + BA)^T$$

g)  $\nabla_A f$ , where  $f = \|A + \lambda B\|_F^2$

$$A, B \in \mathbb{R}^{n \times n}, \lambda \in \mathbb{R}$$

Frobenius norm :  $\|M\|_F^2 = \text{tr}(M^T M)$

$$\begin{aligned}\|A + \lambda B\|_F^2 &= \text{tr}((A + \lambda B)^T (A + \lambda B)) \\ &= \text{tr}(A^T A) + 2\lambda \text{tr}(A^T B) + \lambda^2 \text{tr}(B^T B)\end{aligned}$$

derivative w.r.t. A  $A^T A \rightarrow 2A$

derivative w.r.t. A  $A^T B \rightarrow B$

$$\boxed{\nabla_A f = 2(A + \lambda B)}$$

4. Deriving least squares with matrix derivatives

$\hat{y} = Wx$ , where W is a matrix

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

$$\begin{aligned}J &= \frac{1}{2} \sum_{i=1}^n (y^{(i)} - Wx^{(i)})^T (y^{(i)} - Wx^{(i)}) \\ &= \frac{1}{2} \sum_{i=1}^n (y^{(i)T} y^{(i)} - y^{(i)T} W x^{(i)} - W x^{(i)T} y^{(i)} + W x^{(i)T} W x^{(i)}) \\ &= \frac{1}{2} \sum_{i=1}^n (y^{(i)T} y^{(i)} - 2 y^{(i)T} W x^{(i)} + x^{(i)T} W^T W x^{(i)})\end{aligned}$$

Let,  $X = [x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(n)}] \rightarrow \text{data}$

$y = [y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(n)}] \rightarrow \text{output}$ , be matrix

$$J = \frac{1}{2} (\text{tr}(y^T y) - 2 \text{tr}(y^T W X) + \text{tr}(W X X^T W^T))$$

$$\text{we know, } \frac{\partial \text{tr}(WA)}{\partial W} = A^T, \quad \frac{\partial \text{tr}(WAW^T)}{\partial W} = WA^T + WA$$

$$\nabla_W J = 0 - 2 y^T X^T + y^T X^T + W X X^T = 0$$

$$\nabla_{\mathbf{w}} \mathcal{L} = -\mathbf{y} \mathbf{x}^T + \mathbf{w} \mathbf{x} \mathbf{x}^T = 0$$

$$\mathbf{w} \mathbf{x} \mathbf{x}^T = \mathbf{y} \mathbf{x}^T$$

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

## 5. Regularized least squares

$$\arg \min_{\theta} \frac{1}{2} \sum_{i=1}^N (y^{(i)} - \theta^T \hat{x}^{(i)})^2 + \frac{\lambda}{2} \|\theta\|_2^2$$

$$\mathcal{L} = \frac{1}{2} \sum_{i=1}^N (y^{(i)} - \theta^T \hat{x}^{(i)})^T (y^{(i)} - \theta^T \hat{x}^{(i)}) + \frac{\lambda}{2} \|\theta\|_2^2$$

$$= \frac{1}{2} \sum_{i=1}^N (y^{(i)T} y^{(i)} - y^{(i)T} \theta^T \hat{x}^{(i)} - \theta^T \hat{x}^{(i)T} y^{(i)} + \theta^T \hat{x}^{(i)T} \theta) + \frac{\lambda}{2} \|\theta\|_2^2$$

$$= \frac{1}{2} \sum_{i=1}^N (y^{(i)T} y^{(i)} - 2 y^{(i)T} \hat{x}^{(i)} \theta + \theta^T \hat{x}^{(i)T} \hat{x} \theta) + \frac{\lambda}{2} \|\theta\|_2^2$$

Let,  $\mathbf{X} = [x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(i)}] \rightarrow \text{data}$

$\mathbf{Y} = [y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(i)}] \rightarrow \text{output, be matrix}$

$$= \frac{1}{2} (\mathbf{Y}^T \mathbf{Y} + 2 \mathbf{Y}^T \mathbf{X} \theta + \theta^T \mathbf{X}^T \mathbf{X} \theta) + \frac{\lambda}{2} \|\theta\|_2^2$$

$$\nabla_{\theta} \mathcal{L} = \frac{1}{N} (-\mathbf{Z} \mathbf{X}^T \mathbf{Y} + \mathbf{Z} \mathbf{X}^T \mathbf{X} \theta) + \lambda \theta$$

$$= -\mathbf{X}^T \mathbf{Y} + \mathbf{X}^T \mathbf{X} \theta + \lambda \theta$$

$$\mathbf{X}^T \mathbf{Y} = \theta (\mathbf{X}^T \mathbf{X} + \lambda I)$$

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



## 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 2025, Prof. J.C. Kao, TAs: B. Qu, K. Pang, S. Dong, S. Rajesh, T. Monsoor, X. Yan

```
In [1]: 1 import numpy as np
          2 import matplotlib.pyplot as plt
          3
          4 #allows matlab plots to be generated in line
          5 %matplotlib inline
```

