

1. (25 points) Linear algebra refresher.

- (a) (12 points) Let  $\mathbf{Q}$  be a real orthogonal matrix.
- (3 points) Show that  $\mathbf{Q}^T$  and  $\mathbf{Q}^{-1}$  are also orthogonal.
  - (3 points) Show that  $\mathbf{Q}$  has eigenvalues with norm 1.
  - (3 points) Show that the determinant of  $\mathbf{Q}$  is either +1 or -1.
  - (3 points) Show that  $\mathbf{Q}$  defines a length preserving transformation.

i. for an orthogonal matrix  $\mathbf{Q}$ ,  $\mathbf{Q} \cdot \mathbf{Q}^T = \mathbf{I}$

$\mathbf{Q}^T$  is an orthogonal matrix since  $\mathbf{Q}^T(\mathbf{Q}^T)^T = \mathbf{Q}^T \mathbf{Q} = \mathbf{I}$

$$(\mathbf{Q}^T)^T \mathbf{Q}^T = \mathbf{Q} \mathbf{Q}^T = \mathbf{I} \checkmark$$

$\mathbf{Q}^{-1}$  is orthogonal because  $\mathbf{Q}^{-1} = \mathbf{Q}^T$  since  $\mathbf{Q} \cdot \mathbf{Q}^T = \mathbf{I}$ , and multiplying a matrix by its inverse is the identity matrix

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

ii.  $\lambda$  is the eigenvalue of  $\mathbf{Q}$ ,  $\mathbf{Q}\mathbf{x} = \lambda\mathbf{x}$

$$(\mathbf{Q}\mathbf{x})^T(\mathbf{Q}\mathbf{x}) = (\lambda\mathbf{x})^T(\lambda\mathbf{x})$$

$$\mathbf{x}^T \mathbf{Q}^T \mathbf{Q} \mathbf{x} = |\lambda|^2 \mathbf{x}^T \mathbf{x}$$

$$\therefore |\lambda|^2 = 1 \quad \text{since } \mathbf{Q}^T \mathbf{Q} = \mathbf{I}$$

Thus  $\mathbf{Q}$  has eigenvalues of norm 1

iii.  $\det(\mathbf{Q}) = \pm 1$   
 $\mathbf{Q}$  is orthogonal

$$\begin{aligned} &\rightarrow \det(\mathbf{Q}^T \cdot \mathbf{Q}) = 1 \\ &= \det(\mathbf{Q}^T) \cdot \det(\mathbf{Q}) \\ &1 = (\det(\mathbf{Q}))^2 \end{aligned}$$

$$\det(\mathbf{Q}) = \det(\mathbf{Q}^T)$$

$$\therefore \boxed{\det(\mathbf{Q}) = \pm 1}$$

iv.  $\mathbf{Q}$  preserves lengths  $\rightarrow \|\mathbf{Q}\mathbf{v}\|^2 = \|\mathbf{v}\|^2$

$$\|\mathbf{Q}\mathbf{v}\|^2 = \mathbf{Q}\mathbf{v} \cdot \mathbf{Q}\mathbf{v} = (\mathbf{Q}\mathbf{v})^T \mathbf{Q}\mathbf{v} = \mathbf{v}^T \mathbf{Q}^T \mathbf{Q}\mathbf{v} = \mathbf{v}^T \mathbf{I} \mathbf{v} = \mathbf{v} \cdot \mathbf{v} = \|\mathbf{v}\|^2$$

(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  $\mathbf{A}^T\mathbf{A}$ ?

i.  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$

$$\begin{aligned}\mathbf{A}\mathbf{A}^T &= (\mathbf{U}\Sigma\mathbf{V}^T)(\mathbf{U}\Sigma\mathbf{V}^T)^T \\ &= \mathbf{U}\Sigma\mathbf{V}^T(\mathbf{V}\Sigma^T\mathbf{U}^T) \\ &= \mathbf{U}\Sigma\Sigma^T\mathbf{U}^T\end{aligned}\quad \begin{aligned}\mathbf{A}^T\mathbf{A} &= (\mathbf{U}\Sigma\mathbf{V}^T)^T(\mathbf{U}\Sigma\mathbf{V}^T) \\ &= \mathbf{V}\Sigma^T\mathbf{U}^T\mathbf{U}\Sigma\mathbf{V}^T \\ &= \mathbf{V}\Sigma^T\Sigma\mathbf{V}^T\end{aligned}$$

The columns of  $\mathbf{U}$  are the left singular vectors of  $\mathbf{A}$  and they are the orthonormal eigenvectors of  $\mathbf{A}\mathbf{A}^T$

The columns of  $\mathbf{V}$  are the right singular vectors of  $\mathbf{A}$  and they are the orthonormal eigenvectors of  $\mathbf{A}^T\mathbf{A}$

ii. The singular values of  $\mathbf{A}$  are the square roots of the eigenvalues of  $\mathbf{A}^T\mathbf{A}$  and  $\mathbf{A}\mathbf{A}^T$

$$\sigma_i(\mathbf{A}) = \lambda_i^{1/2}(\mathbf{A}^T\mathbf{A}) = \lambda_i^{1/2}(\mathbf{A}\mathbf{A}^T)$$

