

EE 239 AS HW1

Problem 1

a) Since $AA^T = I$, the rows & columns will be mutually orthogonal $\rightarrow A$ is an orthogonal matrix.

$$\text{i. let } A = \begin{bmatrix} 4/5 & 3/5 \\ 3/5 & -4/5 \end{bmatrix}, |A - 3I| = \begin{vmatrix} 4/5 - 3 & 3/5 \\ 3/5 & -4/5 - 3 \end{vmatrix} = 0$$

$$(4/5 - 3)(-4/5 - 3) - 9/25 = 0$$

$$\frac{-16}{25} - \frac{4}{5}x + \frac{4}{5}x + x^2 - 9/25 = 0 \rightarrow -1 + x^2 = 0, x^2 = 1,$$

Eigenvectors:

$$Av = v \rightarrow \begin{bmatrix} 4/5 & 3/5 \\ 3/5 & -4/5 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} x \\ y \end{bmatrix}$$

$$\boxed{\begin{array}{l} x_1 = 1 \\ x_2 = -1 \end{array}} \quad \boxed{\text{eigenvalues}}$$

$$4/5x + 3/5y = x$$

$$\left(\frac{3}{5}x - 4/5y = y \rightarrow \frac{3}{5}x = 9/5y \rightarrow y = x/3 \right) \quad (x, y \text{ are elements of the eigenvector})$$

$$\rightarrow 4/5x + \left(\frac{3}{5}\right)(x/3) = x \rightarrow x = x \rightarrow \text{unit vector s.t. } y = x/3$$

$$\rightarrow x^2 + \frac{x^2}{9} = 1 \rightarrow \frac{10x^2}{9} = 1, x = \sqrt{\frac{9}{10}} \rightarrow x = \frac{3}{\sqrt{10}}, y = \frac{1}{\sqrt{10}}$$

In [eigenvector corresponding to $\lambda_1 = 1$ is $\boxed{\begin{bmatrix} \frac{3\sqrt{10}}{10} \\ \frac{1}{\sqrt{10}} \end{bmatrix}}$]

the eigenvalues are complements of each other (or complex conjugates if they were complex), the eigenvectors are normal to each other. They also are basis (an orthonormal basis) for \mathbb{R}^2

$$Av = -v$$

$$\rightarrow \begin{bmatrix} 4/5 & 3/5 \\ 3/5 & -4/5 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} -x \\ -y \end{bmatrix}$$

$$4/5x + 3/5y = -x$$

$$3/5x - 4/5y = -y \rightarrow \frac{3}{5}x = -y/5, 3x = -y, \text{ so } x = -y/3$$

$$\rightarrow \frac{3}{5}(-y/3) - \frac{4}{5}y = -y \rightarrow -y = -y \rightarrow \text{unit vec s.t. } x = -y/3$$

$$y^2 + y^2/9 = 1 \rightarrow y = \frac{3}{\sqrt{10}}, x = -\frac{1}{\sqrt{10}}$$

eigenvector corresponding to $\lambda_2 = -1$ is $\boxed{\begin{bmatrix} -\frac{\sqrt{10}}{10} \\ \frac{3\sqrt{10}}{10} \end{bmatrix}}$

Rank 1

a)

~~Since A is an orthogonal matrix,~~

ii. We have $Av = \lambda v$. The L_2 norms on both sides must therefore be equal:

$\|Av\| = \|\lambda v\|$. Since λ is a (potentially complex) scalar, we can take it out:

\rightarrow complex modulus/norm of eigenvalue

$$\|Av\| = \|\lambda\| \|v\|$$

since $\|x\| = \sqrt{x^T x}$, we have

$$\sqrt{(Av)^T Av} = \|\lambda\| \|v\|$$

$$\sqrt{v^T A^T Av} = \|\lambda\| \|v\|$$

$$\sqrt{v^T I v} = \|\lambda\| \|v\|$$

$$\Rightarrow \sqrt{v^T v} = \|\lambda\| \|v\|$$

$$\|v\| = \|\lambda\| \|v\| \text{ so } \|\lambda\| = 1 \text{ to keep this true.}$$

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

So for all possible (real or complex) values of λ_1 and λ_2 , we're shown that $\lambda_1, \lambda_2 \neq 1$. Therefore, our assumption that $x \neq 0$ is incorrect, so $x^T y = 0$.