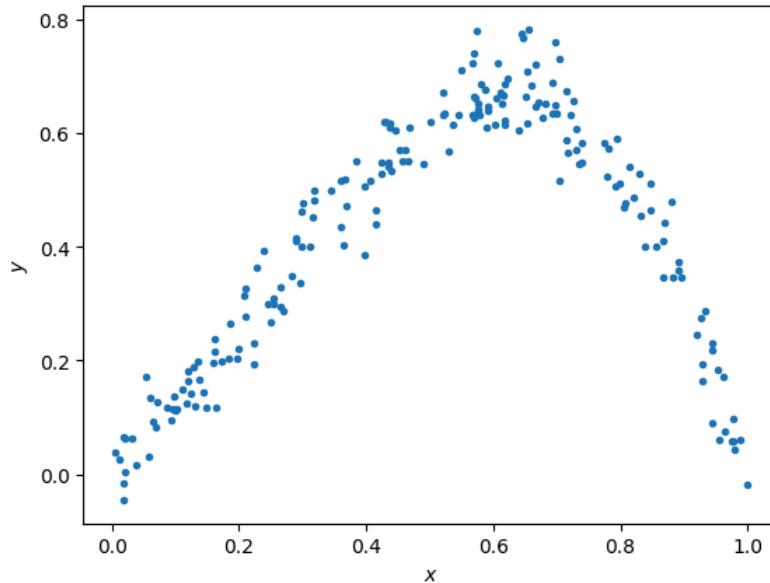
### Data generation

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

```
In [2]: 1 np.random.seed(0) # Sets the random seed,
2 num_train = 200      # Number of training data points
3
4 # Generate the training data
5 x = np.random.uniform(low=0, high=1, size=(num_train,))
6 y = x + 2*x**2 - 3*x**3 + np.random.normal(loc=0, scale=0.05, size=(num_train,))
7 f = plt.figure()
8 ax = f.gca()
9 ax.plot(x, y, '.')
10 ax.set_xlabel('$x$')
11 ax.set_ylabel('$y$')
```

Out [2]: Text(0, 0.5, '\$y\$')

```
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "orientation" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "dpi" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "facecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "edgecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "bbox_inches_restore" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
```



## 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  $\epsilon$ ?

## ANSWERS:

- (1) Values of  $x$  are sampled from a uniform distribution between 0 and 1.
- (2) The additive noise values are derived from a normal distribution with mean 0 and standard deviation 0.05.

## Fitting data to the model (5 points)

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

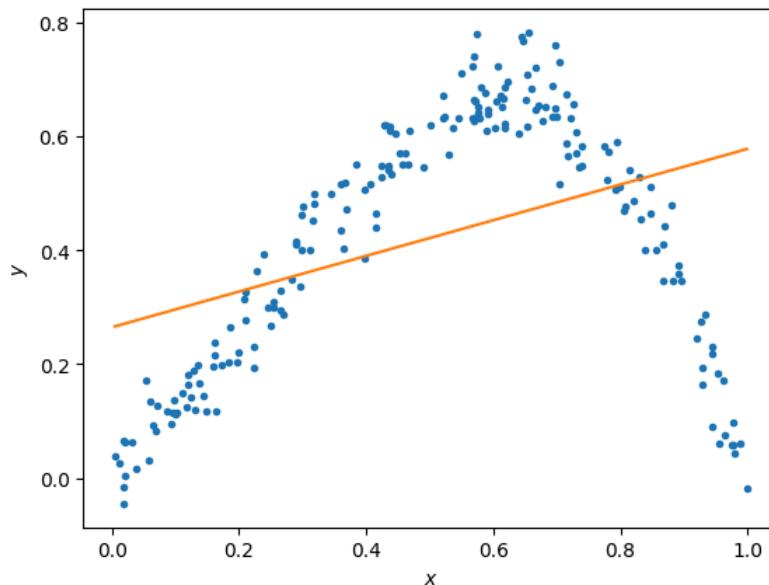
```
In [3]: 1 #xhat = (x, 1)
2 xhat = np.vstack((x, np.ones_like(x)))
3
4 # Projecting y onto the column space of xhat
5 theta = np.linalg.inv(xhat.dot(xhat.T)).dot(xhat).dot(y)
6 print(theta)
7
```