(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.

i. False. should be at most  $n$  distinct eigenvalues

$I_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$  but all eigenvalues are 1  
(only 1 distinct eigenvalue)

ii. False, only if they have the same eigenvalue

Example  $A = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$   $\lambda_1 = 1$ ,  $\lambda_2 = -1$

$A\mathbf{v}_1 = \lambda_1 \mathbf{v}_1$ ,  $A\mathbf{v}_2 = \lambda_2 \mathbf{v}_2$ ,  $\mathbf{v}_1 + \mathbf{v}_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ , which is not an eigenvector of  $A$

iii.  $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$  for all  $\mathbf{x}$ , then eigenvalues are nonnegative

True

$$\begin{aligned} A\mathbf{x} &= \lambda \mathbf{x} \\ \mathbf{x}^T A \mathbf{x} &= \mathbf{x}^T \lambda \mathbf{x} \\ \mathbf{x}^T A \mathbf{x} &= \mathbf{x}^T \mathbf{x} \lambda \geq 0 \end{aligned}$$

$\mathbf{x}^T \mathbf{x}$  is positive, so  $\lambda$  must be positive or 0

iv.  $\text{rank}(A)$  can exceed number of distinct nonzero  $\lambda$

True, the number of nonzero eigenvalues must be  $\leq \text{rank}$

v. True

2. (25 points) Probability refresher.

(a) (5 points) A jar of coins is equally populated with two types of coins. One is type "H50" and comes up heads with probability 0.5. Another is type "H60" and comes up heads with probability 0.6.

- i. (1 points) You take one coin from the jar and flip it. It lands tails. What is the posterior probability that this is an H50 coin?
- ii. (2 points) You put the coin back, take another, and flip it 4 times. It lands T, H, H, H. How likely is the coin to be type H50?
- iii. (2 points) A new jar is now equally populated with coins of type H50, H55, and H60 (with probabilities of coming up heads 0.5, 0.55, and 0.6 respectively). You take one coin and flip it 10 times. It lands heads 9 times. How likely is the coin to be of each possible type?

a. Event T is the coin lands tails, H for heads

Two coins H50, H60

$$\begin{aligned} P(H|H50) &= 0.5 & P(T) &= P(T|H50)P(H50) \\ P(H|H60) &= 0.6 & & + P(T|H60)P(H60) \\ & & & = 0.5 \cdot 0.5 + 0.4 \cdot 0.5 = 0.45 \end{aligned}$$

i.  $P(H50|T) = ?$

$$\begin{aligned} &= \frac{P(T|H50)P(H50)}{P(T)} \\ &= \frac{0.5 \cdot 0.5}{0.45} = \boxed{\frac{5}{9}} \end{aligned}$$

ii.  $P(H50|THHH) = \frac{P(THHH|H50)P(H50)}{P(THHH)}$

$$= \frac{(0.5)^4 \cdot 0.5}{0.5(0.5)^4 + 0.5(0.4 \cdot 0.6^3)}$$

$$\approx \boxed{0.4197}$$

$$\text{iii } P(H50 | 9H, IT) = ?$$

$$= \frac{P(9H, IT | H50) \cdot P(H50)}{P(9H, IT)}$$

$$P(9H, IT) = P(9H, IT | H50) P(H50) + P(9H, IT | H55) P(H55)$$

$$+ P(9H, IT | H60) P(H60)$$

$$= 0.5^{10} \cdot \frac{1}{3} + 0.55^9 \cdot 0.45 \cdot \frac{1}{2} + 0.6^9 \cdot 0.4 \cdot \frac{1}{3}$$

$$\approx 0.00236$$

$$P(H50 | 9H, IT) = \frac{0.5^{10} \cdot \frac{1}{3}}{P(9H, IT)} \approx [0.1379]$$

$$P(H55 | 9H, IT) = \frac{0.55^9 \cdot 0.45 \cdot \frac{1}{3}}{P(9H, IT)} \approx [0.2927]$$

$$P(H60 | 9H, IT) = \frac{0.6^9 \cdot 0.4 \cdot \frac{1}{3}}{P(9H, IT)} = [0.5694]$$

- (b) (5 points) Students at UCLA are from these disciplines: 15% Science, 21% Healthcare, 24% Liberal Arts, and 40% Engineering. (Each student belongs to a unique discipline.) The students attend a lecture and give feedback. Suppose 90% of the Science students liked the lecture, 18% of the Healthcare students liked it, none of the Liberal Arts students liked it, and 10% of the Engineering students liked it. If a student is randomly chosen, and the student liked the lecture, what is the conditional probability that the student is from Science?

