

Margaret Capotz  
 (47 HW1)

i) Linear alg

a) Let  $Q$  be a real orthogonal matrix.

i) Show  $Q^T$  and  $Q^{-1}$  are also orthogonal.

To test if a matrix is orthogonal, see if it satisfies  $AA^T = I$   
 Definition of orthogonal matrix  $A$ :  $A^T = A^{-1}$ .

Pf.

$(Q^{-1})^{-1} = Q$ ,  $(Q^T)^T = Q$  by definition of inverse, transpose.

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

$$\therefore (Q^T)^T = (Q^T)^{-1} \Rightarrow Q^T \text{ is orthogonal.}$$

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

$$\therefore (Q^{-1})^{-1} = (Q^{-1})^T \Rightarrow Q^{-1} \text{ is orthogonal. } \square.$$

ii) Show that  $Q$  has eigenvalues with norm 1.

Pf.

By definition,  $Q$  is comprised of orthonormal vectors as columns & rows.  
 Thus the norms of each orthonormal vector, by definition, is 1.

Consider  $Qx = \lambda x$ , where  $\lambda$  is an eigenvalue and  $x$  is the corresponding eigenvector.

$$\begin{aligned} \|Qx\|^2 &= \|\lambda x\|^2 \\ (Qx)^T(Qx) &= |\lambda|^2 \|x\|^2 \text{ by def of length} \\ x^T Q^T Q x &= |\lambda|^2 \|x\|^2 \text{ by transpose distribution} \\ x^T x &= |\lambda|^2 \|x\|^2 \text{ by def of orthonormal} \\ \|x\|^2 &= |\lambda|^2 \|x\|^2 \text{ by def of length} \\ |\lambda| &= |\lambda|^2 \\ |\lambda| &= \lambda \end{aligned}$$

$\therefore$  all eigenvalues of  $Q$  have norm 1.  $\square$ .

iii) Show the determinant of  $Q$  is either +1 or -1.

Pf.

Take  $QQ^T = I$  by definition of orthogonal.

$\det(QQ^T) = \det I$  by taking the det of both sides

$$\det(Q)\det(Q^T) = 1$$

$$(\det Q)^2 = 1 \quad \text{since } \det(Q) = \det(Q^T)$$

$$\det Q = \pm 1 \quad \text{square rooting both sides}$$

$$\therefore \det Q = \pm 1. \quad \square.$$

length preserving transformation:  $\|Qv\|_2^2 = \|v\|_2^2$  length preserved

- iv) show that  $\alpha$  defines a length preserving transformation.

Pf.

$$\begin{aligned} & \|Qv\|_2^2 \\ &= (Qv)^T(Qv) \\ &= v^T Q^T Q v \\ &= v^T v \\ &= \|v\|_2^2 \end{aligned}$$

$\therefore \alpha$  transformation preserves length.  $\square$ .

CAN WORK ON  
PIZZA

- b) let  $A$  be a matrix.

- i) what is the relationship between the singular vectors of  $A$  and the eigenvectors of  $AA^T$ ? What about  $A^TA$ ?

$$\begin{aligned} A &= U \Sigma V^T \\ &\text{deft right} \\ A A^T &= U \Lambda U^T \\ A^T A &= V \Lambda V^T \\ &\uparrow \quad \nwarrow \text{orthogonal} \\ &\text{eigenvalues} \end{aligned}$$

$$\begin{aligned} AA^T &= U \Sigma V^T (U \Sigma V^T)^T \\ &= U \Sigma V^T V^T \Sigma^T U^T \\ &= U \Sigma V^T V \Sigma^T U^T \\ &= U \Sigma \Sigma^T U^T \\ &\sim \\ &= \Sigma \Sigma^T = \sigma^2 \\ &= U \Sigma^2 U^T \end{aligned}$$

for  $\begin{bmatrix} a & b \\ c & d \end{bmatrix}$ ,  $A = U \Sigma V^T$

singular vectors  
singular values  
eigendecomp:  $A$  square, diag  
SVD basics:  
 $A$  not necessarily square  
 $AA^T$  always  
symmetric

$U$  is the eigenvector of  $AA^T$ .  
 $V$  is the eigenvector of  $A^TA$ .

- ii) what about values of  $A$  & eigenvalues of  $AA^T$ ,  $A^TA$

singular values ( $\sigma$ ) are the square of the eigenvalues ( $\lambda$ ), for both  $AA^T$  and  $A^TA$ .

$$\sigma^2 = \lambda$$

see part (i) for work.

- c) True or false.

- i) every linear operator in an  $n$ -dimensional vector space has  $n$  distinct eigenvalues.  
False, should be at most  $n$  distinct eigenvalues.
- ii) A non-zero sum of two eigenvectors of a matrix  $A$  is an eigenvector.  
False, only if they share an eigenvalue.

- iii) If a matrix  $A$  has the positive semidefinite property, i.e.  $x^T A x \geq 0 \forall x$ , then its eigenvalues must be non-negative.  
 True. Similar to discussion question.
- iv) The rank of a matrix can exceed the number of distinct non-zero eigenvalues.  
 True, if there is repetition of distinct eigenvalues. 1 eigenvalue for  $I$  matrix.
- v) A non-zero sum of eigenvectors of a matrix  $A$  corresponding to the same eigenvalue  $\lambda$  is always an eigenvector.  
 True.

