

$$1) \quad 2) \quad 3) \quad A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$A \cdot A^T = I \quad \checkmark$$

Eigenvalue

$$\det \begin{bmatrix} -\lambda & 1 \\ 1 & -\lambda \end{bmatrix} = \lambda^2 - 1 = (\lambda + 1)(\lambda - 1) = 0$$

$$\lambda = -1, 1$$

Eigenvector:

$$\lambda = -1 \quad \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\vec{x}_1 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$\lambda = 1 \quad \begin{bmatrix} -1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\vec{x}_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

I notice that \vec{x}_1 and \vec{x}_2 are orthogonal to each other and that $\lambda_1 = -\lambda_2$

ii) $Av = \lambda v$ by def of eigenvector, eigenvalue

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

$$(Av)^T(Av) = \lambda^2 |v|^2$$

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

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

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

$$| = \lambda^2 \Rightarrow | = |\lambda|$$

iii) Since A is orthogonal, we have $A^T = A^{-1}$
 and thus $A^{-1}A = I = AA^{-1}$
 $\Rightarrow A^TA = AA^T$

Thus A is normal

Since A is normal, $A = UDV^*$

where D is diagonal and the columns
 are eigenvectors

Thus the eigenvectors are orthogonal

iv) Vector x will be rotated or flipped

b) i) Left singular vectors of A are orthonormal
 eigenvectors of AA^T

Right singular vectors of A are orthonormal
 eigenvectors of A^TA

ii) Nonzero singular values of A are the
 square roots of non-zero eigenvalues of
 both AA^T and A^TA

c) i) False, consider I_2

ii) False, not true in all cases - $A = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$

iii) True

iv) True

v) True

$$\begin{aligned} 2) 2) i) P(\text{coin} = \text{HSD} \mid \text{lands} = \text{tails}) &= \frac{P(l = \text{HSD} \wedge l = \text{tails})}{P(l = \text{tails})} \\ &= P(l = \text{tails} \mid \text{C} = \text{HSD})P(\text{C} = \text{HSD}) \end{aligned}$$

$$\begin{aligned}
 & \overline{P(l=t_{21}l_5 | C=H_{SO})P(C=H_{SO}) + P(l=t_{21}l_5 | C=H_{BO})P(C=H_{BO})} \\
 & = \frac{(0.5)(0.5)}{(0.5)(0.5) + (0.4)(0.5)} \\
 & = 0.555
 \end{aligned}$$

$$\begin{aligned}
 \text{ii)} \quad & P(C=H_{SO} | l_1=T \wedge l_2=H \wedge l_3=H \wedge l_4=H) \\
 & = \frac{P(C=H_{SO} \wedge l_1=T \wedge l_2=H \wedge l_3=H \wedge l_4=H)}{P(l_1=T \wedge l_2=H \wedge l_3=H \wedge l_4=H)} \\
 & \quad P(l_1=T | C=H_{SO})P(l_2=H | C=H_{SO}) \\
 & \quad P(l_3=H | C=H_{SO})P(l_4=H | C=H_{SO}) \\
 & \quad \frac{P(l=H_{SO})}{P(l_1=T | C=H_{SO})P(l_2=H | C=H_{SO})} \\
 & \quad P(l_3=H | C=H_{SO})P(l_4=H | C=H_{SO}) \\
 & \quad P(l=H_{SO}) + \\
 & \quad P(l=H_{SO})P(l_2=H | C=H_{BO}) \\
 & \quad P(l_3=H | C=H_{BO})P(l_4=H | C=H_{BO}) \\
 & \quad P(l=H_{BO}) \\
 & = \frac{(0.5)^4 \cdot 0.5}{(0.5)^4 \cdot 0.5 + (0.4)^4 (0.6)^3 \cdot 0.5} \\
 & = 0.419
 \end{aligned}$$

$$\begin{aligned}
 \text{iii)} \quad & P(l_S=9H/10 | C=H_{SO}) = 0.5^4 \\
 & P(l_S=9H/10 | C=H_{SS}) = 0.55^4 \cdot 0.45 \\
 & P(l_S=9H/10 | C=H_{BO}) = 0.6^4 \cdot 0.40 \\
 & P(C=H_{SO} | l_S=9H/10) = \frac{0.5^4 \cdot 1/3}{0.5^4 \cdot 1/3 + 0.55^4 \cdot 0.45 \cdot 1/3 + 0.6^4 \cdot 0.40 \cdot 1/3} \\
 & = 0.137
 \end{aligned}$$

