

Hw 1

Amirali Omidfar

1 a)

$$i) A = \begin{bmatrix} .6 & .8 \\ .8 & -.6 \end{bmatrix}$$

$$AA^T = \begin{bmatrix} .6 & .8 \\ .8 & -.6 \end{bmatrix} \begin{bmatrix} .6 & .8 \\ .8 & -.6 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$AX = \lambda X$$

$$AX - \lambda X = 0$$

$$(A - \lambda I)X = 0$$

$$\begin{bmatrix} .6 - \lambda & .8 \\ .8 & -.6 - \lambda \end{bmatrix} X = 0$$

$$\det(A - \lambda I) = (.6 - \lambda)(-.6 - \lambda) - .8^2$$

$$\Rightarrow \lambda^2 - 1 = 0 \Rightarrow \lambda = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$\text{if } \lambda = 1 \Rightarrow \begin{bmatrix} -.4 & .8 & 0 \\ .8 & -1.6 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} -.4 & .8 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\text{Eigenvector } V_1 = \begin{bmatrix} x_1 \\ x_1/2 \end{bmatrix}$$

$$\text{if } \lambda = -1 \Rightarrow \begin{bmatrix} 1.6 & .8 & 0 \\ .8 & .4 & 0 \end{bmatrix} \Rightarrow \begin{bmatrix} 1.6 & .8 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\text{Eigenvector } V_2 = \begin{bmatrix} x_1 \\ -2x_1 \end{bmatrix}$$

$\Rightarrow |det A| = 1$ & eigenvectors are orthogonal to each other
 $|\lambda| = 1$

$$ii) \quad Av_i = \lambda_i v_i$$

$$\Rightarrow \|Av_i\|_2^2 = \|\lambda_i v_i\|_2^2$$

$$\Rightarrow \underbrace{v_i^T A^T A v_i}_I = \lambda_i^2 \|v_i\|_2^2$$

$$\Rightarrow v_i^T v_i = \lambda_i^2 \|v_i\|_2^2$$

$$\Rightarrow \|v_i\|_2^2 = \|v_i\|_2^2 \quad \checkmark$$

iii) Given v_1 & v_2 being the eigenvectors for λ_1 & λ_2

$$\Rightarrow v_1 \cdot v_2 = v_2 \cdot v_1$$

$$\Rightarrow A(v_1 \cdot v_2) = A(v_2 \cdot v_1)$$

$$\Rightarrow (Av_1) \cdot v_2 = (Av_2) \cdot v_1$$

$$\Rightarrow (\lambda_1 v_1) \cdot v_2 = (\lambda_2 v_2) \cdot v_1$$

$$\Rightarrow \lambda_1 (v_1 \cdot v_2) - \lambda_2 (v_2 \cdot v_1) = 0$$

$$\Rightarrow (\lambda_1 - \lambda_2) (v_1 \cdot v_2) = 0$$

$$\text{since } \lambda_1 \neq \lambda_2 \Rightarrow \lambda_1 - \lambda_2 \neq 0$$

$$\text{So } v_1 \cdot v_2 = 0 \Rightarrow \underline{v_1 \perp v_2} \quad \checkmark \checkmark$$

iv) Any x is subject to rotation or reflection under AX

when $|\det A| = 1$ in this case its reflection as $\det A = -1$

The length is preserved.

b)

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

$$AA^T = U \Sigma \underbrace{V V^T}_I \Sigma^T U^T = A^T A = V \Sigma^T U^T U \Sigma V^T$$

$$AA^T = \underbrace{U \Sigma \Sigma^T U^T}_{\text{diagonal matrix}} \quad A^T A = \underbrace{V \Sigma^T \Sigma V^T}_{\text{diagonal matrix}}$$

According to Eigendecomposition of singular vectors of A
 are eigenvectors of AA^T & $A^T A$.