iii. We have $x^T y = x^T I y = x^T A^T A y$

$$= (Ax)^T (Ay) = (\lambda_1 x)^T (\lambda_2 y) = \lambda_1 \lambda_2 x^T y$$

so $x^T y = \lambda_1 \lambda_2 x^T y$. Assume $x^T y \neq 0$, this means that $\lambda_1 \lambda_2 \neq 1$. But if $\lambda_1, \lambda_2 \in \mathbb{R}$ and $\lambda_1 \neq \lambda_2 \in \mathbb{R}$, then since $\|\lambda\| = 1$ and $\lambda_1 \neq \lambda_2$ (distinct), $\lambda_1, \lambda_2 \neq 1$, so we have a contradiction and $x^T y = 0$. Now, let $\lambda_1 \in \mathbb{C}$ and $\lambda_2 \in \mathbb{C}$ or $\lambda_1 \in \mathbb{C}$ and $\lambda_2 \in \mathbb{R}$, the product of $\lambda_1, \lambda_2 \neq 1$ in this case, since one λ is complex and the other is not, so $x^T y = 0$ here also. Now, let $\lambda_1 \in \mathbb{C}$ and $\lambda_2 \in \mathbb{C}$. If λ_1 and λ_2 are conjugates, then λ_1, λ_2 will be complex, so $\lambda_1 \lambda_2 \neq 1$ and $x^T y = 0$. If λ_1 and λ_2 are complex conjugates, then let $\lambda_1 = x + yi$ and $\lambda_2 = x - yi$. $\lambda_1 \lambda_2 = x^2 - y^2$, and we know that $\|\lambda\| = \|\lambda_1\| = \|\lambda_2\| = \sqrt{x^2 + y^2} = 1$, so we have $x^2 - y^2 = 1$ and $x^2 + y^2 = 1$, so $x^2 = 1$ and $y = 0$. But if $y = 0$, then $x^T y = 0$ here also, contradicting our assumption that λ_1 and λ_2 were complex.

iv. In general, the inner products are preserved: $(Ax)^T (Ay)$

$$= x^T A^T A y = x^T I y = x^T y, \text{ meaning that the transformation}$$

A corresponds to a rotation or reflection, meaning that the vector x may be rotated by some degree θ or reflected about some axis.

Problem 1

- b)
- Left-singular values of $A = \text{eigenvectors of } AA^T$ (citation: p.43 of Olshausen)
right-singular vectors of $A = \text{eigenvectors of } A^TA$ (citation: p.43 of Olshausen)
 - (Number) singular values of $A = \sqrt{\text{eigenvalues}(A^TA)}$ (citation: p.43 of Olshausen)
 $= \sqrt{\text{eigenvalues}(AA^T)}$
- (SVD: $A = UDV^T$ where cols of U = left-singular vectors of A , cols of V = right-singular vectors, and D = diag matrix where diagonal are singular values of A .)
- c.
- False, there's at most n distinct values
 - False: If $Av_1 = \gamma_1 v_1$ & $Av_2 = \gamma_2 v_2$, then $A(v_1 + v_2) = Av_1 + Av_2 = \gamma_1 v_1 + \gamma_2 v_2$ generally $\neq \gamma_n(v_1 + v_2)$.
 - True, we have $x^T Ax \geq 0 \rightarrow x^T x \geq 0$
 $\rightarrow x^T x \geq 0$, since $x \neq \vec{0}$,
 \checkmark ✓
 - False. By rank-nullity, we have $\text{rank}(A) + \text{nullity}(A) = n$.
If $\text{nullity}(A) = x \geq 0$, then there are $n-x$ eigenvalues zero eigenvalues, and $\text{rank}(A) = n - \text{nullity}(A) = n - x$.
 - True,
If $Av_1 = \gamma_1 v_1$ & $Av_2 = \gamma_2 v_2$
 $A(v_1 + v_2) = Av_1 + Av_2 = \gamma_1 v_1 + \gamma_2 v_2 = \gamma_1(v_1 + v_2)$