$$P(C=HSS \mid LS=9H/10) = \frac{0.5^{10} \cdot 1/3}{0.5^{10} \cdot 1/3 + 0.55^9 \cdot 0.45 \cdot 1/3 + 0.60^9 \cdot 0.40 \cdot 1/3}$$

$$= 0.292$$

$$P(C=HSD \mid LS=9H/10) = \frac{0.5^{10} \cdot 1/3}{0.5^{10} \cdot 1/3 + 0.55^9 \cdot 0.45 \cdot 1/3 + 0.60^9 \cdot 0.40 \cdot 1/3}$$

$$= 0.569$$

b) $P(P=T \mid T=T) = \frac{P(P=T \cap T=T)}{P(T=T)}$

$$= \frac{P(T=T \mid P=T)P(P=T)}{P(T=T \mid P=T)P(P=T) + P(T=T \mid P=F)P(P=F)}$$

$$= \frac{(0.99)(0.01)}{(0.99)(0.01) + (0.10)(0.99)}$$

$$= 0.0909$$

Not most likely b/c false positive is relatively high - if false pos. were lower, could be more confident

c) $E[AX+b] = A E[X] + b$ b/c linearity

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

$$= E[(Ax+b - AE[x]-b) \cdot (Ax+b - AE[x]-b)^T]$$

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

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

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

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

$$= x_1 (a_{11}y_1 + \dots + a_{1m}y_m) + \\ x_2 (a_{21}y_1 + \dots + a_{2m}y_m) + \\ \dots + \\ x_n (a_{n1}y_1 + \dots + a_{nm}y_m)$$

$$\nabla_x x^T A y = \begin{bmatrix} \frac{\partial x^T A y}{\partial x_1} \\ \vdots \\ \frac{\partial x^T A y}{\partial x_n} \end{bmatrix} = A y$$

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

$$= y_1 (x_1 a_{11} + \dots + x_n a_{n1}) + \\ y_2 (x_1 a_{12} + \dots + x_n a_{n2}) + \\ \dots + \\ y_m (x_1 a_{1m} + \dots + x_n a_{nm})$$

$$\nabla_y x^T A y = \begin{bmatrix} \frac{\partial x^T A y}{\partial y_1} \\ \vdots \\ \frac{\partial x^T A y}{\partial y_m} \end{bmatrix} = x^T A$$

c) $\nabla_A x^T A y = \begin{bmatrix} \frac{\partial x^T A y}{\partial a_{11}} & \dots & \frac{\partial x^T A y}{\partial a_{1m}} \\ \vdots & \ddots & \vdots \\ \frac{\partial x^T A y}{\partial a_{n1}} & \dots & \frac{\partial x^T A y}{\partial a_{nm}} \end{bmatrix}$

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

$$= x y^T$$

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

$$\begin{aligned}
&= [x_1, \dots, x_n] \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \\
&\quad + [b_1, \dots, b_n] \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \\
&= x_1(a_{11}x_1 + \dots + a_{1n}x_n) + x_2(a_{21}x_1 + \dots + a_{2n}x_n) \\
&\quad + \dots + x_n(a_{n1}x_1 + \dots + a_{nn}x_n) \\
&\quad + b_1x_1 + \dots + b_nx_n
\end{aligned}$$

$$\begin{aligned}
\nabla_x f &= \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n + \dots + b_1 \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n + \dots + b_n \end{bmatrix} \\
&= A^T x + b
\end{aligned}$$

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

$$\begin{aligned}
&= \text{tr} \left(\begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} b_{11} & \dots & b_{1n} \\ \vdots & \ddots & \vdots \\ b_{n1} & \dots & b_{nn} \end{bmatrix} \right) \\
&= \text{tr} \left(\begin{bmatrix} a_{11}b_{11} + \dots + a_{1n}b_{1n} & \dots & \dots \\ \vdots & \ddots & \vdots \\ a_{n1}b_{n1} + \dots + a_{nn}b_{nn} & \dots & \dots \end{bmatrix} \right)
\end{aligned}$$

