

(1)

1) Linear Algebra Review

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

$$AA^T = I \Rightarrow \begin{bmatrix} a^2+b^2 & ac+bd \\ ac+bd & c^2+d^2 \end{bmatrix} \\ = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$a^2+b^2=1$$

$$ac = -bd$$

$$c^2+d^2=1$$

$$\det A = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix}$$

Eigenvalues & Eigen Vectors

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

$$\det \begin{bmatrix} \frac{1}{\sqrt{2}} - \lambda & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} - \lambda \end{bmatrix} = 0$$

$$+ \left(-\lambda + \frac{1}{\sqrt{2}}\right) \left(\frac{1}{\sqrt{2}} - \lambda\right) - \frac{1}{2} = 0$$

$$\lambda^2 + \frac{\lambda}{\sqrt{2}} - \frac{\lambda}{\sqrt{2}} - \frac{1}{2} - \frac{1}{2} = 0$$

$$\lambda^2 = 1$$

$$\boxed{\lambda = \pm 1}$$

(2)

Eigen Vectors

$$\text{for } \lambda_1 = 1$$

$$(A - \lambda_1 I) x_1 = 0$$

$$(A - I) x_1 = 0$$

$$\begin{bmatrix} \frac{1}{\sqrt{2}} - 1 & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} - 1 \end{bmatrix} x_1 = 0$$

$$\therefore x_1 = \begin{bmatrix} 1 \\ \sqrt{2} - 1 \end{bmatrix}$$

$$\text{for } \lambda_1 = -1$$

$$(A + I) x_2 = 0$$

$$\begin{bmatrix} \frac{1}{\sqrt{2}} + 1 & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} + 1 \end{bmatrix} x_2 = 0$$

$$x_2 = \begin{bmatrix} 1 - \sqrt{2} \\ 1 \end{bmatrix}$$

we notice the following:

$$|\lambda| = 1 \quad \& \quad \text{norm of eigen vector}$$

$|\det A| = 1$. Eigen vectors are ortho.

gonal to each other & can be

normalized.

14

(3)

$$\text{ii) } Av = \lambda v$$

$$\|Av\|^2 = \|\lambda v\|^2 = \|\lambda\|^2 \|v\|^2$$

$$\|Av\|^2 = v^T A^T A v = v^T v = \|v\|^2$$

$$\therefore \|\lambda\|^2 = 1 \Rightarrow \boxed{|\lambda| = 1}$$

iii) if $AA^T = I \Rightarrow A^T A = I$

in this case, $\lambda = \pm 1$.

$$Ax_1 = \lambda_1 x_1$$

$$Ax_2 = \lambda_2 x_2$$

$$Ax_1 = x_1$$

$$Ax_2 = -x_2$$

$$x_1^T = x_1^T A^T$$

$$x_2 = -Ax_2$$

$$x_1^T x_2 = (x_1^T A^T) (-Ax_2)$$

$$= -x_1^T A^T A x_2$$

$$\boxed{x_1^T x_2 = -x_1^T x_2}$$

$$\therefore x_1^T x_2 = 0$$

$\Rightarrow x_1 \perp x_2$ \therefore eigenvectors
corresponding to distinct eigenvalues
are orthogonal.

iv) Any vector x is subject to a rotation or reflection under the transformation Ax .

In this current example, $\det A = 1$

\therefore it is a rotation by 45° in the counter-clockwise direction.

The orientation & the length of x are preserved as is.

by A is a matrix

A can be written as

$$A = UDV^T$$

columns of $U \rightarrow$ left singular vectors

columns of $V \rightarrow$ right singular vectors

elements along diagonal of $D \rightarrow$ singular values.

iv) The left singular vectors of A are the eigenvectors of AA^T .

The Right singular vectors of A are the eigenvectors of A^TA .

ii) The non zero singular values of 'A' (5)
are the square roots of eigen values
of $A^T A$.

The same is true for $A A^T$.

c) True / false

i) False : There are linear operators
with no eigenvalues

ii) False

iii) True

iv) False. If $A = U D V^T \Rightarrow \text{rank}(A) = \text{rank}(D)$
 $= \# \text{ of non zero eigen values}$.

v) True.