[0.31325736 0.26474646]

```
In [4]: 1 # Plot the data and your model fit.
2 f = plt.figure()
3 ax = f.gca()
4 ax.plot(x, y, '.')
5 ax.set_xlabel('$x$')
6 ax.set_ylabel('$y$')
7
8 # Plot the regression line
9 xs = np.linspace(min(x), max(x), 50)
10 xs = np.vstack((xs, np.ones_like(xs)))
11 plt.plot(xs[0,:], theta.dot(xs))
```

Out[4]: [<matplotlib.lines.Line2D at 0x107703e90>]

```
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "orientation" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "dpi" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "facecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "edgecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "bbox_inches_restore" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
```



## QUESTIONS

- (1) Does the linear model under- or overfit the data?

(2) How to change the model to improve the fitting?

## ANSWERS

(1) The linear model underfits the data. Looking at the generated data and the true model ( $y = x + 2x^2 - 3x^3 + \epsilon$ ), we can see the relationship is cubic, but we're trying to fit it with just a straight line. The linear model is too simple to capture the underlying nonlinear relationship.

(2) We should use a higher-order polynomial model that can capture the nonlinear relationship. Since the true relationship is cubic, at minimum a third-order polynomial model would be appropriate. This would allow the model to capture both the linear and nonlinear components of the relationship.

### Fitting data to the model (5 points)

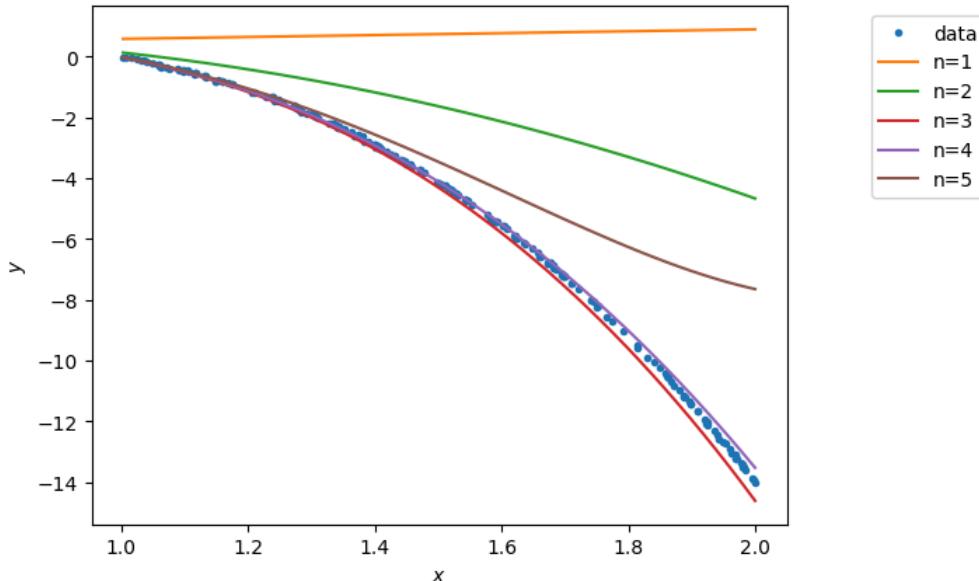
Here, we'll now do regression to polynomial models of orders 1 to 5. Note, the order 1 model is the linear model you prior fit.

```
In [5]:  
1  
2 N = 5  
3 xhats = []  
4 thetas = []  
5  
6 # For each polynomial order 1 to 5  
7 for i in range(N):  
8     # Create feature matrix for current polynomial order  
9     if i == 0:  
10         # First order: [x, 1]  
11         xhat = np.vstack((x, np.ones_like(x)))  
12     else:  
13         # Higher orders: [x^n, x^(n-1), ..., x, 1]  
14         xhat = np.vstack((x**(i+1), xhat))  
15  
16 xhats.append(xhat)  
17  
18 # Model Coefficients  
19 theta = np.linalg.inv(xhat.dot(xhat.T)).dot(xhat).dot(y)  
20 thetas.append(theta)  
21  
22 pass  
23  
24 print(thetas)  
25  
[array([0.31325736, 0.26474646]), array([-2.54077983, 2.81808862, -0.14765329]), array([-3.19269866, 2.29870971, 0.8529672, 0.01631781]), array([ 0.23466729, -3.65889518, 2.59581785, 0.78750808, 0.01958868]), array([ 0.87387653, -1.94094667, -1.73245937, 1.87616209, 0.88956527, 0.01619333])]
```