$$\begin{aligned}
\nabla_A f &= \begin{bmatrix} \frac{\partial f}{\partial a_{11}} & \dots & \frac{\partial f}{\partial a_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial a_{n1}} & \dots & \frac{\partial f}{\partial a_{nn}} \end{bmatrix} \\
&= \begin{bmatrix} b_{11} & \dots & b_{1n} \\ \vdots & \ddots & \vdots \\ b_{n1} & \dots & b_{nn} \end{bmatrix} = B^T
\end{aligned}$$

$$\begin{bmatrix} b_{11} & \dots & b_{mn} \end{bmatrix}$$

q) Derive optimal ω for $\min_{\omega} \sum_{i=1}^n \|y^{(i)} - \omega x^{(i)}\|^2$

$$= \sum_{i=1}^n (y^{(i)} - \omega x^{(i)}) (y^{(i)} - \omega x^{(i)})^\top$$

$$= \text{tr} \left((\Upsilon^\top - X^\top \omega^\top) (\Upsilon - \omega X) \right)$$

where $\Upsilon^\top = \begin{bmatrix} - & y^{(1)\top} & - \\ & \vdots & \\ - & y^{(n)\top} & - \end{bmatrix}$

$$X^\top = \begin{bmatrix} - & x^{(1)\top} & - \\ & \vdots & \\ - & x^{(n)\top} & - \end{bmatrix}$$

$$y^{(i)} \in \mathbb{R}^{m \times 1}$$

$$\omega \in \mathbb{R}^{m \times c}$$

$$x^{(i)} \in \mathbb{R}^{c \times 1}$$

$$\Upsilon^\top \in \mathbb{R}^{n \times m}$$

$$\omega^\top \in \mathbb{R}^{c \times n}$$

$$X^\top \in \mathbb{R}^{n \times c}$$

$$\frac{\partial f}{\partial \omega} = 0 = \frac{\partial}{\partial \omega} \sum_{i=1}^n \text{tr} \left((\Upsilon^\top - X^\top \omega^\top) (\Upsilon - \omega X) \right)$$

$$= \sum_{i=1}^n \frac{\partial}{\partial \omega} \text{tr} \left(\Upsilon^\top \Upsilon - \Upsilon^\top \omega X - X^\top \omega^\top \Upsilon + X^\top \omega^\top \omega X \right)$$

$$= \sum_{i=1}^n \frac{\partial}{\partial \omega} \text{tr} (\Upsilon^\top \Upsilon) - \text{tr} (\Upsilon^\top \omega X) - \text{tr} (X^\top \omega^\top \Upsilon) + \text{tr} (X^\top \omega^\top \omega X)$$

$$= \sum_{i=1}^n \frac{\partial}{\partial \omega} - \text{tr} (\Upsilon^\top \omega X) - \text{tr} (X^\top \omega^\top \Upsilon) + \text{tr} (X^\top \omega^\top \omega X)$$

$$= \sum_{i=1}^n (-2 \Upsilon^\top X + 2 \omega X X^\top)$$

$$= \omega X X^\top - \Upsilon^\top X$$

$$\omega X X^\top (X X^\top)^{-1} = \Upsilon^\top X (X X^\top)^{-1}$$

$$\omega = \Upsilon^\top X (X X^\top)^{-1}$$

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 239AS, Winter Quarter 2019, Prof. J.C. Kao, TAs M. Kleinman and A. Wickstrom and K. Liang and W. Chuang

```
In [26]: import numpy as np
import matplotlib.pyplot as plt

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

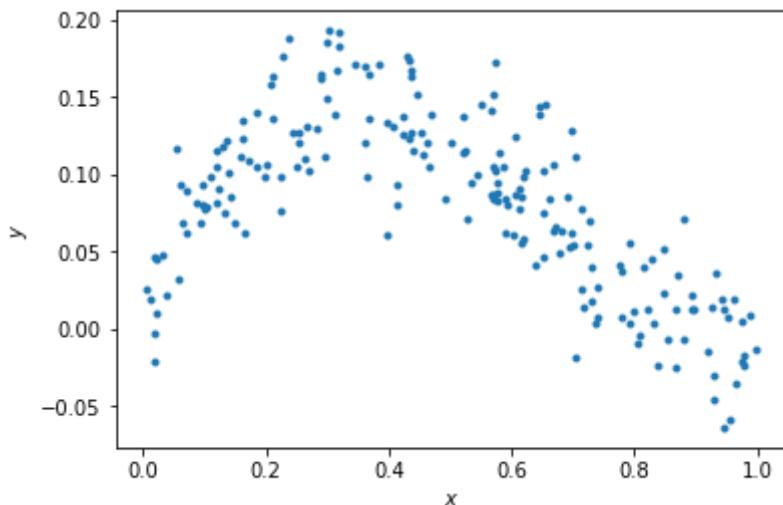
Data generation

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 [27]: 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[27]: 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 ϵ ?

ANSWERS:

- (1) x is from a uniform distribution between 0 and 1
- (2) ϵ is from a normal distribution with mean 0 and stddev 0.03

Fitting data to the model (5 points)

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