Event  $L$  = liked the lecture.

$$P(\text{Science} \mid L) = \frac{P(L \mid \text{Science}) P(\text{Science})}{P(L)} = ?$$

$$\begin{aligned} P(\text{Science}) &= 0.15 \\ P(L \mid \text{Science}) &= 0.9 \end{aligned}$$

$$= \frac{0.9 \cdot 0.15}{0.2128} = 0.6344$$

$$\begin{aligned} P(L) &= P(L \mid S) P(S) + P(L \mid H) P(H) + P(L \mid LA) P(LA) + P(L \mid E) \\ &= 0.9 \cdot 0.15 + 0.18 \cdot 0.21 + 0 + 0.1 \cdot 0.4 \\ &= 0.2128 \end{aligned}$$

(c) (5 points) Consider a pregnancy test with the following statistics.

- If the woman is pregnant, the test returns “positive” (or 1, indicating the woman is pregnant) 99% of the time.
- If the woman is not pregnant, the test returns “positive” 10% of the time.
- At any given point in time, 99% of the female population is not pregnant.

What is the probability that a woman is pregnant given she received a positive test? The answer should make intuitive sense; give an explanation of the result that you find.

“1” is positive, “0” is negative test

P is pregnant

$$P(P|1) = \frac{P(1|P)P(P)}{P(1)} = \frac{0.99 \cdot 0.01}{P(1)} = 0.09$$

$$\begin{aligned} P(1) &= P(1| \text{Not } P)P(\text{Not } P) + P(1|P)P(P) \\ &= 0.1 \cdot 0.99 + 0.99 \cdot 0.01 \\ &= 0.1089 \end{aligned}$$

The test sucks (only 9% correct given positive test)  
There are way many more nonpregnant women than pregnant women (99% not pregnant), and it still fails 10% of the time on nonpregnant women. It fails more on nonpregnant women (10% of 99%) than it gets correct given it is positive (9%).

- (d) (5 points) Let  $x_1, x_2, \dots, x_n$  be identically distributed random variables. A random vector,  $\mathbf{x}$ , is defined as

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

What is  $\mathbb{E}(\mathbf{Ax} + \mathbf{b})$  in terms of  $\mathbb{E}(\mathbf{x})$ , given that  $\mathbf{A}$  and  $\mathbf{b}$  are deterministic?

$$\begin{aligned} \mathbf{x} &= \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} & \mathbb{E}(\mathbf{Ax} + \mathbf{b}) &=? \\ & & \mathbb{E}(\mathbf{Ax}_i) &= \mathbb{E}\left(\sum_{j=1}^n A_{ij} x_j\right) \\ & & &= \sum_{j=1}^n A_{ij} \mathbb{E}(x)_j \\ & & &= [\mathbf{A} \cdot \mathbb{E}(\mathbf{x})]_i \end{aligned}$$

Expectation is linear, so  $\mathbb{E}(\mathbf{Ax}) = \mathbf{A} \mathbb{E}(\mathbf{x})$

$$\boxed{\mathbb{E}(\mathbf{Ax} + \mathbf{b}) = \mathbf{A} \mathbb{E}(\mathbf{x}) + \mathbf{b}}$$

- (e) (5 points) Let

$$\text{cov}(\mathbf{x}) = \mathbb{E}((\mathbf{x} - \mathbb{E}\mathbf{x})(\mathbf{x} - \mathbb{E}\mathbf{x})^T)$$

What is  $\text{cov}(\mathbf{Ax} + \mathbf{b})$  in terms of  $\text{cov}(\mathbf{x})$ , given that  $\mathbf{A}$  and  $\mathbf{b}$  are deterministic?

$$\begin{aligned} \text{cov}(\mathbf{x}) &= \mathbb{E}((\mathbf{x} - \mathbb{E}\mathbf{x})(\mathbf{x} - \mathbb{E}\mathbf{x})^T) \\ \text{cov}(\mathbf{Ax} + \mathbf{b}) &=? \\ &= \mathbb{E}((\mathbf{Ax} + \mathbf{b} - \mathbb{E}(\mathbf{Ax} + \mathbf{b}))(\mathbf{Ax} + \mathbf{b} - \mathbb{E}(\mathbf{Ax} + \mathbf{b}))^T) \\ &= \mathbb{E}((\mathbf{Ax} + \mathbf{b} - \mathbf{A}\mathbb{E}(\mathbf{x}) - \mathbf{b})(\mathbf{Ax} + \mathbf{b} - \mathbf{A}\mathbb{E}(\mathbf{x}) - \mathbf{b})^T) \\ &= \mathbb{E}(\mathbf{A}(\mathbf{x} - \mathbb{E}(\mathbf{x}))(\mathbf{A}(\mathbf{x} - \mathbb{E}(\mathbf{x})))^T) \\ &= \mathbb{E}(\mathbf{A}(\mathbf{x} - \mathbb{E}(\mathbf{x}))((\mathbf{x} - \mathbb{E}(\mathbf{x}))^T \mathbf{A}^T)) \\ &= \mathbf{A} \mathbb{E}((\mathbf{x} - \mathbb{E}(\mathbf{x}))(\mathbf{x} - \mathbb{E}(\mathbf{x}))^T) \mathbf{A}^T \\ &= \boxed{\mathbf{A} \text{cov}(\mathbf{x}) \mathbf{A}^T} \end{aligned}$$

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  $\nabla_{\mathbf{x}} f$ ?