Also λ is for $A^T A$ & AA^T are GS for A (singular values)

$$ii) A \text{ is symmetric} : \therefore AA^T = U \underbrace{\Sigma \Sigma^T}_\Lambda U^T \Rightarrow$$

λ is for AA^T & AA^T are singular values of A squared!!

u)

i) False

$I = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ has one only

ii) False

If eigenvalues are different then it doesn't hold

$$v_3 = v_1 + v_2 \Rightarrow A[v_3] = Av_1 + Av_2 = \lambda_1 v_1 + \lambda_2 v_2 \neq \lambda(v_1 + v_2) \text{ unless } \lambda_1 = \lambda_2$$

iii) True. If λ is an eigenvalue for $A \Rightarrow Ax = \lambda x$

$$\Rightarrow \begin{cases} x^T A x \gg 0 \\ x^T A x = \lambda x^T x \end{cases} \Rightarrow \begin{cases} \lambda x^T x \gg 0 \\ x^T x \gg 0 \text{ is always greater than zero} \end{cases}$$

$\Rightarrow \lambda \gg 0$ so the λ s are non-negative.

iv) False, according to the theorem: $A \in \mathbb{R}^{n \times n} \Rightarrow n = \text{Rank } A + \text{nullity of } A$

\Rightarrow if v is in the nullity of $A \Rightarrow Av = 0$ so v is eigenvector for

0 eigenvalue. Since there exists an eigenspace of degree

at least one for each eigenvalue, the number of

non-zero eigenvalues must be less than or equal to

the rank.

v) True, Assume: $A \in \mathbb{R}^{n \times n}$ & $Av = \lambda v$
 $u, v \in \mathbb{R}^{n \times 1}$

$$\Rightarrow C = A(v + u') \Rightarrow C_{ik} = \sum_{j=1}^n a_{ijk} (v_j + u'_j) = \sum a_{ijk} v_j + \sum a_{ijk} u'_j$$

$$\Rightarrow C = Av + Au' = \lambda v + \lambda u' = \lambda(v + u') \Rightarrow A(v + u') = \lambda(v + u')$$

$$2. P(H|H50) = 0.5, P(H|H60) = 0.6$$

$$i) P(H50|T) = \frac{P(T|H50)P(H50)}{P(T|H50)P(H50) + P(T|H60)P(H60)}$$

$$= \frac{(0.5)(0.5)}{(0.5)(0.5) + (0.4)(0.5)} = \frac{5}{9}$$

ii)

$$S = THHH$$

$$P(H50|S) = \frac{P(S|H50)P(H50)}{P(S|H50)P(H50) + P(S|H60)P(H60)}$$

$$= \frac{(0.5)^4(0.5)}{(0.5)^4(0.5) + (0.4)(0.6)^3(0.5)} = 0.4197$$

iii) $S = 9$ heads in 10 attempts

$$P(S|H50) = \binom{10}{9} (0.5)^9 (0.5) = 0.097$$

$$P(S|H55) = \binom{10}{9} (0.55)^9 (0.45) = 0.0207$$

$$P(S|H60) = \binom{10}{9} (0.6)^9 (0.4) = 0.0403$$

$$P(H_i|S) = \frac{P(S|H_i)P(H_i)}{\sum_{H_i=H50}^{H60} P(S|H_i)P(H_i)}$$

$$\Rightarrow P(H50|S) = \frac{(0.097)(1/3)}{0.097(1/3) + 0.0207(1/3) + 0.0403(1/3)} = 0.1372$$

$$P(H55|S) = \frac{(0.0207)(1/3)}{0.097(1/3) + 0.0207(1/3) + 0.0403(1/3)} = 0.2928$$

$$P(H60|S) = \frac{(0.0403)(1/3)}{0.097(1/3) + 0.0207(1/3) + 0.0403(1/3)} = 0.5700$$

$$b) P(+|Pr) = 0.99$$

$$P(+|NPr) = 0.1$$

$$P(NPr) = 0.99$$

$$\begin{aligned} P(Pr|+) &= \frac{P(+|Pr)P(Pr)}{P(+|Pr)P(Pr) + P(+|NPr)P(NPr)} \\ &= \frac{(0.99)(0.01)}{(0.99)(0.01) + (0.1)(0.99)} \\ &= \frac{1}{11} = 0.09 \end{aligned}$$

It means the test is not reliable when only in 9% of the time the women are pregnant given the test being positive.