Date: _____

Problem 2

a)

$$\text{i. } p(H50 \mid \text{tail}) = \frac{p(\text{tails} \mid H50) p(H50)}{p(\text{tail})} = \frac{\left(\frac{1}{2}\right)^2}{\cancel{p(T, H50) + p(T, H60)}} \\ = \left(\frac{1}{2}\right)^2$$

$$\frac{p(T \mid H50) p(H50) + p(T \mid H60) p(H60)}{\left(\frac{1}{2}\right)^2 + \left(\frac{2}{5}\right)\left(\frac{1}{2}\right)} = \frac{\frac{1}{4}}{\frac{1}{4} + \frac{2}{5}} = \frac{\left(\frac{1}{4}\right)}{\left(\frac{9}{20}\right)} = \boxed{\frac{5}{9}}$$

$$\text{ii. } p(H50 \mid \overset{a}{T} \overset{b}{H} \overset{c}{H} \overset{d}{H})$$

$$= \frac{p(T \overset{a}{H} \overset{b}{H} \overset{c}{H} \overset{d}{H} \mid H50) p(H50)}{p(T \overset{a}{H} \overset{b}{H} \overset{c}{H} \overset{d}{H})} = \frac{p(T \mid H50) (p(H \mid H50))^3 p(H50)}{p(T \overset{a}{H} \overset{b}{H} \overset{c}{H} \overset{d}{H} \mid H50) p(H50) + p(T \overset{a}{H} \overset{b}{H} \overset{c}{H} \overset{d}{H} \mid H60) p(H60)}$$

$$= \frac{p(T \mid H50) p(H \mid H50)^3 p(H50)}{p(T \mid H50) p(H \mid H50)^3 p(H50) + p(T \mid H60) p(H \mid H60)^3 p(H60)}$$

↳ independence of flips

$$= \frac{(0.5)(0.5)^3(0.5)}{\frac{1}{2^3} + \left(\frac{2}{5}\right)\left(\frac{3}{5}\right)^3\left(\frac{1}{2}\right)} = \frac{\frac{1}{8}}{\frac{1}{8} + \left(\frac{1}{5}\right)\left(\frac{27}{125}\right)}$$

$$= \frac{1}{32} \quad \boxed{0.41974}$$

$$\left(\frac{1}{32} + \frac{27}{625}\right) =$$

Problem 2

a) iii. $p(H_50 | H_9 T_1) = \frac{p(H_9 T_1 | H_50) p(H_50)}{p(H_9 T_1)}$

~~$$P(H_50 | H_9 T_1) = \frac{(0.5)^{10} \left(\frac{1}{3}\right)}{p(H_9 T_1) p(H_50)}$$~~

$$= \frac{(0.5)^{10} \left(\frac{1}{3}\right)}{\sum_i p(H_9 T_1 | i) p(i)} = \frac{(0.5)^{10} \left(\frac{1}{3}\right)}{(0.5)^{10} \left(\frac{1}{3}\right) + \left(\frac{11}{20}\right)^9 \left(\frac{9}{20}\right) \left(\frac{1}{3}\right) + \left(\frac{12}{20}\right)^9 \left(\frac{8}{20}\right) \left(\frac{1}{3}\right)}$$

i: H_{50}, H_{55}, H_{60}

$$= \frac{\left(\frac{1}{2^{10}}\right) \left(\frac{1}{3}\right)}{\boxed{0.1379 = p(H_50 | H_9 T_1)}}$$

define as

$$p(H_9 T_1) \leftarrow \left[\left(\frac{1}{2^{10}}\right) \left(\frac{1}{3}\right) + \left(\frac{11}{20}\right)^9 \left(\frac{9}{20}\right) \left(\frac{1}{3}\right) + \left(\frac{3}{5}\right)^9 \left(\frac{2}{5}\right) \left(\frac{1}{3}\right) \right]$$

$$\hookrightarrow 0.00236 \quad \boxed{0.1379 = p(H_50 | H_9 T_1)}$$

$$p(H_{55} | H_9 T_1) = \frac{p(H_9 T_1 | H_{55}) p(H_{55})}{p(H_9 T_1)} = \frac{\left(\frac{11}{20}\right)^9 \left(\frac{9}{20}\right) \left(\frac{1}{3}\right)}{p(H_9 T_1)}$$

$$p(H_{60} | H_9 T_1) = \frac{p(H_9 T_1 | H_{60}) p(H_{60})}{p(H_9 T_1)} = \frac{\left(\frac{3}{5}\right)^9 \left(\frac{2}{5}\right) \left(\frac{1}{3}\right)}{p(H_9 T_1)} = \boxed{0.5694 = p(H_{60} | H_9 T_1)}$$

check: $0.1379 + 0.2927 + 0.5694 = 1$



b) Problem 2

b) define I =test positive, P =present.