```
In [28]: # 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

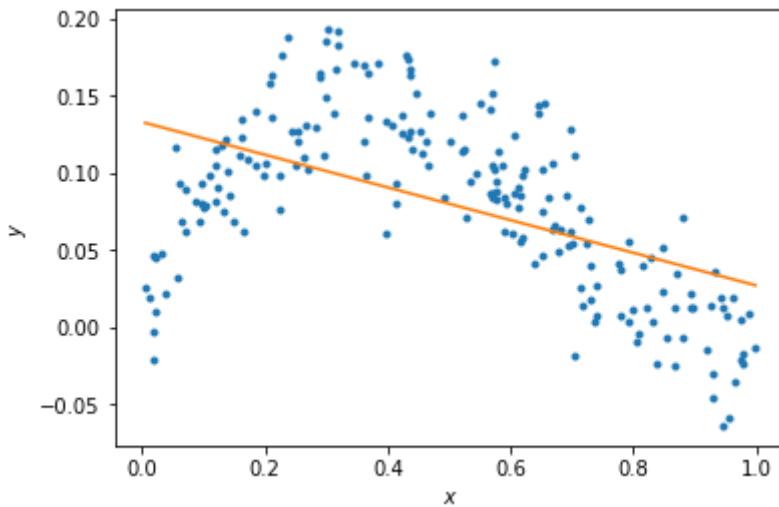
theta = np.linalg.inv((xhat).dot(xhat.T)).dot(y)

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

```
In [29]: # 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[29]: [`<matplotlib.lines.Line2D at 0x113054710>`]



QUESTIONS

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

ANSWERS

- (1) Underfits the data
- (2) Increase complexity of the model (add x^2 and x^3 terms)

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 [30]:

```
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
# i.e., thetas[0] is equivalent to theta above.
# i.e., thetas[1] should be a length 3 np.array with the coefficients of
# ... etc.
for i in np.arange(N):
    xhats.append(xhat if i == 0 else np.vstack((x**(i+1), xhats[i-1])))
    thetas.append(np.linalg.inv((xhats[i]).dot(xhats[i].T)).dot(xhats[i]).dc)

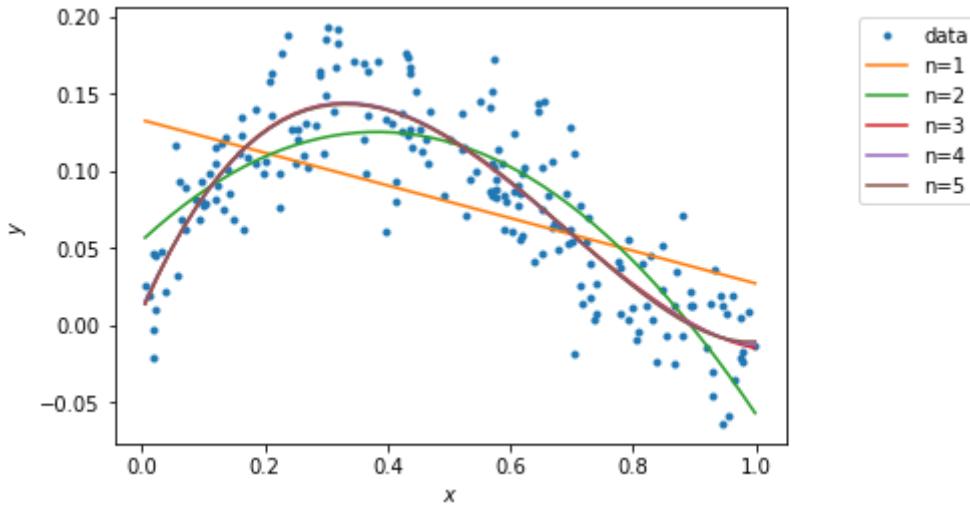
# ===== #
# END YOUR CODE HERE #
# ===== #
```