$$c) E(X) = \sum_x x P(x)$$

$$\Rightarrow E(AX+b) = \sum_x (AX+b) P(x)$$

$$= A \sum_x x P(x) + b \sum_x P(x)$$

$$= A \sum_x x P(x) + b$$

$$= AE(X) + b$$

Using scaling & additivity properties
 \Rightarrow linearity of expectation with respect to one variable

$$\begin{aligned}
 2. d) \quad \text{Cov}(x) &= E((x - E(x))(x - E(x))^T) \\
 \text{Cov}(Ax+b) &= E((Ax+b - E(Ax+b))(Ax+b - E(Ax+b))^T) \\
 &= E((Ax+b - AE(x)-b)(Ax+b - AE(x)-b)^T) \\
 &= E((Ax - AE(x))(Ax - AE(x))^T) \\
 &= E([A(x - E(x))][A(x - E(x))]^T) \\
 &= E(A \underbrace{(x - E(x))(x - E(x))^T}_{\text{Cov}(x)} A^T) \\
 &= A \text{Cov}(x) A^T
 \end{aligned}$$

$$3. a) \quad x \in \mathbb{R}^n, y \in \mathbb{R}^m, A \in \mathbb{R}^{n \times m} \Rightarrow \nabla_x x^T A y?$$

$$f = \nabla_x x^T A y \Rightarrow \in \mathbb{R}^{n \times 1} \mid f_i = \frac{\partial}{\partial x_i} x^T A y$$

$$f = x^T A y = [x_1 \dots x_n] \begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nm} \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

$$f = \sum_{j=1}^m \left(\sum_{i=1}^n x_i a_{ij} \right) y_j$$

$$\frac{\partial f}{\partial x_k} = \frac{\partial}{\partial x_k} \sum_{j=1}^m (x_k a_{kj}) y_j = \frac{\partial}{\partial x_k} x_k \sum_{j=1}^m (a_{kj} y_j) = \sum_{j=1}^m a_{kj} y_j$$

$$\Rightarrow \nabla_x f = \begin{bmatrix} \sum_{j=1}^m a_{1j} y_j \\ \vdots \\ \sum_{j=1}^m a_{nj} y_j \end{bmatrix} = A y$$

b) $\nabla_y x^T A y$?

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

$$\nabla_y f = \begin{bmatrix} \sum_{i=1}^n x_i a_{i1} \\ \vdots \\ \sum_{i=1}^n x_i a_{im} \end{bmatrix} = \boxed{A^T x}$$

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

$$f = \sum_{j=1}^m \left(\sum_{i=1}^n x_i a_{ij} \right) y_j$$

$$\nabla_A f = \begin{bmatrix} \frac{\partial f}{\partial a_{11}} & \dots & \frac{\partial f}{\partial a_{1m}} \\ \vdots & & \vdots \\ \frac{\partial f}{\partial a_{n1}} & \dots & \frac{\partial f}{\partial a_{nm}} \end{bmatrix} \Rightarrow \frac{\partial f}{\partial a_{kl}} = x_k y_l$$

$$= \begin{bmatrix} x_1 y_1 & \dots & x_1 y_m \\ \vdots & & \vdots \\ x_n y_1 & \dots & x_n y_m \end{bmatrix} = \boxed{x y^T}$$

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

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

$$\nabla_x b^T x = \begin{bmatrix} \frac{\partial \sum b_i x_i}{\partial x_1} \\ \vdots \\ \frac{\partial \sum b_i x_i}{\partial x_n} \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix} = \boxed{b}$$

$$x^T A x = \sum_{j=1}^n x_j \left(\sum_{i=1}^n a_{ij} x_i \right)$$

\Rightarrow linearity of expectation with respect to one variable

$$\frac{\partial x^T A x}{\partial x_k} = \frac{\sum_i x_i (\sum_j a_{ij} x_j)}{\partial x_k} = \sum_i x_i (a_{ik}) + \sum_i x_i (a_{ki})$$

$$= A^T x + A x$$

$$\Rightarrow \nabla_x (x^T A x + b^T x) = \underline{A^T x + A x + b}$$

e) $f = \text{tr}(AB) \Rightarrow \nabla_A f?$

$$A \in \mathbb{R}^{n \times m}, B \in \mathbb{R}^{m \times n} \rightarrow C = AB, C \in \mathbb{R}^{n \times n}$$