Q2) Probability Refresher

a)

$$\text{ii)} \quad P(H50|T) = \frac{P(T|H50) \cdot P(H50)}{P(T)} \quad \text{--- (1)}$$

$$P(T) = P(T|H50) \cdot H50 + P(T|H60) \cdot P(H60)$$

$$= \frac{1}{2} \cdot \frac{1}{2} + \frac{2}{5} \cdot \frac{1}{2}$$

$$P(T) = \frac{9}{20}$$

substituting in (1)

$$\therefore P(H50|T) = \frac{1/4}{9/20} = \frac{5}{9} = 0.5555$$

\therefore The posterior probability that the tail is an H50 win is 0.5555.

$$\text{iii)} \quad P(H50|T_H H H) = \frac{P(T_H H H | H50) \cdot P(H50)}{P(T_H H H)}$$

$$P(T_H H H) = P(T_H H H | H50) P(H50) + P(T_H H H | H60) P(H60)$$

$$= \left(\frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \right) \left(\frac{1}{2} \right) + \left(\frac{2}{5} \times \frac{3}{5} \times \frac{3}{5} \times \frac{3}{5} \right) \left(\frac{1}{2} \right)$$

$$= \frac{1}{16} + \frac{9}{125}$$

$$\therefore P(H50 | THHH) = \frac{1/16}{1/16 + \frac{54}{1250}} \quad (7)$$

$$\boxed{P(H50 | THHH) = \cancel{0.4646} \quad 0.4197}$$

iii)

$$P(H50 | 9 \text{ Heads, } 1 \text{ Tail}) = \frac{P(9 \text{ heads, } 1 \text{ Tail} | H50) \cdot P(H50)}{P(9 \text{ heads, } 1 \text{ tail})}$$

$$\Rightarrow P(9 \text{ heads, } 1 \text{ tail}) = P(9 \text{ heads } 1 \text{ Tail} | H50) P(H50) \\ + P(9 \text{ heads } 1 \text{ Tail} | H55) P(H55) \\ + P(9 \text{ heads } 1 \text{ Tail} | H60) P(H60)$$

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

$$= 2.36 \times 10^{-3}$$

$$\therefore P(H50 | 9 \text{ heads } 1 \text{ tail}) = \frac{3.25 \times 10^{-4}}{2.36 \times 10^{-3}}$$

$$\boxed{P(H50 | 9 \text{ heads } 1 \text{ tail}) = \underline{\underline{0.1379}}}$$

$$P(H55|9HIT) = \frac{P(9HIT|H55) P(H55)}{P(9HIT)}$$

$$P(9HIT) = P(9HIT|H50) P(H50) + P(9HIT|H55) P(H55) + P(9HIT|H60) P(H60)$$

$$= 2.36 \times 10^{-3}$$

$$P(H55|9HIT) = \frac{6.90 \times 10^{-4}}{2.36 \times 10^{-3}}$$

$$P(H55|9HIT) = 0.2927$$

similarly,

$$P(H60|9HIT) = \frac{P(9HIT|H60) P(H60)}{P(9HIT)}$$

$$= \frac{1.3436 \times 10^{-3}}{2.36 \times 10^{-3}}$$

$$P(H60|9HIT) = 0.5694$$

2/6/

9

	Positive	Negative
Pregnant	0.99	0.01
Not Pregnant	0.10	0.90

$$P(\text{Pregnant} | \text{Positive}) = \frac{P(\text{Positive} | \text{Pregnant}) P(\text{Pregnant})}{P(\text{Positive})}$$

$$\begin{aligned} P(\text{Positive}) &= P(\text{Positive} | \text{Preg.}) * P(\text{Preg.}) + \\ &\quad P(\text{Positive} | \text{Not Preg.}) * P(\text{Not Preg.}) \\ &= (0.99 \times 0.01) + (0.10 \times 0.99) \\ &= 0.1089 \end{aligned}$$

$$\begin{aligned} \therefore P(\text{Preg.} | \text{Positive}) &= \frac{0.99 \times 0.01}{0.1089} \\ &= \frac{99 \times 10^{-3}}{0.1089} \end{aligned}$$

$$\boxed{\therefore P(\text{Preg} | \text{Positive}) = 0.09091}$$

Intuitively, the answer says that the pregnancy test results are not reliable since a woman is pregnant only 9% of the times when the test results are positive.