$$a) \nabla_{\mathbf{x}} \underset{n \times 1}{\mathbf{x}}^T \mathbf{A} \mathbf{y} = \partial \underset{\mathbf{x}}{\mathbf{x}^T \mathbf{A} \mathbf{y}} / \partial \mathbf{x} = [\mathbf{A} \mathbf{y}]$$

$$b) \nabla_{\mathbf{y}} \mathbf{x}^T \mathbf{A} \mathbf{y} = \partial \underset{\mathbf{y}}{\mathbf{x}^T \mathbf{A} \mathbf{y}} / \partial \mathbf{y} = (\mathbf{x}^T \mathbf{A})^T = \underset{m \times 1}{[\mathbf{A}^T \mathbf{x}]}$$

$$c) \nabla_{\mathbf{A}} \mathbf{x}^T \mathbf{A} \mathbf{y} = ?$$

$$\mathbf{x}^T \mathbf{A} \mathbf{y} = \sum_{i=1}^n \sum_{j=1}^m x_i y_j A_{ij} \quad \nabla_{\mathbf{A}} \underset{n \times 1 \times m}{\sum_{i=1}^n \sum_{j=1}^m x_i y_j A_{ij}} = \underset{n \times m}{[\mathbf{x} \mathbf{y}^T]}$$

$$d) f = \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} \quad \text{From class, } \nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$$

$$\nabla_{\mathbf{x}} f = ? = \boxed{(\mathbf{A} + \mathbf{A}^T) \mathbf{x} + \mathbf{b}}$$

(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  $\nabla_{\mathbf{A}} f$ ?

e)  $\mathbf{A} \in \mathbb{R}^{n \times n}, \mathbf{B} \in \mathbb{R}^{n \times n}, f = \text{tr}(\mathbf{AB})$

$\nabla_{\mathbf{A}} f = ?$

$$\text{tr}(\mathbf{AB}) = \sum_{i=1}^N (\mathbf{AB})_{ii} = \sum_{i=1}^N a_{1i} b_{1i} + \sum_{i=1}^N a_{2i} b_{12} + \dots$$

$$\frac{\partial \text{tr}(\mathbf{AB})}{\partial a_{ij}} = b_{ji}, \text{ so } \boxed{\nabla_{\mathbf{A}} f = \mathbf{B}^T} \quad \begin{matrix} \leftarrow \text{Matrix} \\ \text{cookbook} \\ \text{eq. 100} \end{matrix}$$

f)  $f = \text{tr}(\mathbf{BA} + \mathbf{A}^T \mathbf{B} + \mathbf{A}^2 \mathbf{B}), \nabla_{\mathbf{A}} f = ?$

$$\text{tr}(\mathbf{BA}) = \sum_{i=1}^N (\mathbf{BA})_{ii} = \sum_{i=1}^N b_{1i} a_{1i} + \sum_{i=1}^N b_{2i} a_{12} + \dots$$

$$\frac{\partial \text{tr}(\mathbf{BA})}{\partial a_{ij}} = b_{ji}, \text{ so } \nabla_{\mathbf{A}} \text{tr}(\mathbf{BA}) = \mathbf{B}^T \quad \checkmark, \text{ since } \text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$$

$$\text{tr}(\mathbf{A}^T \mathbf{B}) = \sum_{i=1}^N \sum_{j=1}^M a_{ij} b_{ij}$$

$$\frac{\partial \text{tr}(\mathbf{A}^T \mathbf{B})}{\partial a_{ij}} = b_{ij} = \mathbf{B} \quad \begin{matrix} \text{Matrix cookbook eq 103} \\ \text{eq. 103} \end{matrix}$$

$$\frac{\partial \text{tr}(\mathbf{A}^2 \mathbf{B})}{\partial \mathbf{A}} = (\mathbf{AB} + \mathbf{BA})^T \quad \begin{matrix} \text{matrix cookbook} \\ \text{equation 107} \end{matrix}$$

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

$$\boxed{\nabla_{\mathbf{A}} f = \mathbf{B}^T + \mathbf{B} + (\mathbf{AB} + \mathbf{BA})^T}$$

(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  $\nabla_{\mathbf{A}} f$ ?

$$f = \|(\mathbf{A} + \lambda\mathbf{B})\|_F^2 \quad \nabla_{\mathbf{A}} f = ?$$

$$f = \text{tr}[(\mathbf{A} + \lambda\mathbf{B})^T(\mathbf{A} + \lambda\mathbf{B})]$$

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

Drop terms with no dependence on  $\mathbf{A}$

$$f = \text{tr}(\mathbf{A}^T\mathbf{A}) + \lambda \text{tr}(\mathbf{A}^T\mathbf{B}) + \lambda \text{tr}(\mathbf{B}^T\mathbf{A})$$

$$\text{tr}(\mathbf{A}^T\mathbf{B}) = \text{tr}(\mathbf{B}^T\mathbf{A})$$

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

$$\nabla_{\mathbf{A}} f = \frac{\partial \text{tr}(\mathbf{A}^T\mathbf{A})}{\partial \mathbf{A}} + 2\lambda \frac{\partial \text{tr}(\mathbf{A}^T\mathbf{B})}{\partial \mathbf{A}}$$

matrix cookbook

eq 103

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

↑

matrix cookbook

eq 115

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:

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

$$4. \min_w \frac{1}{2} \sum_{i=1}^n \|y^{(i)} - w x^{(i)}\|^2 \quad \frac{\partial \text{tr}(WA)}{\partial w} = A^T$$

$$W = ?$$

$$\frac{\partial \text{tr}(WAw^T)}{\partial w} = WA^T + WA$$

$$\|A\|^2 = \text{tr}(A^T A)$$

$$f(w) = \frac{1}{2} \sum_{i=1}^n \text{tr}((y_i - w x_i)^T (y_i - w x_i)) \text{ Ignore terms without } w$$

$$= \frac{1}{2} \text{tr}((Y - W X)^T (Y - W X)) \quad \downarrow$$

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

$$= -\text{tr}(Y^T W X) + \frac{1}{2} \text{tr}(X^T W^T W X) \quad \text{tr}(AB) = \text{tr}(BA)$$

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

$$F(w) = -\text{tr}(W X Y^T) + \frac{1}{2} \text{tr}(W X X^T W^T)$$

$$\frac{\partial F(w)}{\partial w} = -Y X^T + \frac{1}{2} (W X X^T + W X X^T) = 0$$

$$= -Y X^T + W X X^T = 0$$

minimize

$$W X X^T = Y X^T$$

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

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  $\theta^*$ .

$$5. \underset{\theta}{\operatorname{argmin}} \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2 + \frac{\lambda}{2} \|\theta\|_2^2$$

$$\theta^* = ?$$

From lecture/discussion

$$\frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2 = \frac{1}{2} (Y - X\theta)^T (Y - X\theta)$$

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

$$X \in \mathbb{R}^{n \times 2}$$

$$\underset{\theta}{\operatorname{argmin}} \frac{1}{2} (Y - X\theta)^T (Y - X\theta) + \frac{\lambda}{2} \|\theta\|_2^2$$

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

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

$$= \frac{1}{2} Y^T Y - Y^T X\theta + \frac{1}{2} [\theta^T (X^T X + \lambda I) \theta]$$

$$\frac{\partial L(\theta)}{\partial \theta} = -X^T Y + (X^T X + \lambda I) \theta = 0$$

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

# 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 2023, Prof. J.C. Kao, TAs: T.M, P.L, R.G, K.K, N.V, S.R, S.P, M.E

```
In [1]: import numpy as np
import matplotlib.pyplot as plt

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

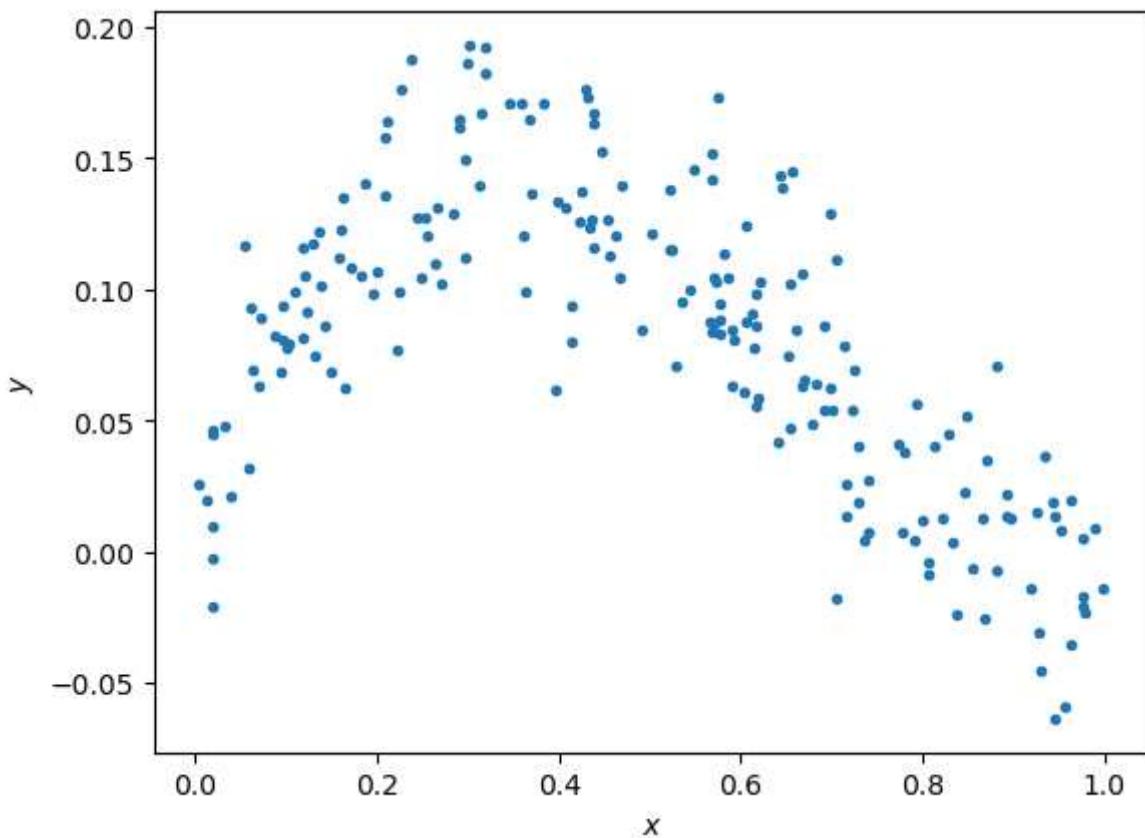
## Data generation

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