$$C_{ke} = \sum_{i=1}^n a_{ki} b_{ie}$$

$$\text{tr}(AB) = \text{tr}(C) = \sum_{j=1}^m c_{jj} = \sum_{j=1}^m \left(\sum_{i=1}^n a_{ji} b_{ij} \right)$$

$$\nabla_A f = \begin{bmatrix} \frac{\partial \text{tr}(AB)}{\partial a_{11}} & \dots & \frac{\partial \text{tr}(AB)}{\partial a_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \text{tr}(AB)}{\partial a_{n1}} & \dots & \frac{\partial \text{tr}(AB)}{\partial a_{nn}} \end{bmatrix}$$

$$\Rightarrow \frac{\partial f}{\partial a_{ij}} = b_{ji} \quad \Rightarrow \nabla_A f = B^T$$

4 $\min_W \frac{1}{2} \sum_{i=1}^n \|y^i - W x^i\|^2$

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

$$\min_W \frac{1}{2} (Y - X W^T)^T (Y - X W^T)$$

$$\min_W \frac{1}{2} (Y^T - W X^T) (Y - X W^T)$$

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

$$\nabla_W L = \frac{1}{2} (0 - 2 Y^T X + 2 W X^T X) = 0 \quad \begin{aligned} W X^T X &= Y^T X \\ W &= Y^T X (X^T X)^{-1} \end{aligned}$$

linear_regression

January 20, 2020

0.1 Linear regression workbook

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

ECE C147/C247 Winter Quarter 2020, Prof. J.C. Kao, TAs W. Feng, J. Lee, K. Liang, M. Kleinman, C. Zheng

```
In [58]: import numpy as np
import matplotlib.pyplot as plt

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

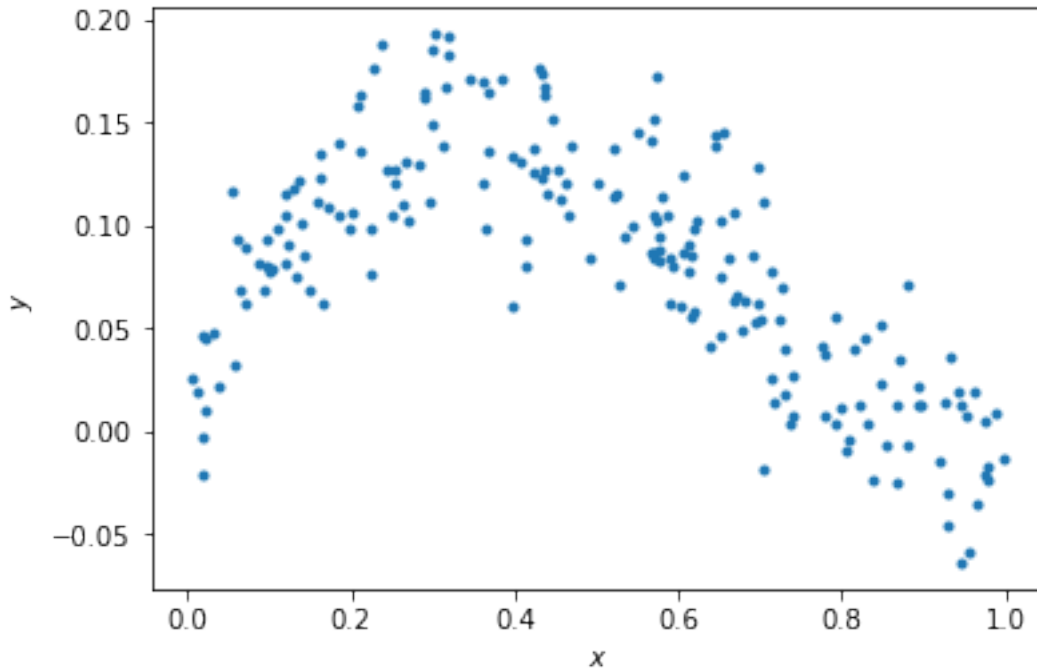
0.1.1 Data generation

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