2/cp

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

$$E(X) = \sum_n x p(x)$$

$$\begin{aligned} \therefore E(AX+b) &= \sum_n (AX+b) p(x) \\ &= \sum_n AX p(x) + b \sum_n p(x) \\ &= A \sum_n x p(x) + b(1) \end{aligned}$$

$$\boxed{\therefore E(AX+b) = A E(X) + b}$$

This proves linearity of expectation with respect to one variable.

2/dp

$$\text{Cov}(X) = E((X - E(X))(X - E(X))^T)$$

$$\text{Cov}(AX+b) = E[(AX+b - E(AX+b))(AX+b - E(AX+b))^T]$$

$$= E[(AX+b - AE(X) - b)(AX+b - AE(X) - b)^T]$$

{ from part c }

$$= E[A(X - E(X))A(X - E(X))^T]$$

$$= E \left[A (x - E(x)) (x - E(x))^T A^T \right]$$

$$= A E \left[(x - E(x)) (x - E(x))^T \right] A^T$$

$$\boxed{\therefore \text{Cov}(Ax+b) = A \text{Cov}(x) A^T}$$

Q/ Problem 3: Multivariate Derivative

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

Let,

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

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

ay $\nabla_x x^T A y$

$$\frac{\partial f}{\partial x_i} = \sum_{j=1}^m a_{ij} y_j = A y$$

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

b)

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

$$\nabla_{\mathbf{y}} (\mathbf{x}^T \mathbf{A} \mathbf{y}) = \frac{\partial f}{\partial y_j} = \sum_{i=1}^n a_{ij} x_i = \mathbf{A}^T \mathbf{x}$$

$$\boxed{\therefore \nabla_{\mathbf{y}} (\mathbf{x}^T \mathbf{A} \mathbf{y}) = \mathbf{A}^T \mathbf{x}}$$

c)

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

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

$$= \frac{\partial f}{\partial A} = \begin{bmatrix} \frac{\partial f}{\partial a_{11}} & \dots & \frac{\partial f}{\partial a_{1m}} \\ \frac{\partial f}{\partial a_{n1}} & \dots & \frac{\partial f}{\partial a_{nm}} \end{bmatrix}$$

$$= \begin{bmatrix} x_1 y_1 & \dots & x_1 y_m \\ x_n y_1 & \dots & x_n y_m \end{bmatrix}$$

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

$$\boxed{\therefore \nabla_{\mathbf{A}} (\mathbf{x}^T \mathbf{A} \mathbf{y}) = \mathbf{x} \mathbf{y}^T}$$

$$Q \quad f = x^T A x + b^T x$$

13

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

$$\nabla_x f = \nabla_x (x^T A x) + \nabla_x (b^T x) \quad \text{--- (1)}$$

$$\text{Consider } y = x^T A x$$

$$y = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j$$

$$\frac{\partial y}{\partial x_i} = 2a_{ii} x_i + \sum_{j \neq i}^n a_{ij} x_j + \sum_{i \neq j}^n a_{ji} x_j$$

$$= \sum_{j=1}^n a_{ij} x_j + \sum_{i=1}^n a_{ji} x_j$$

$$\frac{\partial y}{\partial x_i} = (A x)_i + (A^T x)_i$$

$$\therefore \nabla_x (x^T A x) = (A + A^T) x \quad \text{--- (2)}$$

$$\nabla_x (b^T x) = b^T \quad \text{--- (3)}$$

substituting (2), (3) in (1)

$$\boxed{\nabla_x f = (A + A^T) x + b^T}$$

3) ex $A = n \times m$
 $B = m \times n$

Let $A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1m} \\ \vdots & & & \\ a_{n1} & \dots & \dots & a_{nm} \end{bmatrix}$ $\{ B = \begin{bmatrix} b_{11} & \dots & b_{1n} \\ \vdots & & \vdots \\ b_{m1} & \dots & b_{mn} \end{bmatrix}$