## 2) Probability

a) Jar of coins. 'H50' coins :  $P(H) = 0.5$   
 'H60' coins :  $P(H) = 0.4$

i) Flip a coin, lands tails. Posterior probability it is an H50 coin?

$$\begin{aligned} P(H50) &= 0.5 \\ P(H60) &= 0.5 \\ P(H | H50) &= 0.5 \\ P(H | H60) &= 0.4 \end{aligned}$$

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

$$P(H50 | T) = \frac{P(H50) P(T | H50)}{P(T)} = \frac{0.5 \cdot 0.5}{0.45} = \boxed{0.5556}$$

ii) Flip a coin 4 times, landing THTH.  
 How likely is the coin to be type H50?

$$\begin{aligned} P(THTH) &= P(THTH \cap H50) + P(THTH \cap H60) \\ &= P(H50) P(THTH | H50) + P(H60) P(THTH | H60) \\ &= 0.5 \cdot (0.5)^4 + 0.5 \cdot (0.4 \cdot 0.6^3) \\ &= 0.07445 \end{aligned}$$

$$\begin{aligned} P(H50 | P(THTH)) &= \frac{P(THTH | H50) P(H50)}{P(THTH)} = \frac{(0.5)^5}{0.07445} = \boxed{0.4197} \end{aligned}$$

iii) New jar, with H50, H55, H60.

Take 1 coin, flip it 10 times. It lands heads 9 times.

How likely is the coin of each type?

Let's call the event in question to be A.

Note:

$\binom{10}{9} = 10$  ways to have 9 heads among 10 flips.

$$P(A) = P(A \cap H50) + P(A \cap H55) + P(A \cap H60)$$

$$= P(H50)P(A|H50) + P(H55)P(A|H55) + P(H60)P(A|H60)$$

$$= \frac{1}{3} \cdot 10 \cdot (\frac{1}{2})^{10} + \frac{1}{3} \cdot 10 \cdot \frac{45}{100} \cdot \left(\frac{55}{100}\right)^9 + \frac{1}{3} \cdot 10 \cdot \frac{40}{100} \cdot \left(\frac{60}{100}\right)^9$$

$$= 0.0236$$

$$P(H50|A) = \frac{P(A|H50)P(H50)}{P(A)} = \frac{\frac{1}{3} \cdot 10 \cdot (\frac{1}{2})^{10}}{0.0236} = [0.1379]$$

$$P(H55|A) = \frac{P(A|H55)P(H55)}{P(A)} = \frac{\frac{1}{3} \cdot 10 \cdot \frac{45}{100} \cdot \left(\frac{55}{100}\right)^9}{0.0236} = [0.2927]$$

$$P(H60|A) = \frac{P(A|H60)P(H60)}{P(A)} = \frac{\frac{1}{3} \cdot 10 \cdot \frac{40}{100} \cdot \left(\frac{60}{100}\right)^9}{0.0236} = [0.5694]$$

b) UCLA students:

15% science	90% science like lecture
21% healthcare	18% healthcare "
24% liberal arts	0% liberal arts "
40% engineering	10% engr "

Suppose a student liked elec. What prob they are sci?

Want:  $P(\text{science} | \text{liked lecture})$ , or  $P(S|L)$

Define:  
 S as science  
 H as healthcare  
 A as liberal arts  
 E as engr  
 L as liked lecture

$$\begin{aligned} P(L) &= P(L \cap S) + P(L \cap H) + P(L \cap A) + P(L \cap E) \\ &= P(S)P(L|S) + P(H)P(L|H) + P(A)P(L|A) + P(E)P(L|E) \\ &= 0.15 \cdot 0.9 + 0.21 \cdot 0.18 + 0.24 \cdot 0 + 0.4 \cdot 0.1 \\ &= 0.2128 \end{aligned}$$

$$P(S|L) = \frac{P(L|S)P(S)}{P(L)} = \frac{0.9 \cdot 0.15}{0.2128} = [0.634]$$

c) Pregnancy test: what is the probability a woman is pregnant given a pos. test?  
 Give an explanation.

Define +: positive  
 P: pregnant

$$\begin{aligned} P(+|P) &= 0.99 & P(+) &= P(+ \cap P) + P(+ \cap P') \\ P(+|P') &= 0.1 & &= P(P)P(+|P) + P(P')P(+|P') \\ P(P) &= 0.01 & &= 0.01 \cdot 0.99 + 0.99 \cdot 0.1 \\ P(P') &= 0.99 & &= 0.1089 \end{aligned}$$

$$P(P|+) = \frac{P(+|P)P(P)}{P(+)} = \frac{0.99 \cdot 0.01}{0.1089} = \boxed{0.09}$$

There is a 9% probability a woman is pregnant given a positive pregnancy test. This is intuitive since the rate of false positives is high at 10%.

The majority of women who take the test receive a false positive, because very little women are actually pregnant.

- d) Let  $x_1, x_2, \dots, x_n$  be identically distributed random vars.

A random vector,  $x$ , is defined as

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

What is  $E(Ax + b)$  in terms of  $E(x)$ , given that  $A, b$  are deterministic?

$E$  is expectation  $E(x)$  is the mean of  $x$

$$y = Wx + b$$

at every single point, you know the exact value

$$E(b) = b, E(A) = A, E$$
 is linear

$$\begin{aligned} E(Ax + b) &= A E(x) + b \\ &= E(A) E(x) + E(b) \end{aligned}$$

$$\begin{bmatrix} a_{11} & & \\ & \ddots & \\ & & a_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} =$$

show it is linearly separable

Pf.

Let  $A \in \mathbb{R}^{m \times n}, x \in \mathbb{R}^n$

$$\text{then } Ax = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n \\ \vdots \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n \end{bmatrix}$$

$$\begin{aligned} \text{then } Ax + b &= \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n \\ \vdots \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n \end{bmatrix} + \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix} \\ &= \begin{bmatrix} a_{11} \\ \vdots \\ a_{m1} \end{bmatrix} x_1 + \begin{bmatrix} a_{12} \\ \vdots \\ a_{m2} \end{bmatrix} x_2 + \dots + \begin{bmatrix} a_{1n} \\ \vdots \\ a_{mn} \end{bmatrix} x_n + \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix} \end{aligned}$$

Thus  $Ax + b$  is linearly separable.

$$\begin{aligned} \therefore E(Ax + b) &= E(A) E(x) + E(b) \\ &= A E(x) + b \quad \text{as } A, b \text{ are deterministic.} \end{aligned}$$

$$\text{Thus } E(Ax + b) = A E(x) + b.$$

□.