```
In [59]: 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 [59]: Text(0,0.5,'$y$')
```

0.1.2 QUESTIONS:

Write your answers in the markdown cell below this one:

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

0.1.3 ANSWERS:

- (1) It is a uniform distribution from (including) 0 to (excluding) 1. (200 data points)
- (2) It is a normal (gussain) distribution with mean=0 and standard deviation = 0.33. (200 samples)

0.1.4 Fitting data to the model (5 points)

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

```
In [60]: # 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.zeros(2) # please modify this line
theta = np.linalg.inv((xhat).dot(xhat.T)).dot(xhat.dot(y))
print(theta)
print(theta.shape)

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

```

```

[-0.10599633  0.13315817]
(2,)

```

```

In [33]: # Plot the data and your model fit.
f = plt.figure()
ax = f.gca()
ax.plot(x, y, '. ')
ax.set_xlabel('$x$')
ax.set_ylabel('$y$')

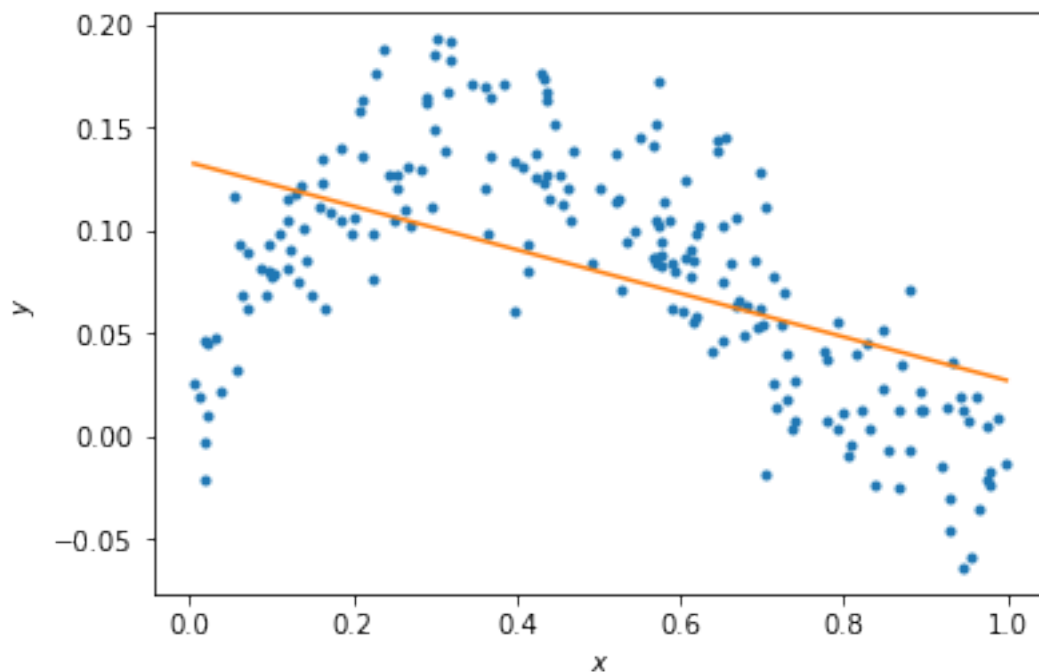
# Plot the regression line
xs = np.linspace(min(x), max(x), 50)
xs = np.vstack((xs, np.ones_like(xs)))
plt.plot(xs[0,:], theta.dot(xs))

```

```

Out[33]: [<matplotlib.lines.Line2D at 0x10dfa1ac8>]

```



0.1.5 QUESTIONS

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

0.1.6 ANSWERS

- (1) The linear model underfits
- (2) To improve the model, we can add more parameters to theta to better fit the curve. This way we'll have an equation of higher order and increase the dimensions of theta and X correspondingly.

0.1.7 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.