$$AB = \begin{bmatrix} (a_{11}b_{11} + a_{12}b_{21} + \dots + a_{1m}b_{m1}) & \dots & \dots \\ \vdots & & \vdots \\ (a_{21}b_{12} + a_{22}b_{22} + \dots + a_{2m}b_{m2}) & \dots & \dots \\ \vdots & & \vdots \\ \vdots & & \vdots \\ (a_{n1}b_{1n} + a_{n2}b_{2n} + \dots + a_{nm}b_{mn}) \end{bmatrix}$$

$$\therefore \text{tr}(AB)_i = \sum_{i=1}^n \sum_{j=1}^m a_{ij} b_{ji}$$

$$\text{let } y = \sum_{i=1}^n \sum_{j=1}^m a_{ij} b_{ji}$$

$$\frac{\partial y}{\partial A} = \begin{bmatrix} \frac{\partial y}{\partial a_{11}} & \dots & \frac{\partial y}{\partial a_{1m}} \\ \vdots & & \vdots \\ \frac{\partial y}{\partial a_{n1}} & \dots & \frac{\partial y}{\partial a_{nm}} \end{bmatrix}$$

$$= \begin{bmatrix} b_{11} & \dots & b_{1m} \\ \vdots & & \vdots \\ b_{n1} & \dots & b_{nm} \end{bmatrix} = B^T$$

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

4) Deriving Least Squares

(15)

In the given problem, we have to minimize the error.

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

∴ our cost funcⁿ is

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

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

$$L = \frac{1}{2} (Y - XW^T)^T (Y - XW^T)$$

$$L = \frac{1}{2} (Y^T Y - \underbrace{YXW^T - WX^T Y}_{}) + WX^T X W^T)$$

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

to find optimum w ,

$\nabla_w L = 0$. using the given hints, we have,

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

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

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

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

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

$$\nabla_x f = \nabla_x (x^T A x) + \nabla_x (b^T x) \quad \text{--- (1)}$$

consider $y = x^T A x$

$$y = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j$$

$$\frac{\partial y}{\partial x_i} = 2a_{ii} x_i + \sum_{j \neq i}^n a_{ij} x_j + \sum_{i \neq j}^n a_{ji} x_i$$

$$= \sum_{j=1}^n a_{ij} x_j + \sum_{i=1}^n a_{ji} x_i$$

$$\frac{\partial y}{\partial x_i} = (A x)_i + (A^T x)_i$$

$$\therefore \nabla_x (x^T A x) = (A + A^T) x \quad \text{--- (2)}$$

$$\nabla_x (b^T x) = b \quad \text{--- (3)}$$

substituting (2), (3) in (1)

$$\boxed{\nabla_x f = (A + A^T) x + b}$$

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

```
In [134]: import numpy as np
import matplotlib.pyplot as plt

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

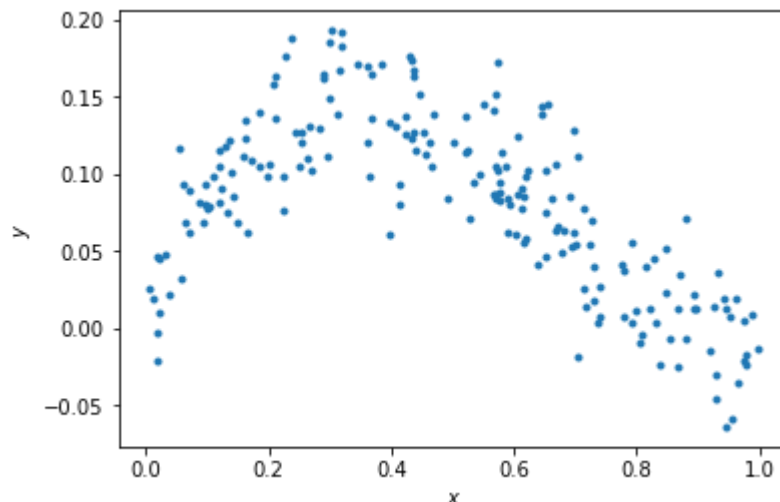
Data generation

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 [135]: 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[135]: <matplotlib.text.Text at 0x1126bf850>



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) The generating distribution of x is Uniform
- (2) The distribution of the additive noise ϵ is Normal

Fitting data to the model (5 points)

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

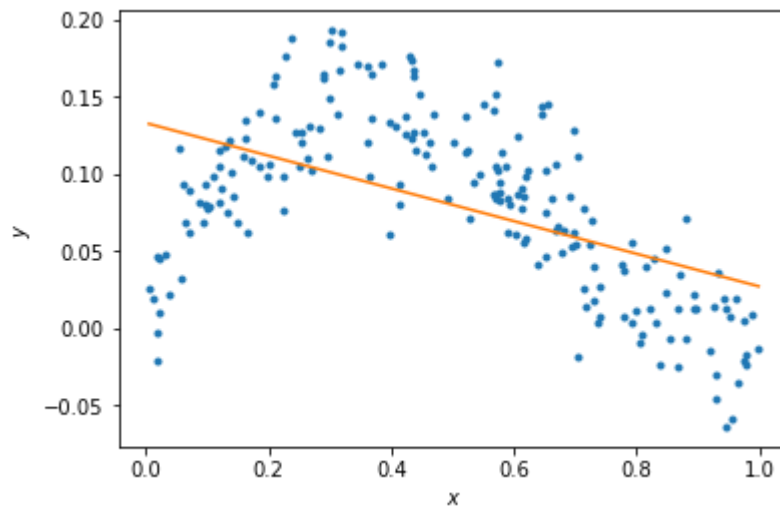
```
In [136]: # xhat = (x, 1)
          xhat = np.vstack((x, np.ones_like(x)))