- e) Let  $\text{cov}(x) = E((x - E(x))(x - E(x))^T)$

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

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

$$\begin{aligned}
&= \mathbb{E}((Ax + b - A\mathbb{E}(x) - b)(Ax + b - A\mathbb{E}(x) - b)^T) \\
&= \mathbb{E}((A(x - \mathbb{E}(x)))(A(x - \mathbb{E}(x))^T)) \\
&= A \cdot \mathbb{E}((x - \mathbb{E}(x))(x - \mathbb{E}(x))^T) \cdot A^T \\
&= A \operatorname{cov}(X) A^T
\end{aligned}$$

Ref me #

3) Multivariate derivatives.

a) Let  $x \in \mathbb{R}^n$ ,  $y \in \mathbb{R}^m$ ,  $A \in \mathbb{R}^{n \times m}$ . What is  $\nabla_x x^T A y$ ?

$$\nabla_x x^T A y = \nabla_x x^T (Ay) = Ay \text{ by the matrix cookbook}$$

$$(\nabla_x x^T A = A)$$

Alternatively,

$$\begin{aligned}
\frac{\partial f(x,y)}{\partial x_i} &= a_{1i} y_1 + \sum_{j=2}^m a_{ij} y_j + 0 \\
&= \sum_{j=1}^m a_{ij} y_j \\
&= \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \\
&= Ay
\end{aligned}$$

same!

$$f(x,y) = x^T A y = \sum_{i=1}^n \sum_{j=1}^m a_{ij} \cdot x_i \cdot y_j$$

$$i=1, j=1$$

$$i=1, j \neq 1$$

$$a) \nabla_x x^T A y = \nabla_x x^T (Ay) = Ay$$

$$b) \nabla_y x^T A y = \nabla_y (x^T A) y = (x^T A)^T = A^T x$$

$$i \neq 1, j=1$$

$$i \neq 1, j \neq 1$$

$$\frac{\partial a_{11} x_1 y_1}{\partial x_i} = a_{11} y_1, \quad \frac{\partial \left( \sum_{j=2}^m a_{ij} \cdot x_i \cdot y_j \right)}{\partial x_i} = \sum_{j=2}^m a_{ij} \cdot y_j \quad \frac{\partial \sum_{i=2}^n (a_{ii} \cdot x_i \cdot y_i)}{\partial x_i} = 0$$

b)

Let  $x \in \mathbb{R}^n$ ,  $y \in \mathbb{R}^m$ ,  $A \in \mathbb{R}^{n \times m}$ . What is  $\nabla_y x^T A y$ ?

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

Alternatively,

$$(\nabla_x A x = A)$$

$$i \neq 1, j=1$$

$$\begin{aligned}
\frac{\partial f(x,y)}{\partial y_1} &= a_{11} x_1 + \sum_{i=2}^n a_{i1} \cdot x_i + 0 \\
&= \sum_{i=1}^n a_{i1} \cdot x_i \\
&= \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix}^T \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \\
&= A^T x
\end{aligned}$$

same!

$$i \neq 1, j \neq 1$$

$$f(x,y) = x^T A y = \sum_{i=1}^n \sum_{j=1}^m a_{ij} \cdot x_i \cdot y_j$$

$$i=1, j=1$$

$$i=1, j \neq 1$$

$$\frac{\partial \sum_{i=2}^n (a_{ii} \cdot x_i \cdot y_i)}{\partial y_1} = \sum_{i=2}^n a_{ii} \cdot x_i$$

$$\frac{\partial a_{11} x_1 y_1}{\partial y_1} = a_{11} x_1$$

$$\frac{\partial \left( \sum_{j=2}^m a_{ij} \cdot x_i \cdot y_j \right)}{\partial y_1} = 0$$

c) What is  $\nabla_A x^T A y$ ?

By the matrix cookbook,

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

d) Let  $x \in \mathbb{R}^n$ ,  $A \in \mathbb{R}^{n \times n}$ ,  $f = x^T A x + b^T x$ . What is  $\nabla_x f$ ?

$$\begin{aligned}\nabla_x f &= \nabla_x x^T A x + \nabla_x b^T x, \text{ we have } \nabla_x x^T A x = Ax + A^T x \text{ and } \nabla_x b^T x = b, \text{ so} \\ &= Ax + A^T x + b \quad \text{using substitutions we showed in class.}\end{aligned}$$

e) Let  $A, B \in \mathbb{R}^{n \times n}$  and  $f = \text{tr}(AB)$ . What is  $\nabla_A f$ ?

$$\nabla_A \text{tr}(AB) = B^T \quad \text{by the matrix cookbook.}$$

f) Let  $A, B \in \mathbb{R}^{n \times n}$ ,  $f = \text{tr}(BA + A^T B + A^2 B)$ . What is  $\nabla_A f$ ?

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

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

$$= B^T + B + (AB + BA)^T \quad \text{by the matrix cookbook.}$$

g) Let  $A, B \in \mathbb{R}^{n \times n}$ ,  $f = \|A + \lambda B\|_F^2$ . What is  $\nabla_A f$ ?

matrix cookbook:

$$\begin{aligned}\nabla_A \|A + \lambda B\|_F^2 &= \nabla_A \text{Tr}((A + \lambda B)(A + \lambda B)^T) \\ &= \nabla_A \text{Tr}((A + \lambda B)(A^T + B^T \lambda^T)) \\ &= \nabla_A \text{Tr}(AA^T + AB^T \lambda^T + \lambda B A^T + \lambda B B^T \lambda^T) \\ &= \nabla_A \text{Tr}(AA^T) + \nabla_A \text{Tr}(AB^T \lambda^T) + \nabla_A \text{Tr}(\lambda B A^T) + \nabla_A (\lambda B B^T \lambda^T) \\ &= 2A + (B^T \lambda^T)^T + \lambda B + 0 \\ &= 2A + \lambda B + \lambda B + 0 = 2A + 2\lambda B\end{aligned}$$

$$\frac{\partial}{\partial X} \|X\|_F^2 = \frac{\partial}{\partial X} \text{Tr}(XX^*) = 2X$$

$$\|A\|_F^2 = \text{Tr}(AA^H) \quad \begin{array}{l} H: \text{conjugate} \\ \text{transpose} \end{array}$$

hints:  
F: Frobenius

For real matrices,  $A \in \mathbb{R}^{m \times n}$ , the conjugate transpose is just the transpose,  $A^H = A^T$

4) Deriving least-squares with matrix derivatives.