In [12]:

```
1 # Plot the data
2 f = plt.figure()
3 ax = f.gca()
4 ax.plot(x, y, '.')
5 ax.set_xlabel('$x$')
6 ax.set_ylabel('$y$')
7
8 # Plot the regression lines
9 plot_xs = []
10 for i in np.arange(N):
11     if i == 0:
12         plot_x = np.vstack((np.linspace(min(x), max(x), 50), np.ones(50)))
13     else:
14         plot_x = np.vstack((plot_x[-2]***(i+1), plot_x))
15     plot_xs.append(plot_x)
16
17 for i in np.arange(N):
18     ax.plot(plot_xs[i][-2,:], thetas[i].dot(plot_xs[i]))
19
20 labels = ['data']
21 [labels.append('n={}'.format(i+1)) for i in np.arange(N)]
22 bbox_to_anchor=(1.3, 1)
23 lgd = ax.legend(labels, bbox_to_anchor=bbox_to_anchor)
```

```
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "orientation" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "dpi" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "facecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "edgecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "bbox_inches_restore" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
```



### Calculating the training error (5 points)

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

```
In [13]: 1 training_errors = []
2
3
4 for i in np.arange(N):
5     y_pred = thetas[i].dot(xhats[i])
6     mse = np.mean((y - y_pred) ** 2)
7
8     training_errors.append(mse)
9
10    pass
11
12 print ('Training errors are: \n', training_errors)
```

```
Training errors are:
[54.246317716012236, 18.911159217529747, 0.08693862570546057, 0.03042966869910436, 5.954356339341897]
```

## QUESTIONS

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

## ANSWERS

- (1) The fifth order polynomial has the best training error.
- (2) Using a higher-order polynomial helps the model capture more relationships compared to lower order ones. Therefore, we get a low training error for the fifth order polynomial.

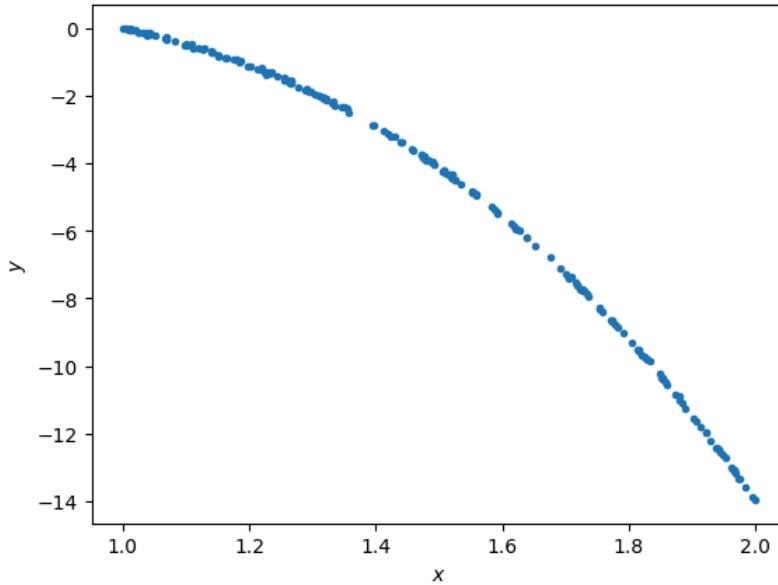
## Generating new samples and testing error (5 points)

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

```
In [14]: 1 x = np.random.uniform(low=1, high=2, size=(num_train,))
2 y = x + 2*x**2 - 3*x**3 + np.random.normal(loc=0, scale=0.03, size=(num_train,))
3 f = plt.figure()
4 ax = f.gca()
5 ax.plot(x, y, '.')
6 ax.set_xlabel('$x$')
7 ax.set_ylabel('$y$')
```

Out[14]: Text(0, 0.5, '\$y\$')

```
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "orientation" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "dpi" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "facecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "edgecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "bbox_inches_restore" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
```