```
In [61]: 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.

         for i in np.arange(N):
             if i == 0:
                 thetas.append(theta)
                 xhats.append(xhat)
             else:
                 xhat = np.vstack((x**(i+1), xhat))
                 xhats.append(xhat)
                 thetas.append(np.linalg.inv((xhats[i]).dot(xhats[i].T)).dot(xhats[i].dot(y)))

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

```

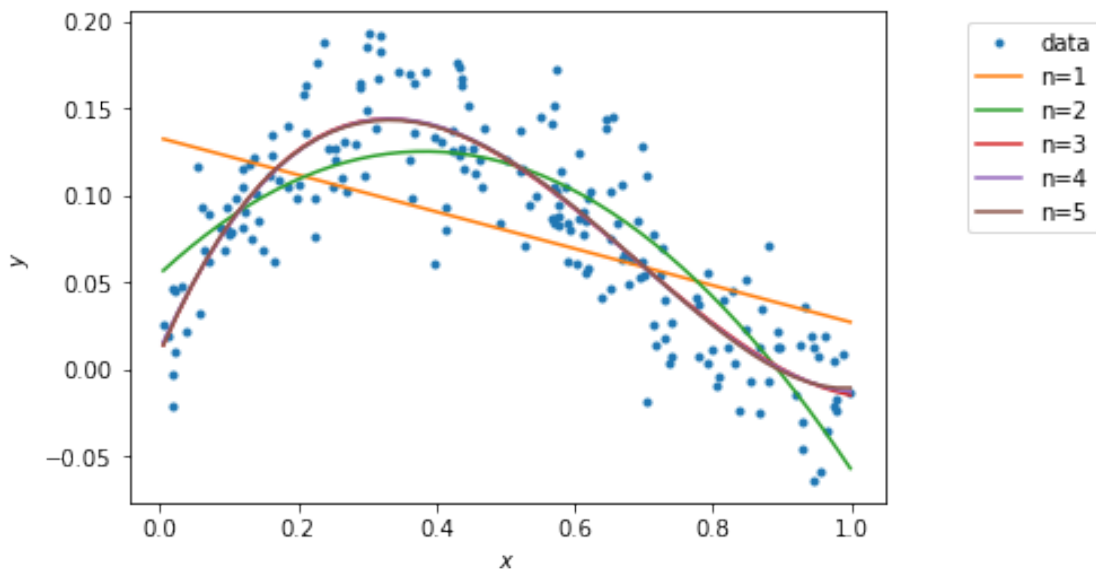
In [62]: # Plot the data
f = plt.figure()
ax = f.gca()
ax.plot(x, y, '.')
ax.set_xlabel('$x$')
ax.set_ylabel('$y$')

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

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

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

```



0.1.8 Calculating the training error (10 points)

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

$$L(\theta) = \frac{1}{2} \sum_j (\hat{y}_j - y_j)^2$$


```

In [73]: 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 np.arange(N):
    training_errors.append( (1/2) * (np.linalg.norm(y - thetas[i].dot(xhats[i]))**2))

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

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

```

Training errors are:

```
[80.861651845505861, 213.19192445058013, 3.1256971083304963, 1.187076519711066, 214.910218105]
```

0.1.9 QUESTIONS

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

0.1.10 ANSWERS

- (1) The polynomial with 5 (4+1) degree has the best error. (The least training error)
- (2) The higher the degree goes the model picks more parameters to fit the curve and according to the theorem it would at least perform as good as the one degree below itself if not better. Therefore it is expected to have lower training error as the degree goes up.