In least squares, we seek to estimate multivariate output  $y$  via model

$$\hat{y} = Wx$$

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

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

Derive the optimal  $W$ .

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

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

Take  $\frac{1}{2} \sum_{i=1}^n \|y^{(i)} - Wx^{(i)}\|^2$   
 "Vectorization"  $= \frac{1}{2} \sum_{i=1}^n (y^{(i)} - Wx^{(i)})^T (y^{(i)} - Wx^{(i)})$  by definition of 2-norm  
 $= \frac{1}{2} \left( \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix} - \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n)} \end{bmatrix} W^T \right)^T \left( \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix} - \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n)} \end{bmatrix} W^T \right)$  change formula to match dims  
 $\quad \quad \quad Y \in \mathbb{R}^{n \times m} \quad X \in \mathbb{R}^{n \times k} \quad W \in \mathbb{R}^{k \times m}$   
 $= \frac{1}{2} \|Y - XW^T\|_F^2$   
 $= \frac{1}{2} \text{Tr}((Y - XW^T)(Y - XW^T)^T)$   
 $= \frac{1}{2} \text{Tr}((Y - XW^T)(Y^T - W^T X))$   
 $= \frac{1}{2} \text{Tr}(YY^T - YW^T X^T - XW^T Y^T + XW^T W^T X^T)$   
 $= \frac{1}{2} \text{Tr}(YY^T - (XW^T Y^T)^T - XW^T Y^T + XW^T W^T X^T)$  same since scalar  $^T$  are equal  
 $= \frac{1}{2} \text{Tr}(YY^T - 2YW^T X^T + XW^T W^T X^T)$ , now take the derivative  
 $\frac{\partial}{\partial W} \left( \frac{1}{2} \text{Tr}(YY^T - 2YW^T X^T + XW^T W^T X^T) \right)$   
 $= \frac{1}{2} \left( \frac{\partial}{\partial W} \text{Tr}(YY^T) - 2 \frac{\partial}{\partial W} \text{Tr}(YW^T X^T) + \frac{\partial}{\partial W} \text{Tr}(XW^T W^T X^T) \right)$   
 $= \frac{1}{2} (0 - 2 \frac{\partial}{\partial W} \text{Tr}(W X^T Y) + \frac{\partial}{\partial W} \text{Tr}(W X^T W X^T))$   
 $= \frac{1}{2} (-2(X^T Y)^T + W(X^T X)^T + W(X^T X))$   
 $= -(X^T Y)^T + \frac{1}{2} W(X^T X + X^T X)$   
 $= -(X^T Y)^T + W(X^T X)$   
 Setting this equal to zero,  
 $0 = -(X^T Y)^T + W(X^T X)$   
 $(X^T Y)^T = W(X^T X)$   
 $W = (X^T X)^T (X^T X)^{-1}$

## 5) Regularized least-squares

In lecture, we worked through the regularized least squares problem

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

Regularization addresses overfitting, we'll consider ridge regularization aka regularized least-squares problem.  
solve the following optimization problem:

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

where  $\lambda$  is a tunable regularization param.

Observe we are seeking the least-squares sol w/ the smaller 2-norm.  
 Derive the solution, i.e. find  $\theta^*$

## Hints:

For  $\sum_{i=1}^N (y^{(i)} - \Theta^T \vec{x}^{(i)})^2 \rightarrow$  express as a vector

not  $y^{(c)}$  is a scalar  $\therefore Y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$  is a vector

for the second term  $\frac{1}{2} \|\Theta\|_F^2$ , try expressing it as the product of  $\Theta$  and  $\Theta^T$

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

by vectorization,

$$= \frac{1}{2} \left( \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix} - \theta^T \begin{bmatrix} \hat{x}^{(1)} \\ \hat{x}^{(2)} \\ \vdots \\ \hat{x}^{(N)} \end{bmatrix} \right)^2 + \frac{\lambda}{2} \|\theta\|_2^2$$

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

$\theta^T \mathbb{X} \in \mathbb{R}^{N \times n}$   
 $N \times N \quad N \times n$

the square of a matrix  
 can be expressed as its innerproduct

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

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

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

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

then, take the derivative w.r.t.  $\theta$

$\frac{\partial}{\partial \theta} \left( \frac{1}{2} (\mathbb{Y}^T \mathbb{Y} - \mathbb{Y}^T \theta^T \hat{\mathbb{X}} - \hat{\mathbb{X}}^T \theta \mathbb{Y} + \theta^T \hat{\mathbb{X}} \hat{\mathbb{X}}^T \theta) + \frac{\lambda}{2} \|\theta\|_2^2 \right)$

$$\begin{aligned}
&= \frac{1}{2} (0 - \cancel{YX^T} - \cancel{XY^T} + (\cancel{XX^T})(\cancel{XX^T}) + \lambda 2\theta) \\
&= \frac{1}{2} (-\cancel{YX^T} - (\cancel{YX^T})^T + (\cancel{XX^T})(\cancel{XX^T}) + \lambda 2\theta) \\
&= \frac{1}{2} (-\cancel{YX^T} - (\cancel{YX^T})^T + (\cancel{XX^T})^2 + 2\lambda\theta) \\
&= -\frac{1}{2} (\cancel{YX^T} + (\cancel{YX^T})^T) + \cancel{XX^T} + \lambda\theta
\end{aligned}$$

Try Again.