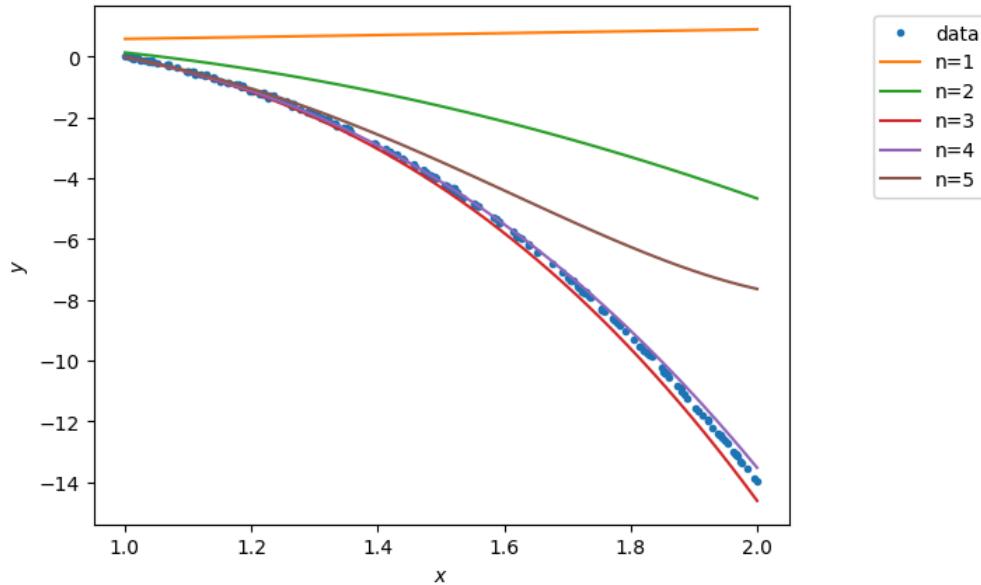


```
In [15]: 1 xhats = []
2 for i in np.arange(N):
3     if i == 0:
4         xhat = np.vstack((x, np.ones_like(x)))
5         plot_x = np.vstack((np.linspace(min(x), max(x), 50), np.ones(50)))
6     else:
7         xhat = np.vstack((x***(i+1), xhat))
8         plot_x = np.vstack((plot_x[-2]***(i+1), plot_x))
9
10    xhats.append(xhat)
```

In [16]:

```
1 # Plot the data
2 f = plt.figure()
3 ax = f.gca()
4 ax.plot(x, y, '.')
5 ax.set_xlabel('$x$')
6 ax.set_ylabel('$y$')
7
8 # Plot the regression lines
9 plot_xs = []
10 for i in np.arange(N):
11     if i == 0:
12         plot_x = np.vstack((np.linspace(min(x), max(x), 50), np.ones(50)))
13     else:
14         plot_x = np.vstack((plot_x[-2]***(i+1), plot_x))
15     plot_xs.append(plot_x)
16
17 for i in np.arange(N):
18     ax.plot(plot_xs[i][-2,:], thetas[i].dot(plot_xs[i]))
19
20 labels = ['data']
21 [labels.append('n={}'.format(i+1)) for i in np.arange(N)]
22 bbox_to_anchor=(1.3, 1)
23 lgd = ax.legend(labels, bbox_to_anchor=bbox_to_anchor)
```

```
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "orientation" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "dpi" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "facecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "edgecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "bbox_inches_restore" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
```



```
In [17]: 1 testing_errors = []
2
3
4
5 for i in np.arange(N):
6
7     y_pred = thetas[i].dot(xhats[i])
8
9     mse = np.mean((y - y_pred) ** 2)
10
11    testing_errors.append(mse)
12
13 pass
14
15
16
17 print ('Testing errors are: \n', testing_errors)
```

```
Testing errors are:
[51.595934377138384, 17.88767990466363, 0.08568137176745275, 0.02624925686033019, 5.438380112067071]
```

## QUESTIONS

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

## ANSWERS

- (1) Lower-order polynomials (e.g., orders 1 or 2) may underfit, leading to high training and testing errors. Higher-order polynomials (e.g., order 5) tend to overfit, capturing noise in the training data, which results in poor testing performance. The fourth order polynomial has the best testing error.
- (2) Polynomial models of order 5 may not generalize well because of overfitting, which occurs when the model becomes too flexible and fits the noise in the training data rather than the underlying trend.