$$P(I|P) = 0.99$$

$$P(I|\sim P) = 0.1$$

$$P(P) = 0.01$$

- find $P(P|I)$.

$$P(P|I) = \frac{P(I|P)P(P)}{P(I)}$$

$$P(I) = P(I, P) + P(I, \sim P)$$

$$= P(I|P)P(P) + P(I|\sim P)P(\sim P)$$

$$= \frac{P(I|P)P(P)}{P(I|P)P(P) + P(I|\sim P)P(\sim P)}$$

$$= \frac{(0.99)(0.01)}{(0.99)(0.01) + (0.1)(0.99)} = 0.0909$$

This surprisingly low probability does make sense if we consider the high false positive rate, 10%. This would mean that 10% of 99% of the population would get a false positive, a pretty large amount. If our false positive rate were lower, such as 0.1%, then our probability would go up to about 91%.

$$c) E[Ax+b] = E[Ax] + E[b] = E[Ax] + b$$

$$= AE[x] + b \quad (\text{due to linearity of expectation})$$

$$d) \text{cov}(Ax+b) = E((Ax+b - E(Ax+b))(Ax+b - E(Ax+b))^T)$$

$$= E((Ax+b - AE(x)+b)(Ax+b - AE(x)+b)^T)$$

$$= E((Ax - AE(x))(Ax - AE(x))^T) = E(A(x - E(x))(A(x - E(x))^T)$$

$$= E(A(x - E(x))(x - E(x))^TA^T) \rightarrow \text{linearity of expectation} \rightarrow AE((x - E(x))(x - E(x))^TA^T)$$

$$= A \text{cov}(x) A^T$$

$$2^5 \quad (5 \text{ n.s. } - x) \quad \frac{\overbrace{1}^{2^5} + \left(\frac{1}{5}\right)\left(\frac{27}{125}\right)}{\overbrace{1}^{2^5} + \frac{27}{625}} = \frac{(3^2)}{32} = 0.4197$$

Problem 3

3) a) $\nabla_x \times^T A y$

$$\cancel{A y} + \cancel{x^T A} \cancel{0} = \cancel{A y}$$

$$\rightarrow \nabla_x \times^T A y = \cancel{x}(Ay)^T x = \boxed{Ay}$$

b) $\nabla_y \times^T A y = (\cancel{x^T A})^T = \boxed{A^T x}$

c) $\nabla_A \times^T A y = \cancel{y} \cancel{x^T} \cancel{x^T y} + \cancel{x^T y} = \cancel{2x^T y^T}$

$$\nabla_A \times^T A y = \cancel{0} \cancel{y} \cancel{x^T y}$$

$$x^T A y = \sum_{i=1} \sum_{j=1} x_i y_j a_{ij} \rightarrow \nabla_{a_{ij}} \sum_{i=1} \sum_{j=1} x_i y_j a_{ij} \cancel{= 0}$$

$\cancel{= 0}$ (cancel)
 $\cancel{= x^T y}$ (cancel)
 $\hookrightarrow x_i y_j$ (match)
 $A: d \times m$

Problem 3

$$d) f = \mathbf{x}^T A \mathbf{x} + b^T \mathbf{x}$$

$$\nabla_{\mathbf{x}} \mathbf{x}^T A \mathbf{x}$$

$$\rightarrow \sum_i \sum_j x_i A_{ij} x_j$$

$$\frac{d}{dx_i} \cdot \sum_i \sum_j x_i A_{ij} x_j$$

when $i=j=1 \rightarrow x_1 A_{11} x_1$
 $\rightarrow 2 A_{11} x_1$ is the alternative

when $i=1, j \neq 1, x_1 A_{1j} x_j$

$$d/dx_i = A_{1j} x_j$$

when $i \neq 1, j=1, \frac{d}{dx_i} (x_i A_{i1} x_1) = x_i A_{i1}$

$$\rightarrow \sum_{j \neq 1} A_{ij} x_j + \sum_{i \neq 1} x_i A_{ij} + 2 A_{11} x_1$$

$\rightarrow \sum_j A_{ij} x_j + \sum_i x_i A_{ij}$ is the general term,

so we have $A \mathbf{x} + A^T \mathbf{x}$

$$\nabla_{\mathbf{x}} b^T \mathbf{x} \rightarrow \sum_i b_i x_i \rightarrow \frac{d}{dx_i} = b_i, \text{ so } \nabla_{\mathbf{x}} b^T \mathbf{x} = b$$