$$\begin{aligned}
&\frac{1}{2} \sum_{i=1}^N (y^{(i)} - \theta^T \hat{x}^{(i)})^2 + \frac{\lambda}{2} \|\theta\|_2^2 \\
&\quad \text{by vectorization} \\
&= \frac{1}{2} \left( \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} - \theta^T \begin{bmatrix} \hat{x}^{(1)} \\ \vdots \\ \hat{x}^{(n)} \end{bmatrix} \right)^2 + \frac{\lambda}{2} \|\theta\|_2^2 \\
&\quad Y \in \mathbb{R}^{1 \times n} \quad \theta^T \in \mathbb{R}^{1 \times m} \\
&\quad X \in \mathbb{R}^{m \times n} \\
&\quad \theta^T X \in \mathbb{R}^{1 \times n} \\
&= \frac{1}{2} (Y - X\theta)^T (Y - X\theta) + \frac{\lambda}{2} \|\theta\|_2^2 \\
&= \frac{1}{2} ((Y^T - \theta^T X^T)(Y - X\theta) + \lambda \|\theta\|_2^2) \\
&= \frac{1}{2} (Y^T Y - Y^T X\theta - \theta^T X^T Y + \theta^T X^T X\theta + \lambda \|\theta\|_2^2) \\
&\quad \nabla_X X^T A X = A X + A^T X \\
&\quad \nabla_X A X = A^T \\
&\quad \nabla_X X^T A = A \\
&\quad \text{Now, let's take the derivative w.r.t. } \theta \\
&\frac{\partial}{\partial \theta} \left( \frac{1}{2} (Y^T Y - \underbrace{Y^T X\theta}_{\text{scalars so can combine as}} - \underbrace{\theta^T X^T Y}_{\text{scalars so can combine as}} + \theta^T X^T X\theta + \lambda \|\theta\|_2^2) \right) \\
&= \frac{1}{2} \frac{\partial}{\partial \theta} (Y^T Y - 2 \underbrace{Y^T X\theta}_{\text{scalars so can combine as}} + \underbrace{\theta^T X^T X\theta}_{\text{scalars so can combine as}} + \lambda \|\theta\|_2^2) \\
&= \frac{1}{2} (0 - 2 X^T Y + (X^T X + (X^T X)^T)\theta + \lambda 2\theta) \\
&= -X^T Y + \frac{1}{2} (X^T X + X^T X)\theta + \lambda\theta \\
&= -X^T Y + X^T X\theta + \lambda I\theta
\end{aligned}$$

Now, setting this equal to zero, we can solve for  $\theta$

$$\theta = -X^T Y + (X^T X + \lambda I)\theta$$

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

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

$$\begin{aligned}
&= \frac{1}{2} (0 - (\cancel{YX^T})^T - (\cancel{X^T Y}) + (X^T X)(\cancel{X^T X}) + \lambda 2\theta) \\
&= \frac{1}{2} (-\cancel{X^T Y} - \cancel{X^T Y} + (X^T X)(\cancel{X^T X}) + 2\lambda\theta) \\
&= -X^T Y + \frac{1}{2} (X^T X)(X^T X) + \lambda\theta \\
&\quad \text{Solving for } \theta \text{ when setting this equal to zero,} \\
&\theta = \frac{1}{\lambda} \left( X^T Y - \frac{1}{2} (X^T X)(X^T X) \right)
\end{aligned}$$

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

$y^{(i)} \in \mathbb{R}^{M \times 1}$   
 $x^{(i)} \in \mathbb{R}^{n \times 1}$

$= \frac{1}{2} \sum_{i=1}^n (y^{(i)} - Wx^{(i)})^T (y^{(i)} - Wx^{(i)})$   
 ↑ flip transpose order for scalar

"vectorization" ↴

$= \frac{1}{2} \left( \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} - \begin{bmatrix} x^{(1)} \\ \vdots \\ x^{(n)} \end{bmatrix}^T W^T \right)^T \left( \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} - \begin{bmatrix} x^{(1)} \\ \vdots \\ x^{(n)} \end{bmatrix}^T W^T \right)$

$= \frac{1}{2} \|Y - X^T W^T\|_F^2$

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

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

$= \frac{1}{2} \text{Tr}(Y Y^T - \underbrace{Y^T W X}_{\text{same since scalar}^T \text{ are equal}} - X^T W^T Y + X^T W^T W X)$

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

Now take the derivative

$$\begin{aligned} & \frac{\partial}{\partial W} \left( \frac{1}{2} \text{Tr}(Y Y^T - 2 Y^T W X + X^T W^T W X) \right) \\ &= 0 - \frac{\partial}{\partial W}(Y^T W X) + \frac{1}{2} \frac{\partial}{\partial W}(X^T W^T W X) \\ &= -\frac{\partial}{\partial W}(W X Y^T) + \frac{1}{2} \frac{\partial}{\partial W}(W X X^T W) \\ &= -(X Y^T)^T + \frac{1}{2}(W(X X^T)^T + W(X X^T)) \\ &= -Y X^T + W(X X^T) \\ & \boxed{W = Y X^T (X X^T)^{-1}} \end{aligned}$$