```
In [2]: 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('$y$')
```

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



## 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) The generating distribution of  $x$  is a uniform distribution from 0 to 1.
- (2) The distribution of the additive noise  $\epsilon$  is normal (Gaussian) distribution; the mean is 0 and standard deviation is 0.03.

## 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]: # 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)

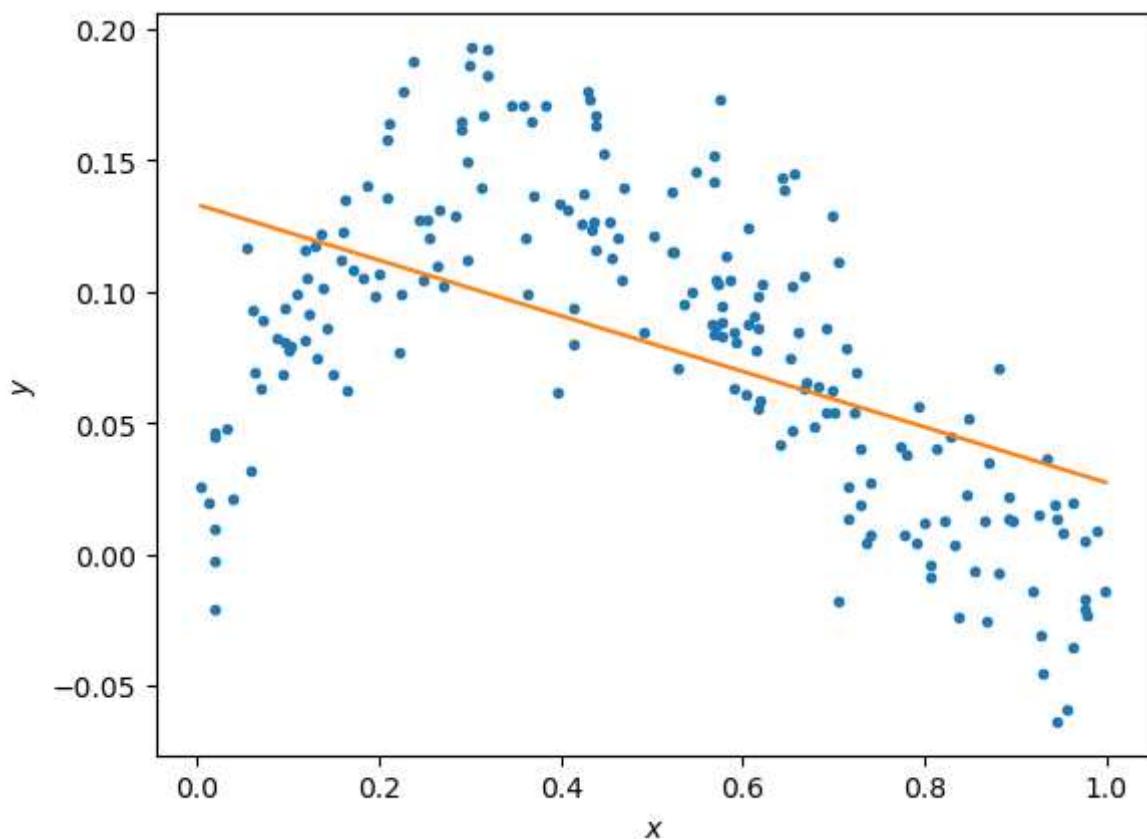
# ===== #
# END YOUR CODE HERE #
# ===== #
```