0.1.11 Generating new samples and testing error (5 points)

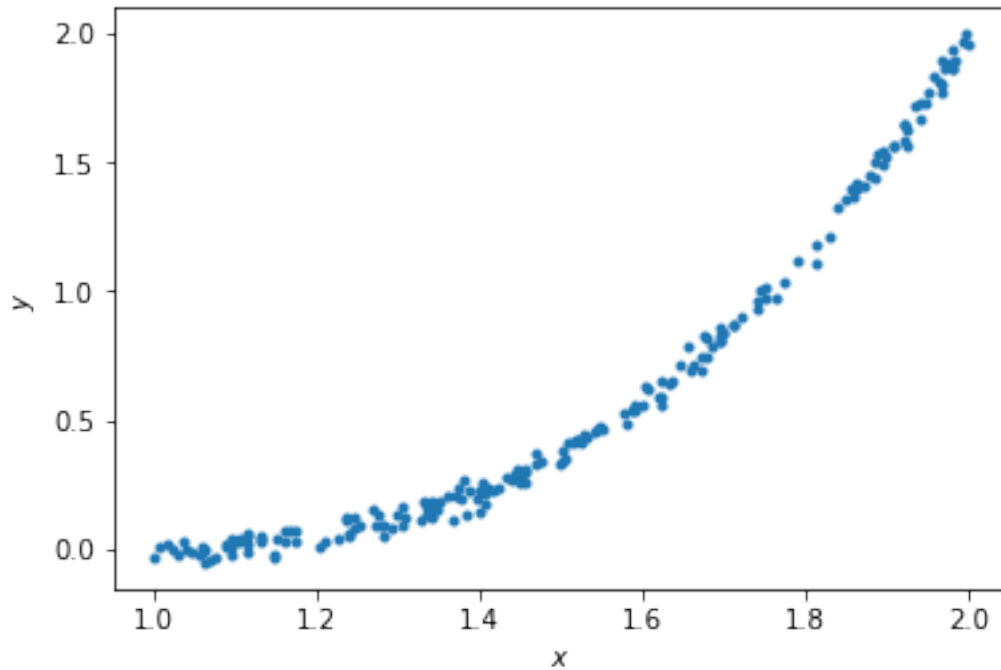
Here, we'll now generate new samples and calculate the testing error of polynomial models of orders 1 to 5.

```

In [66]: 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 [66]: Text(0,0.5,'$y$')
```



```
In [67]: 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 [68]: # 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))
```

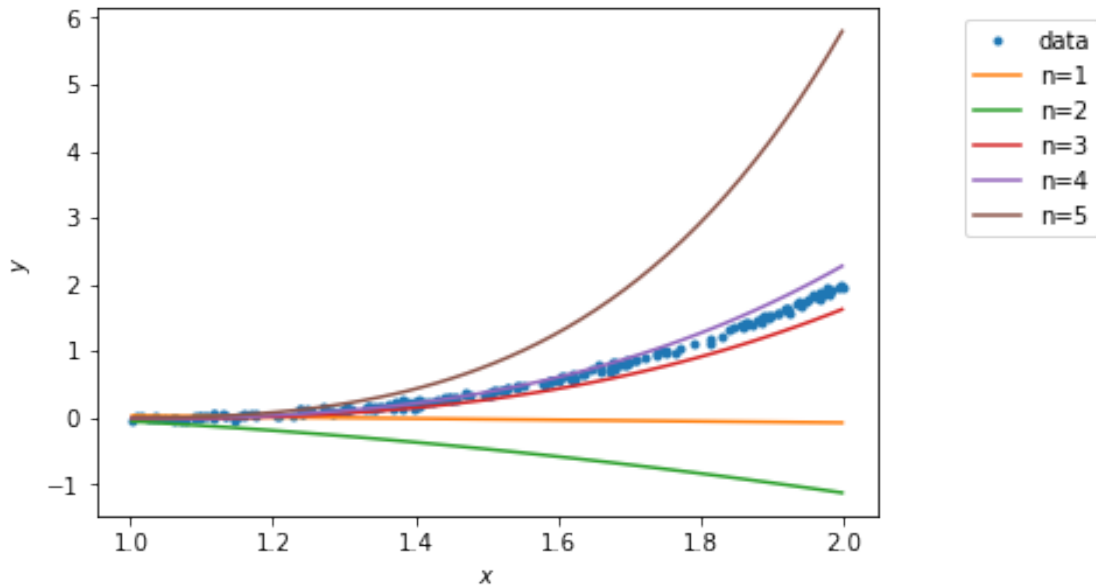
```

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 [75]: 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 np.arange(N):
    testing_errors.append((1/2)* (np.linalg.norm(y - thetas[i].dot(xhats[i]))**2))

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

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

```


Testing errors are:

[80.861651845505861, 213.19192445058013, 3.1256971083304963, 1.187076519711066, 214.910218105]

0.1.12 QUESTIONS

- (1) Which polynomial model has the best testing error?
- (2) Why does the order-5 polynomial model not generalize well?

0.1.13 ANSWERS

- (1) The polynomial with degree 4 does the best (the lowest) testing error
- (2) This is a very good example of overfitting. The training error with $n=5$ tries to cover the maximum number of that training set which does not necessarily generalize to the best error rate for all test cases.