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

```
In [137]: # 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[137]: [



QUESTIONS

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

ANSWERS

- (1) The model is underfitting.
- (2) We can improve the fitting by using an equation of higher order or degree i.e., to increase the dimensions of theta and X correspondingly.

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 [138]: 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.

          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].d
ot(y)))

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

```

```

In [139]: # Plot the data
f = plt.figure()
ax = f.gca()
ax.plot(x, y, '.')
ax.set_xlabel('$x$')
ax.set_ylabel('$y$')

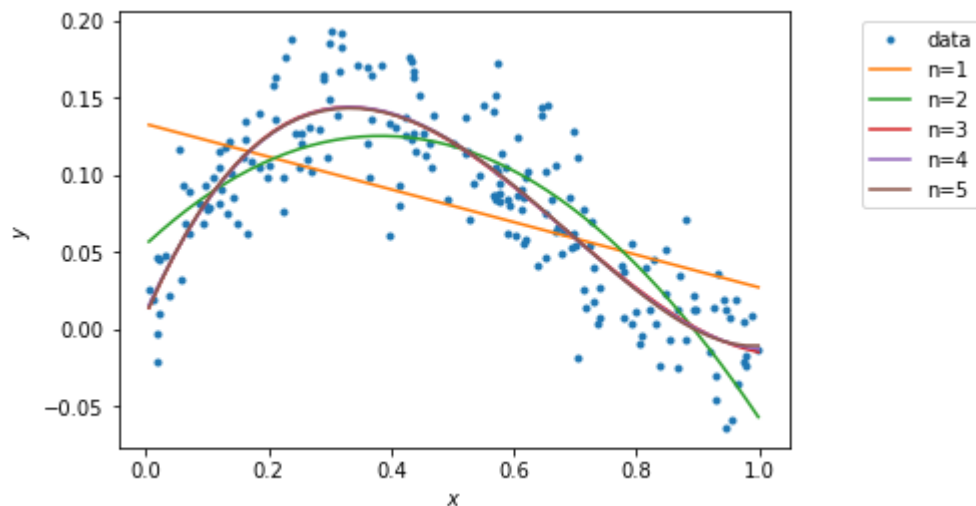
# Plot the regression lines
plot_xs = []
for i in np.arange(N):
    if i == 0:
        plot_x = np.vstack((np.linspace(min(x), max(x), 50), np.ones(50)))
    else:
        plot_x = np.vstack((plot_x[-2]**(i+1), plot_x))

    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.