In [4]: *# 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))
```

Out[4]: [`<matplotlib.lines.Line2D at 0x240a907b160>`]



## QUESTIONS

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

## ANSWERS

- (1) This model underfits the data.

(2) We need to use a model of higher order; this means adding more parameters to theta to increase the complexity of the model.

## 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]: 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
# 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,
# ... etc.

xhats = [xhat]
thetas = [theta]

for i in range(1, N):
    xhat = np.vstack((x***(i+1), xhat))
    theta = np.linalg.inv(xhat @ xhat.T) @ (xhat @ y)
    xhats.append(xhat)
    thetas.append(theta)

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

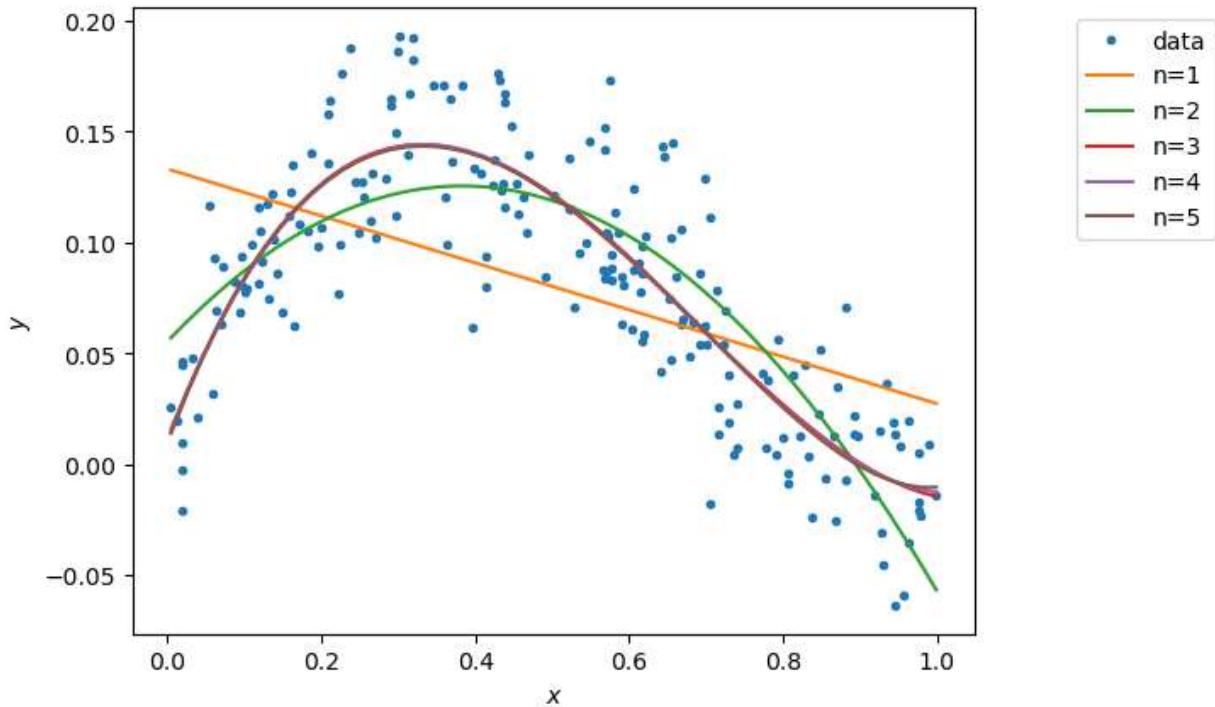
```
In [6]: # Plot the data
f = plt.figure()
ax = f.gca()
ax.plot(x, y, '.')
ax.set_xlabel('$x$')
ax.set_ylabel('$y$')

# 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)
```



## Calculating the training error (5 points)

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

```
In [7]: 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(N):
    y_pred = y - thetas[i] @ xhats[i]
    training_errors.append(y_pred.T @ y_pred / num_train)

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

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

Training errors are:  
[0.0023799610883627007, 0.0010924922209268528, 0.0008169603801105372, 0.000816535373529698, 0.0008161479195525292]

## QUESTIONS

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

## ANSWERS

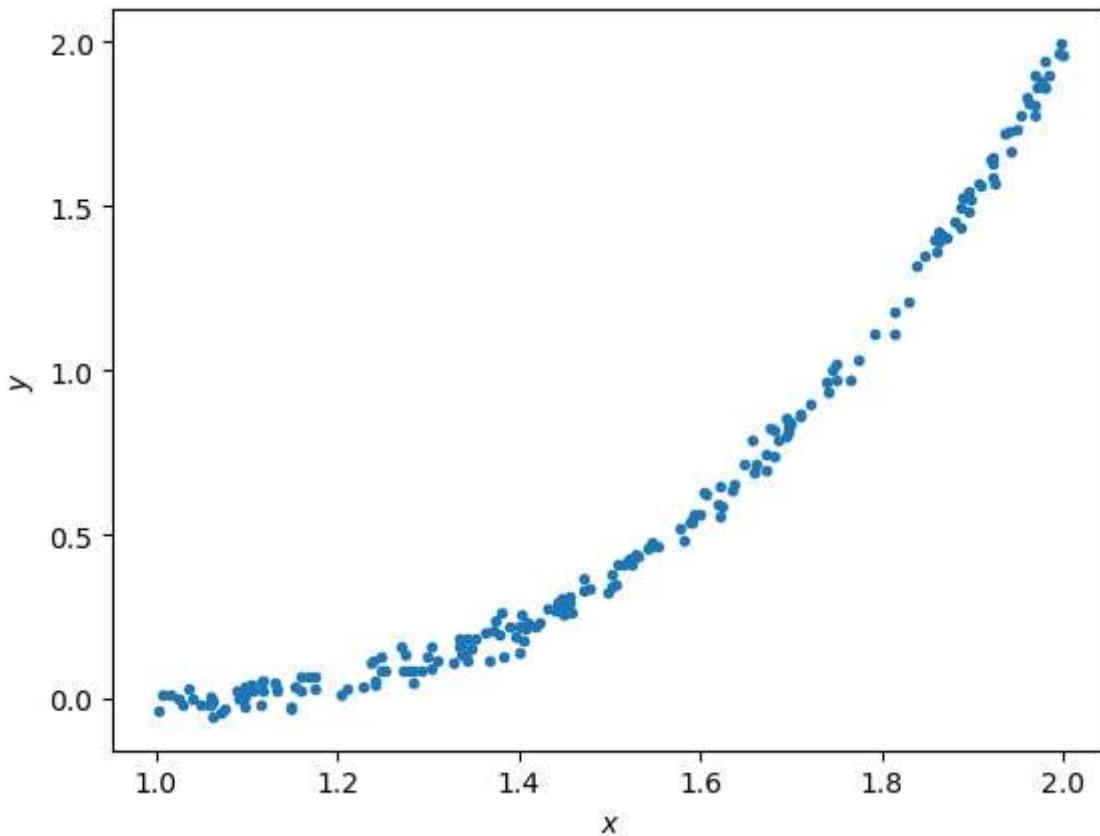
- (1) The polynomial of degree 5 has the lowest training error (best).
- (2) This is expected because higher-order polynomials have more parameters and dimensions to fit the data best (barring any overfitting). Higher order models perform at least as well as lower order models, so we expect there to be lower training error with more complex models.