# 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 [ ]: import numpy as np
import matplotlib.pyplot as plt

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

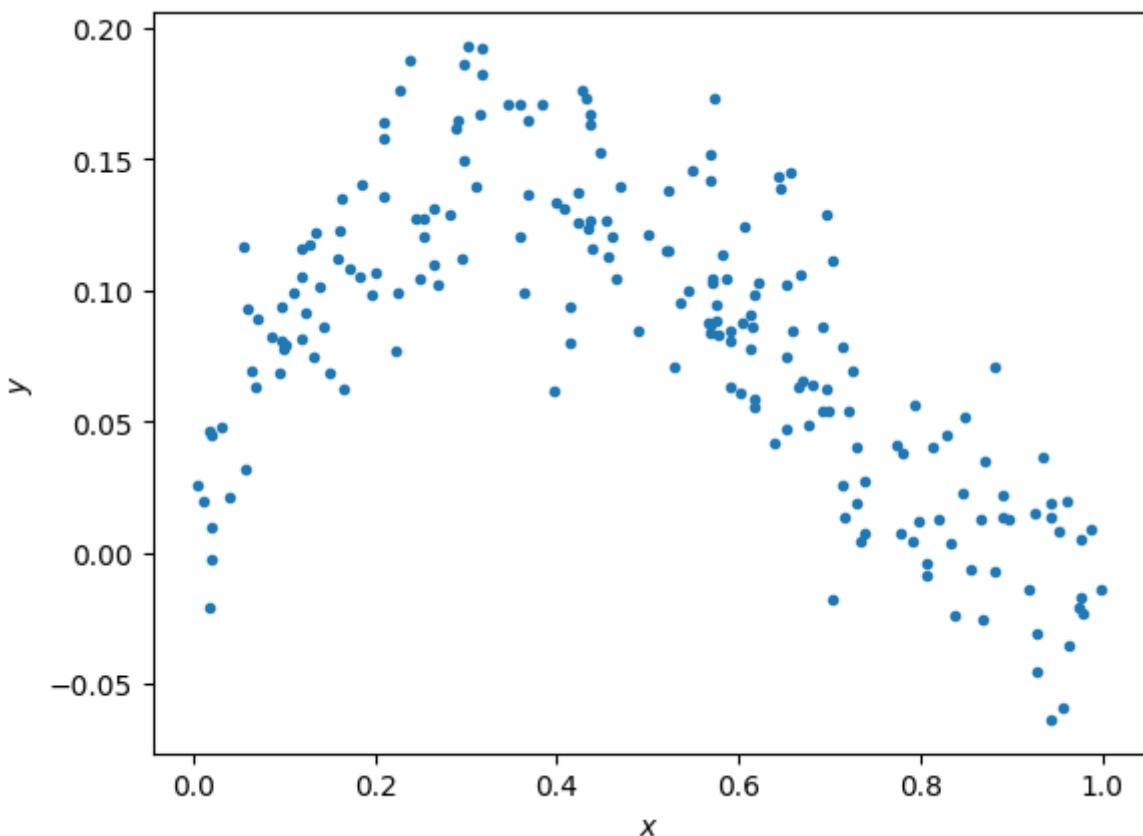
## Data generation

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 [ ]: 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[ ]: 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) Uniform probability distribution.
- (2) Normal (Gaussian) distribution.

## Fitting data to the model (5 points)

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

```
In [ ]: # 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,
theta = np.zeros(2) # please modify this line
```

```

theta = ((np.linalg.inv((xhat).dot(xhat.T)).dot(xhat.dot(y)))).T
# print(xhat.dot(y))
# print(xhat.shape, y.shape)
print(theta)

# use the transpose of what we got in class because in class x was horizontally
# try to match dims

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

[-0.10599633  0.13315817]