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

~~e) let $C = AB$. Then $C_{ij} = \sum_k A_{ik} B_{kj} \rightarrow \nabla_{A_{ik}} \sum_k A_{ik} B_{kj} = B_{kj}$~~

~~e) let $C = AB$. Then~~

~~$C_{ij} = \sum_k A_{ik} B_{kj}$. Since we're taking~~

~~tr(C), we have the diagonal elements~~

~~of C: $C_{nn} = \sum_k A_{nk} B_{kn}$.~~

Now, $\nabla_{A_{nk}} \sum_k A_{nk} B_{kn} = B_{kn}$, so each element will be B_{kn} . So

$$\nabla_A \text{tr}(AB) = \boxed{B^T}$$

~~$\nabla_{A_{nk}} \sum_k A_{nk} B_{kn} = B_{kn}$~~

(B is the transpose of B)
~~(dim of B is n)~~
~~matrix form of~~
~~(A)~~

4)

$$\min_w \frac{1}{2} \sum_{i=1}^n \|y^i - w x^i\|^2$$

(Frobenius Norm from above)

First, replace $\|A\|^2$ w/ $\text{tr}(A A^T)$:

$$\min_w \frac{1}{2} \sum_{i=1}^n \text{tr}[(y^i - w x^i)(y^i - w x^i)^T]$$

$$= \frac{1}{2} \sum_{i=1}^n \text{tr}[y^i y^{iT} - \underbrace{y^i x^{iT} w^T}_{\text{combine}} - w x^i y^{iT} + w x^i (w x^i)^T]$$

$$= \frac{1}{2} \sum_{i=1}^n \text{tr}[y^i y^{iT} - 2 w x^i y^{iT} + w x^i (w x^i)^T]$$

$$= \frac{1}{2} \sum_{i=1}^n \text{tr}(y^i y^{iT}) - \sum_{i=1}^n (w x^i y^{iT}) + \frac{1}{2} \sum_{i=1}^n (w x^i (w x^i)^T)$$

$$\frac{d}{dw} (\dots)$$

$$= \frac{d}{dw} \left(- \sum_{i=1}^n \text{tr}(w x^i y^{iT}) + \frac{d}{dw} \left(\frac{1}{2} \sum_{i=1}^n \text{tr}(w x^i (w x^i)^T) \right) \right)$$

$$= - \sum_{i=1}^n (x^i y^{iT})^T + \frac{1}{2} \sum_{i=1}^n (W x^i x^{iT} + W x^i x^{iT})$$

$$= - \sum_{i=1}^n (x^i y^{iT})^T + \frac{1}{2} \sum_{i=1}^n 2 W x^i x^{iT}$$

$$\rightarrow \sum_{i=1}^n (x^i y^{iT})^T = \sum_{i=1}^n W x^i x^{iT}$$

$$\sum_{i=1}^n (x^i y^{iT})^T = w \sum_{i=1}^n x^i x^{iT}$$

$$\rightarrow \sum_{i=1}^n y^i x^{iT} = w \sum_{i=1}^n x^i x^{iT}$$

$$\rightarrow \boxed{w = \left(\sum_{i=1}^n x^i x^{iT} \right)^{-1} \sum_{i=1}^n y^i x^{iT}}$$

*extending
or using matrix X = design matrix where
row i is the ith example,
and y is a vector of labels, we have*