```
In [31]: # 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 [32]: 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
for i in np.arange(N):
    training_errors.append(0.5*np.sum(np.square(y - (xhats[i].T).dot(thetas)

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

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

```
Training errors are:
[0.2379961088362701, 0.10924922209268531, 0.08169603801105374, 0.0816535
3735296982, 0.08161479195525295]
```

QUESTIONS

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

ANSWERS

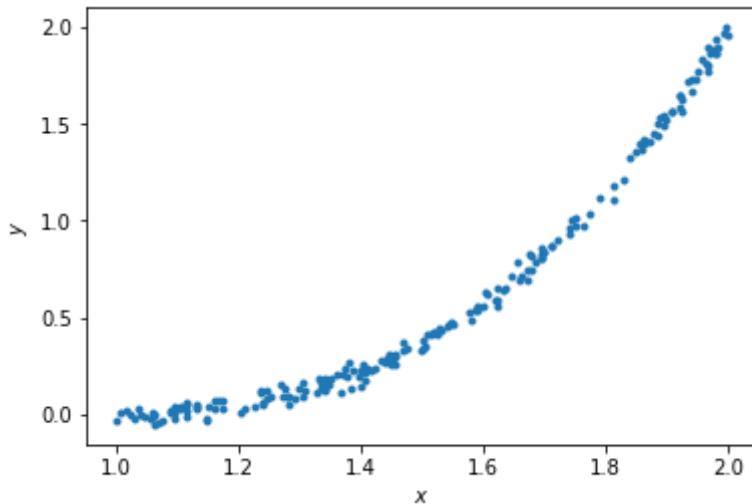
- (1) Polynomial with degree 5
- (2) There are more than 5 datapoints, and so training error is strictly decreasing as polynomial degree increases (because polynomial degree 5 encapsulates polynomial degree 4, and so on, thus, it can use the x^5 to decrease training error even further)

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 [33]: 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[33]: Text(0, 0.5, '$y$')
```



```
In [34]: 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))

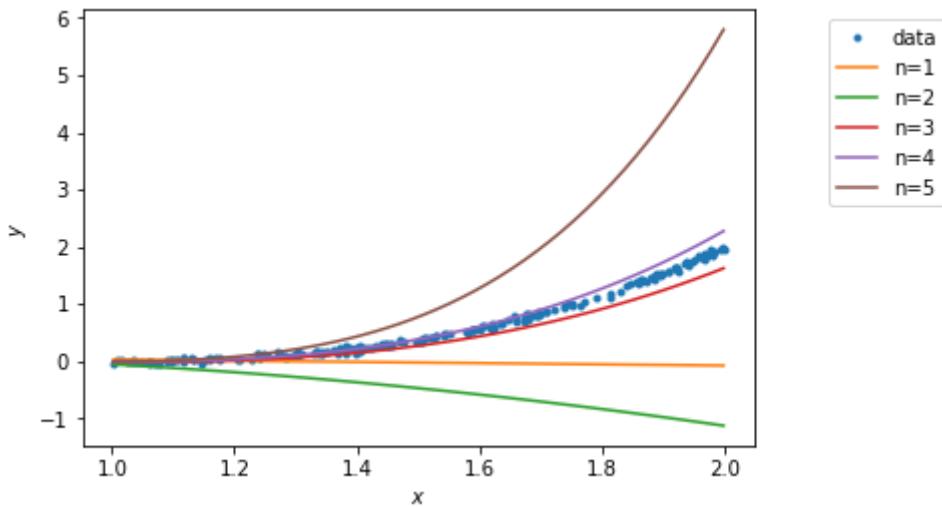
    xhats.append(xhat)
```

```
In [35]: # 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 [36]: 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
for i in np.arange(N):
    testing_errors.append(0.5*np.sum(np.square(y - (xhats[i].T).dot(thetas[0:i+1]))))

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

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

```
Testing errors are:
[80.86165184550586, 213.19192445057894, 3.1256971082763925, 1.1870765189
474703, 214.91021817652626]
```

QUESTIONS

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

ANSWERS

- (1) Polynomial degree 4
- (2) Polynomial degree 5 had overfit the training data which came from an X uniformly distributed over 0 and 1, whereas the testing data has an X uniformly distributed over 1 and 2. The x^5 term then begins to cause very high error, because the testing data seems unlike the training data, despite the underlying function being identical.