```

In [ ]: *# Plot the data and your model fit.*

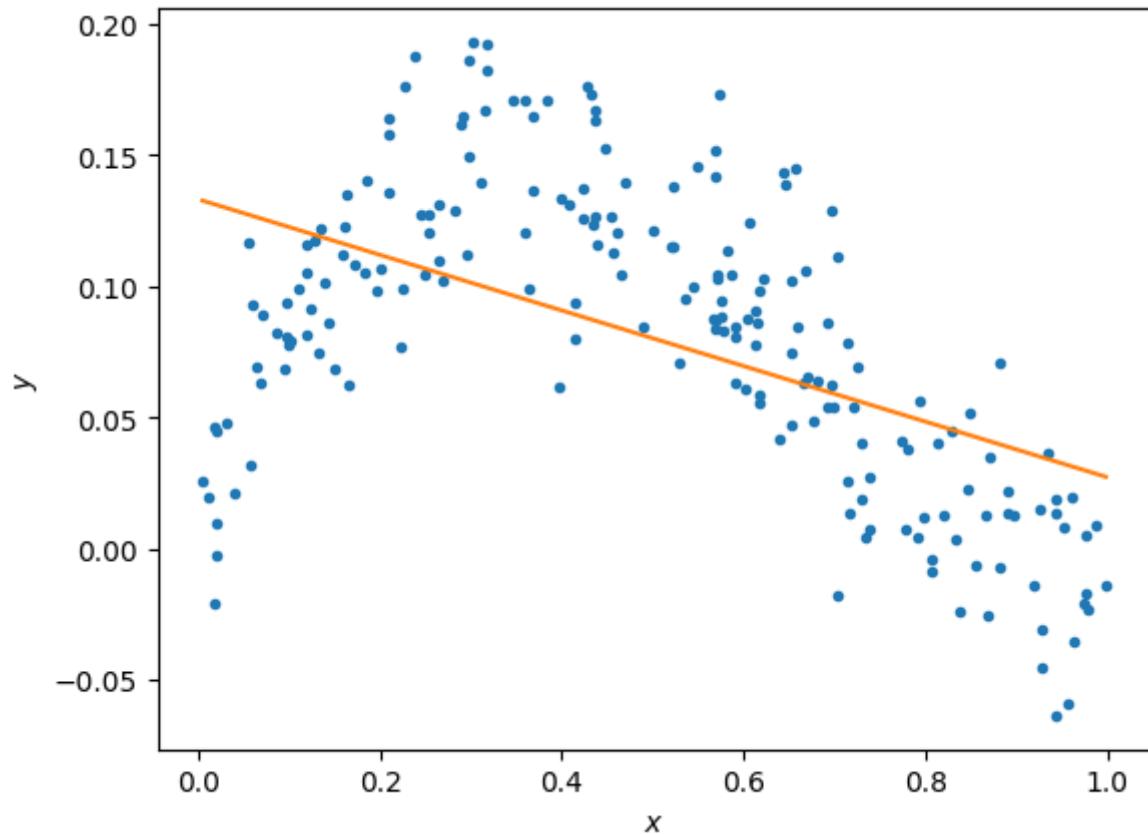
```

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[ ]: [`<matplotlib.lines.Line2D at 0x12c8a99f0>`]



## QUESTIONS

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

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

## ANSWERS

(1) Underfit.

(2) Change the order of the polynomial, as a higher order polynomial model will better fit the data due to its curvature.

## 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 [ ]: N = 5
xhats = []
thetas = []

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

# GOAL: create a variable thetas.
# thetas is a list, where thetas[i] are the model parameters for the polynomial
# 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
# ... etc.

# hint: add to the vstack

# vstack xhats with x^3, will add to the back
# vstack(x^n, xhat)

# 2 then 3 then 4 then 5 using a loop

xhats.append(xhat)
temp = np.zeros(2)
temp = ((np.linalg.inv((xhat).dot(xhat.T)).dot(xhat.dot(y)))).T
thetas.append(temp)

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

    print(i, thetas[i])
    # print("xhat shape", xhats.shape)
    print("theta shape", thetas[i].shape)

pass

# 0 : x, 1
# 1: x^2, x, 1
# 2: x^3, x^2, x, 1
# 3: x^4, x^3, x^2, x, 1
```

```
# 4: x^5, x^4, x^3, x^2, x, 1
# ===== #
# END YOUR CODE HERE #
# ===== #

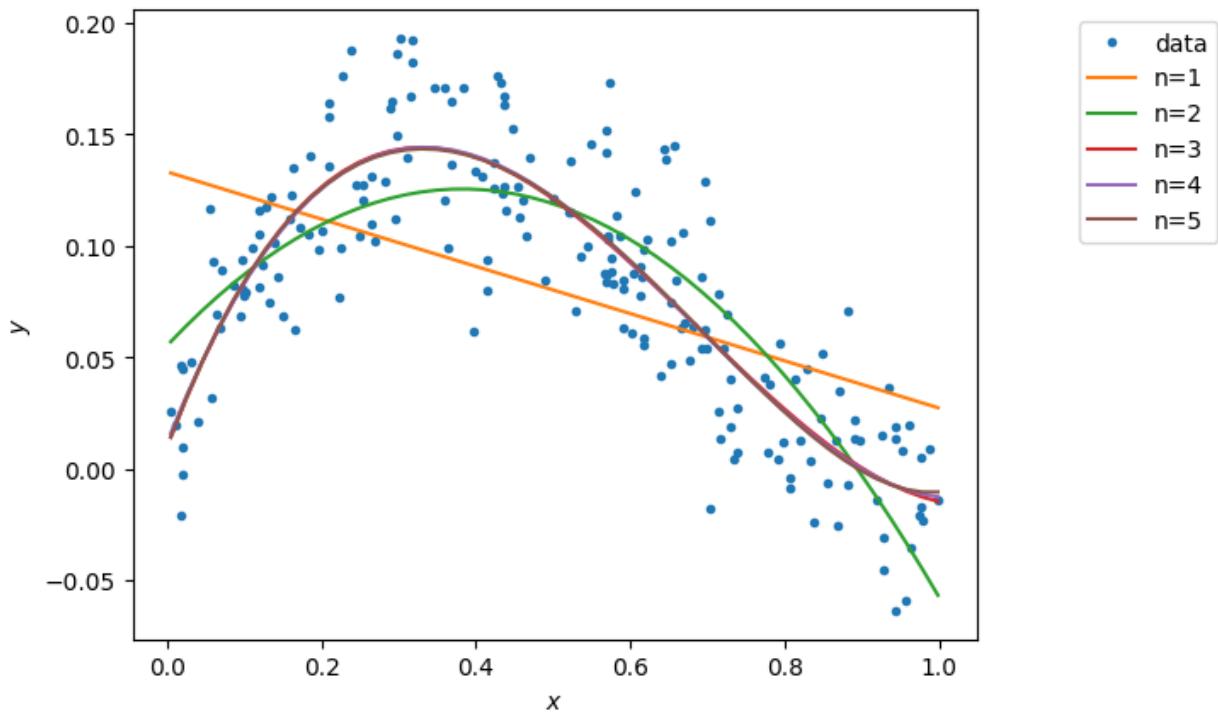
1 [-0.48023061  0.36743967  0.05521084]
theta shape (3,)
2 [ 0.8843808 -1.82077417  0.91178032  0.00979068]
theta shape (4,)
3 [ 0.14080037  0.60466289 -1.64250929  0.87250485  0.01175321]
theta shape (5,)
4 [ 0.52432592 -1.164568      1.76052438 -2.07430275  0.93373916  0.009716  ]
theta shape (6,)
```

```
In [ ]: # 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 [ ]: # ===== #
# 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
training_errors = [0] * 5
# MSE

for i in range(5):
    training_errors[i] = 0.5*(y.T.dot(y) - 2*(y.T).dot(xhats[i].T).dot(thetas[i]))

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

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

Training errors are:  
[0.23799610883626965, 0.10924922209268251, 0.08169603801102243, 0.08165353735294645, 0.08161479195520172]

## QUESTIONS

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

## ANSWERS

(1) Polynomial with degree 5.

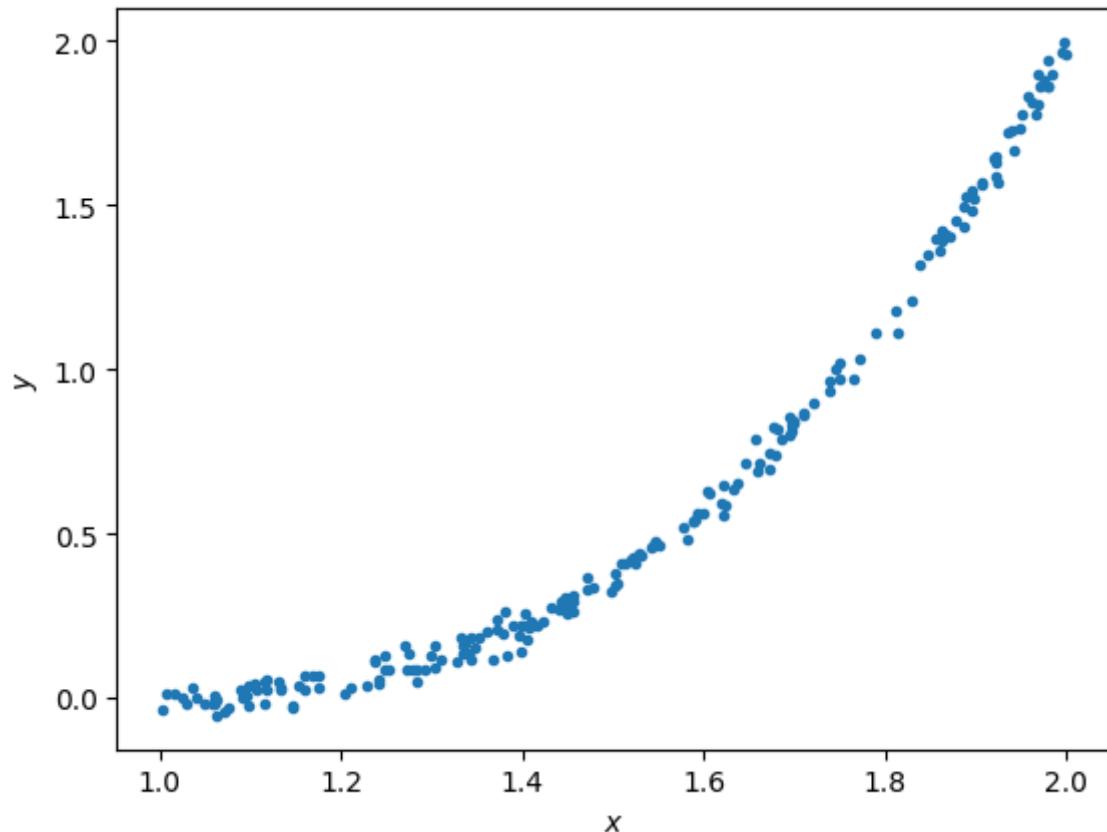
(2) As the polynomial order increases, the model can better fit the data, eventually passing through each point to make the training data equal to zero. However, this may lead to overfitting.

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