$$\boxed{\hat{w} = (X^T X)^{-1} X^T y}$$

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 239AS, Winter Quarter 2018, Prof. J.C. Kao, TAs C. Zhang and T. Xing

```
import numpy as np
import matplotlib.pyplot as plt

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

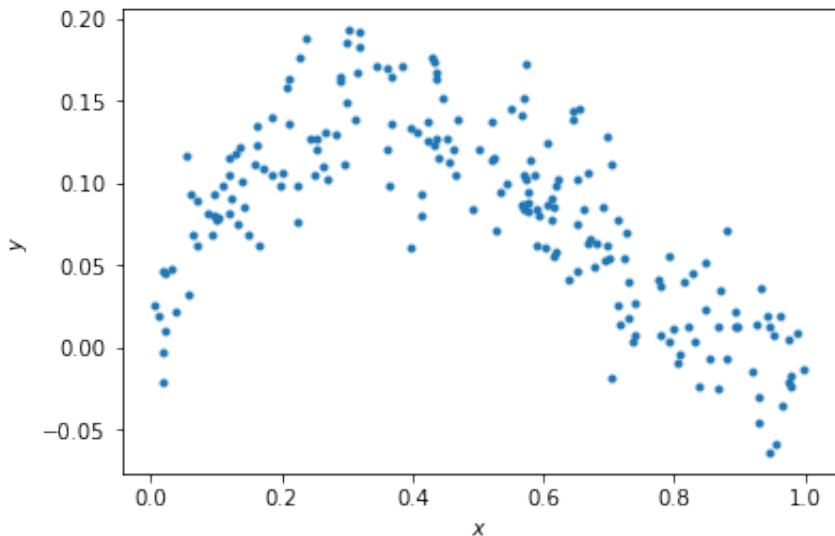
Data generation

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

```
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('$y$')
```

```
<matplotlib.text.Text at 0x10621fd68>
```



QUESTIONS:

Write your answers in the markdown cell below this one:

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

ANSWERS:

- (1) x is drawn from a uniform distribution with parameters $a = 0$ and $b = 1$.
- (2) ϵ is drawn from a Gaussian with 0 mean and standard deviation 0.03 .

Fitting data to the model (5 points)

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

```

# 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 = (X^TX)^_1 (X^Ty)
x_ = xhat.T
# asserts to make sure the shapes are as expected
assert x_.T.shape[1] == x_.shape[0] and x_.T.shape[1] == y.shape[0]
theta = np.linalg.inv(x_.T.dot(x_)).dot(x_.T.dot(y))
assert theta.shape[0] == 2
# ===== #
# END YOUR CODE HERE #
# ===== #

```

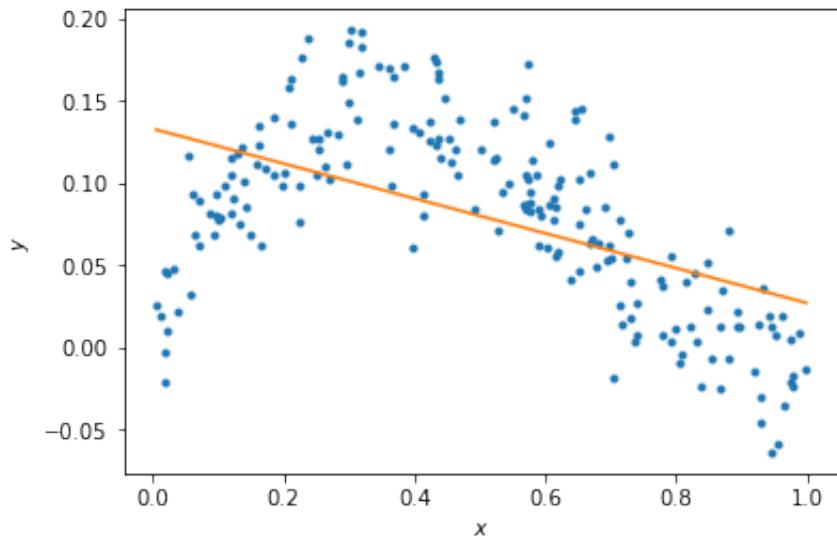
```

# 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))

```

[<matplotlib.lines.Line2D at 0x1063d1550>]



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 mode underfits the data.
- (2) We can reduce underfitting py fitting a polynomial of higher degree, instead of a degree 1 polynomial (linear function).

Fitting data to the model (10 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.

```

N = 5
xhats = []
thetas = []

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

# GOAL: create a variable thetas.
# thetas is a list, where theta[i] are the model parameters for the polynomial
# fit of order i+1.
# i.e., thetas[0] is equivalent to theta above.
# i.e., thetas[1] should be a length 3 np.array with the coefficients of the
# x^2, x, and 1 respectively.
# ... etc.
# cur matrix will hold the features generated up to the current highest degree
# polynomial.
cur_matrix = []
cur_matrix.insert(0, np.ones_like(x))
for i in range(5):
    # fit a polynomial of degree i + 1
    # first, generate the features: x^(i + 1) down to x^0 (basically the bias
    # units).
    cur_matrix.insert(0, np.array(x**(i + 1)))
    x_ = np.array(cur_matrix).T
    # verify shapes are as expected
    assert x_.T.shape[1] == x_.shape[0] and x_.T.shape[1] == y.shape[0]
    # least squares to find theta
    cur_theta = np.linalg.inv(x_.T.dot(x_)).dot(x_.T.dot(y))
    if i == 0:
        assert cur_theta.all() == theta.all() # the ax + b model should match
        the previous result
    thetas.append(cur_theta)
    # keep track of the cur_matrix in xhats for later error checking
    xhats.append(np.array(cur_matrix))