```
In [141]: 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/float(len(y))) * np.linalg.norm(y - thetas[i].dot(xhats[i])))

pass

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

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

('Training errors are: \n', [0.0034496094622164849, 0.0023371908575540567, 0.0020210892856459082, 0.0020205635025033215, 0.0020200840571032307])
```

QUESTIONS

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

ANSWERS

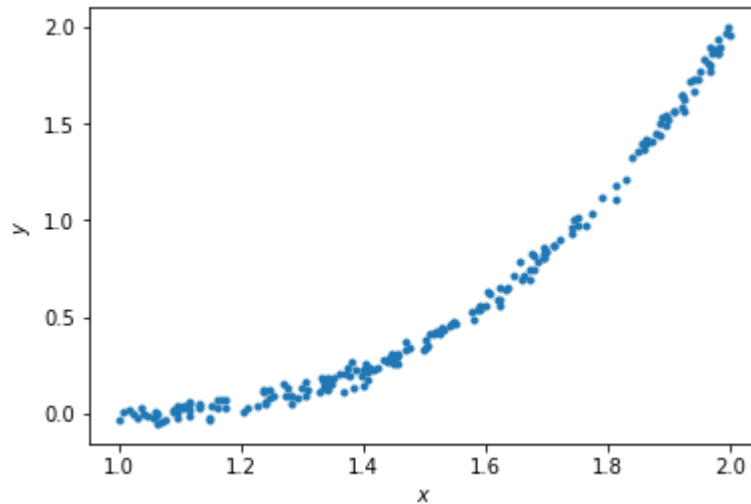
- (1) The polynomial of degree 5 has the best training error (i.e., the least training error.)
- (2) This was expected because a polynomial of higher degree covers all the points in the training data thus decreasing the training error to its least value.

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 [144]: 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[144]: <matplotlib.text.Text at 0x112a28e50>



```
In [145]: 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 [146]: # Plot the data
f = plt.figure()
ax = f.gca()
ax.plot(x, y, '.')
ax.set_xlabel('$x$')
ax.set_ylabel('$y$')

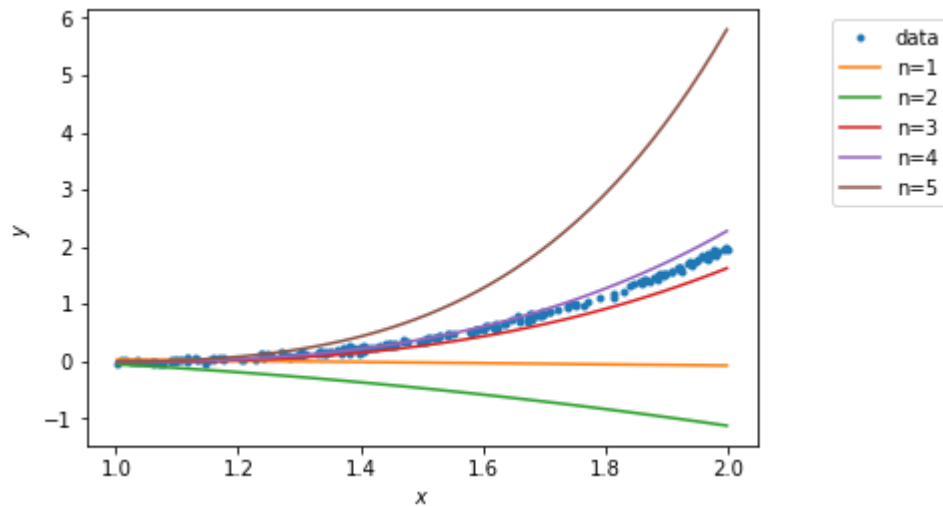
# Plot the regression lines
plot_xs = []
for i in np.arange(N):
    if i == 0:
        plot_x = np.vstack((np.linspace(min(x), max(x), 50), np.ones(50)))
    else:
        plot_x = np.vstack((plot_x[-2]**(i+1), plot_x))

    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 [147]: 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/float(len(y))) * np.linalg.norm(y - thetas[i].dot
(xhats[i])))

pass

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

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

('Testing errors are: \n', [0.063585238792311649, 0.10324532058417471, 0.0125
01394138930177, 0.0077041434297982813, 0.10366055616645076])
```

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 has the best Testing Error
- (2) The polynomial of degree 5 is over fitting the training data. Although the training error is the least, the 5th degree polynomial does not generalize well as it tries to cover maximum number points in the training data.