```
In [ ]: 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)

```

```

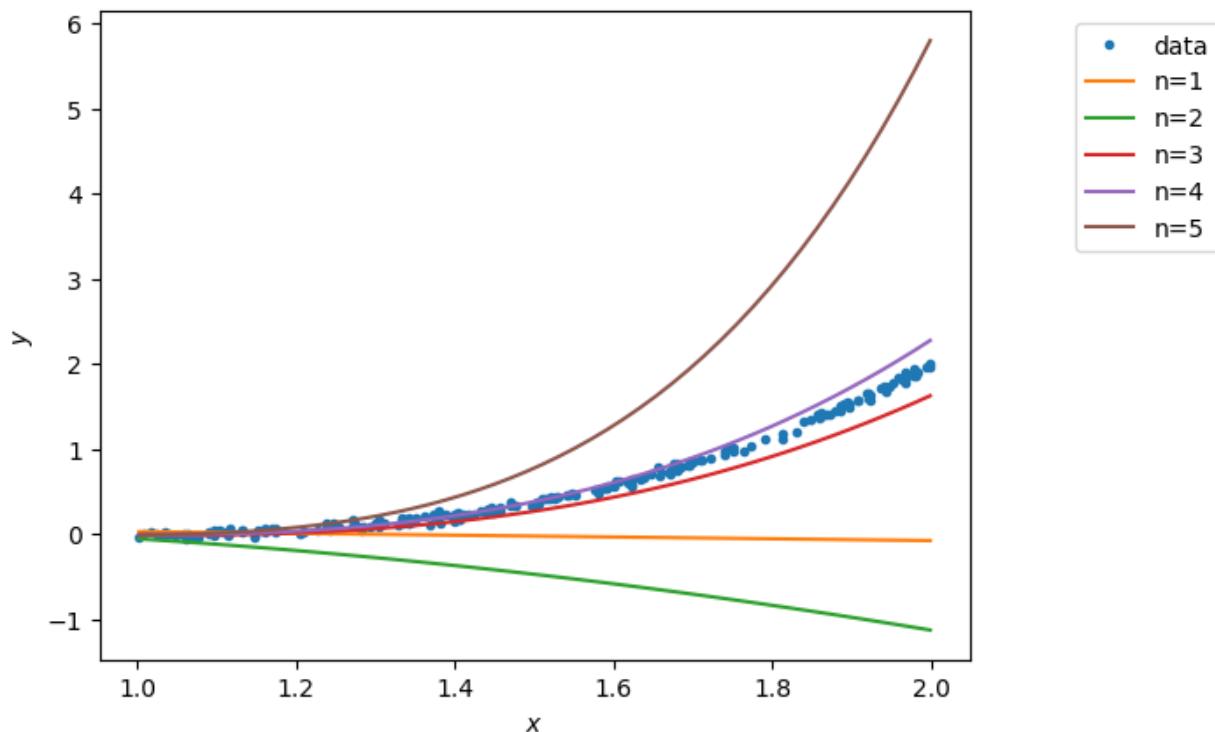
In [ ]: # Plot the data
f = plt.figure()
ax = f.gca()
ax.plot(x, y, '.')
ax.set_xlabel('$x$')
ax.set_ylabel('$y$')

# 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 [ ]: 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
testing_errors = [0] * 5
# MSE

for i in range(5):
    testing_errors[i] = 0.5*(y.T.dot(y) - 2*(y.T).dot(xhats[i].T).dot(thetas[i]))

pass

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

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

Testing errors are:

```
[80.86165184550593, 213.19192445058508, 3.125697108313929, 1.187076519555475
3, 214.91021831758638]
```

## QUESTIONS

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

## ANSWERS

- (1) Polynomial of degree 4.
- (2) Order 5 polynomial models do not generalize well because such a high order polynomial is prone to overfitting.