for idx, val in enumerate(thetas):
    assert len(val) == idx + 2
# ===== #
# END YOUR CODE HERE #
# ===== #

```

```

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

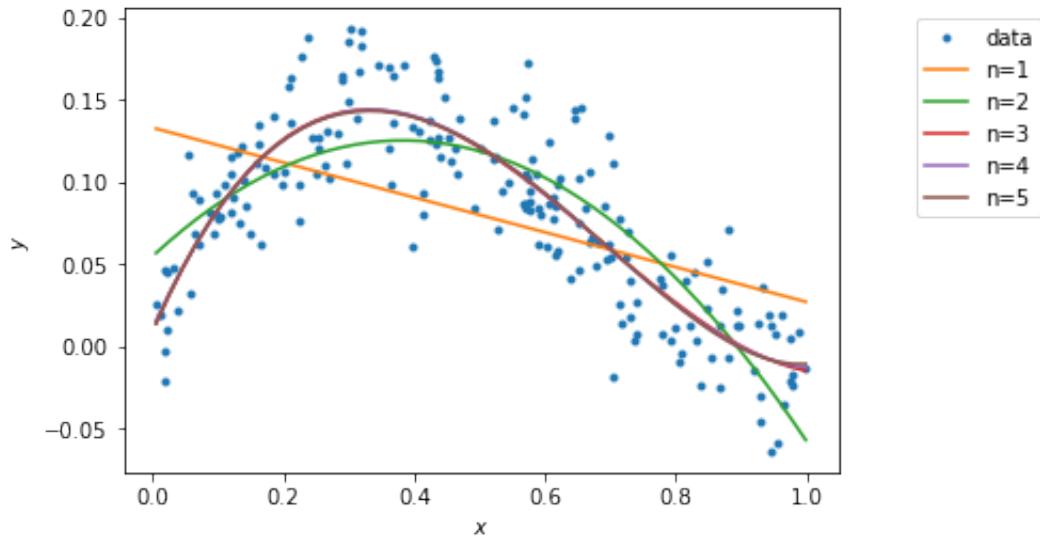
# Plot the regression lines
plot_xs = []
for i in np.arange(N):
    if i == 0:
        plot_x = np.vstack((np.linspace(min(x), max(x), 50), np.ones(50)))
    else:
        plot_x = np.vstack((plot_x[-2]***(i+1), plot_x))

    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 (10 points)

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

```

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 theta, x_ in zip(thetas, xhats):
    # get the predictions and calculate MSE
    predictions = x_.T.dot(theta)
    mse = sum((predictions - y)**2)/len(y)
    training_errors.append(mse)

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

# using MSE loss: 1/N sum(y - y_pred)^2
print ('Training errors are: \n', training_errors)

```

Training errors are:
[0.0023799610883627007, 0.001092492220926853, 0.00081696038011053683,
0.00081653537352969758, 0.00081614791955252942]

QUESTIONS

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

ANSWERS

- (1) The polynomial of degree $n = 5$ has the best (lowest) training error.
- (2) Higher degree polynomials have more free parameters to fit to the dataset, so they can better fit the data. Intuitively, a higher degree allows the polynomial to "wiggle" to a much greater extent than for example a degree 1 polynomial (which is just a line), leading to a better fit on the training data.