## 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 [8]: 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('$y$')

Out[8]: Text(0, 0.5, '$y$')
```



```
In [9]: 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)
```

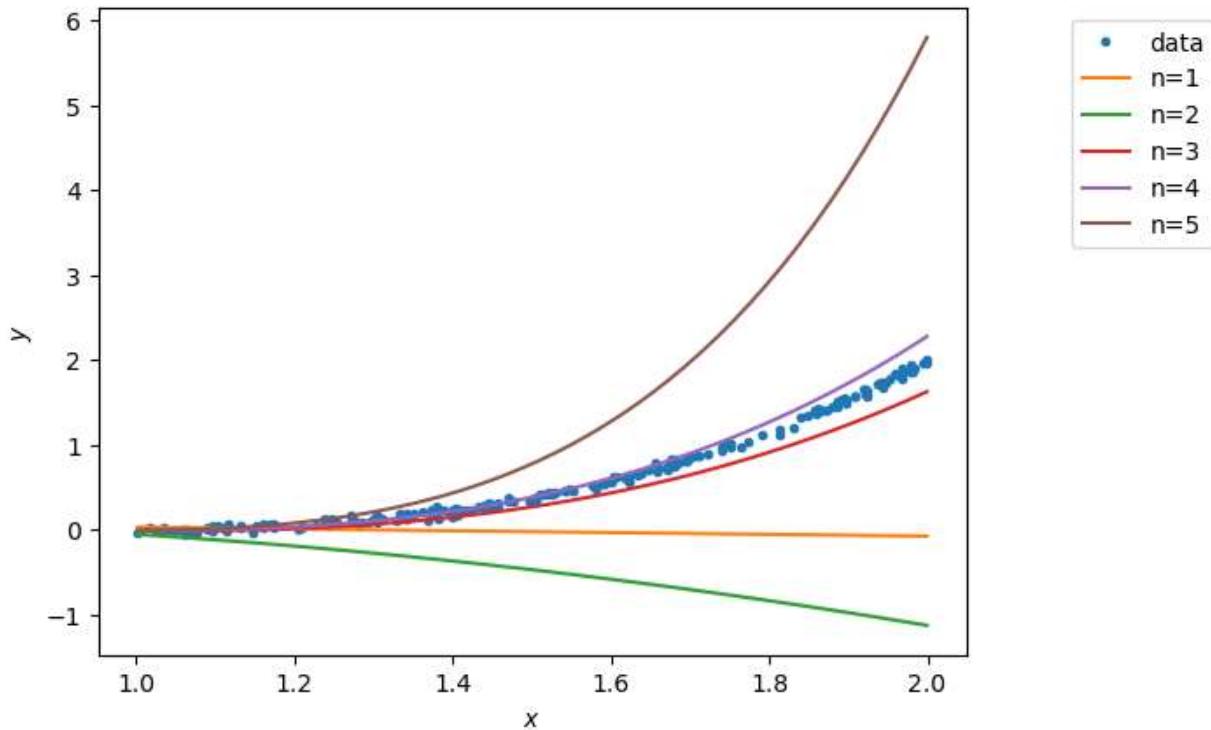
```
In [10]: # Plot the data
f = plt.figure()
ax = f.gca()
ax.plot(x, y, '.')
ax.set_xlabel('$x$')
ax.set_ylabel('$y$')

# 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)
```



```
In [11]: 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(N):
    y_pred = y - thetas[i] @ xhats[i]
    testing_errors.append(y_pred.T @ y_pred / num_train)

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

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

Testing errors are:  
[0.8086165184550589, 2.131919244505802, 0.03125697108274475, 0.011870765189429009, 2.149102183459647]

## QUESTIONS

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

## ANSWERS

- (1) The fourth order polynomial (degree 4) has the lowest testing error (best).
- (2) The fifth order polynomial does not generalize well because it overfits the data. It fits very well to the training data, but given new unseen data (the test data), it does not generalize well because it learns too much from the noise of the training data.