Generating new samples and testing error (5 points)

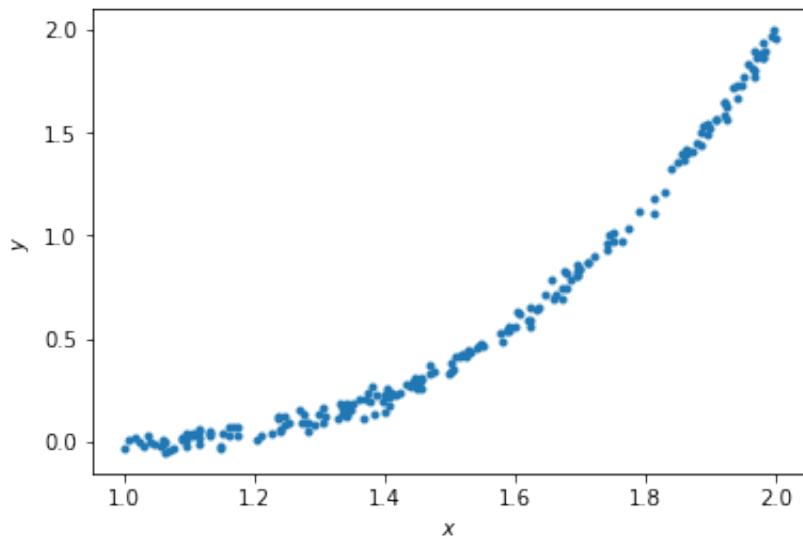
Here, we'll now generate new samples and calculate testing error of polynomial models of orders 1 to 5.

```

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('$y$')

```

<matplotlib.text.Text at 0x1066274e0>



```

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)

```

```

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

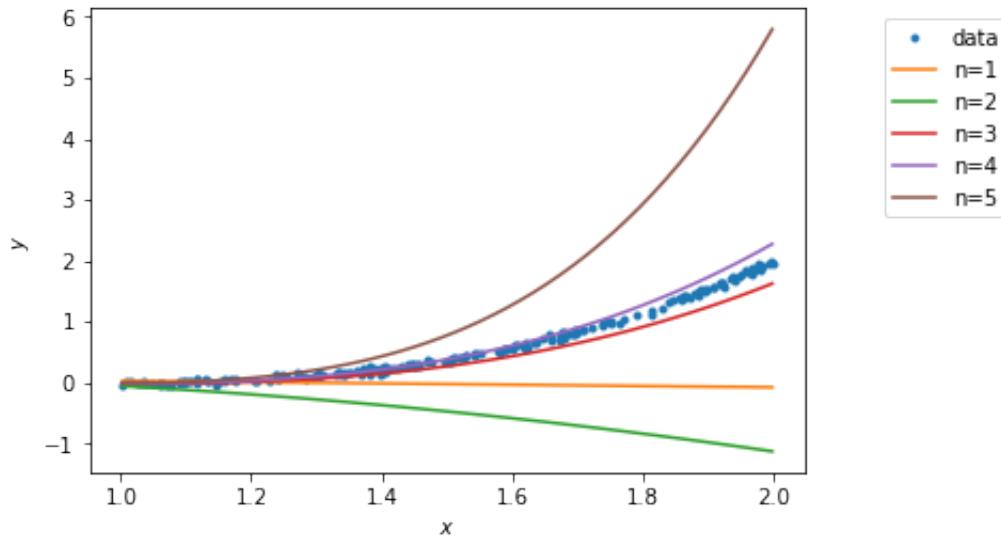
# Plot the regression lines
plot_xs = []
for i in np.arange(N):
    if i == 0:
        plot_x = np.vstack((np.linspace(min(x), max(x), 50), np.ones(50)))
    else:
        plot_x = np.vstack((plot_x[-2]***(i+1), plot_x))

    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)

```



```

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 theta, x_ in zip(thetas, xhats):
    # get the preds and calc MSE
    mse = sum((y - theta.dot(x_))**2)/len(y)
    testing_errors.append(mse)
# ===== #
# END YOUR CODE HERE #
# ===== #

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

```

```

Testing errors are:
[0.80861651845505822, 2.1319192445057884, 0.031256971083289641,
0.01187076519660833, 2.1491021807712438]

```

QUESTIONS

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

ANSWERS

- (1) The polynomial with degree $n = 4$.
- (2) The polynomial model of degree $n = 5$ did not generalize well because it overfit the training data. This means that it started to learn the noise in the training data, instead of the overall pattern of the data generating distribution from which x was drawn from. Since the polynomial was overfit on our training dataset, it performed worse when tested on data it had never seen before, indicating that the fit was very dependent on the data that was provided to it - if we trained it on a slightly different dataset, the model would have been quite different. This indicates that the degree $n = 5$ model had high variance. High variance models do not generalize well, because they tend to learn the intricacies of the data that they are trained on, but the actual testing data does not have these exact intricacies.

