

①

$$AA^T = I$$

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

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

$$AA^T = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} a & c \\ b & d \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

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

$$ac + bd = 0$$

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

$$a = \cos \theta$$

$$b = \sin \theta$$

$$c = -\sin \theta$$

$$d = \cos \theta$$

$$bd = -ac$$

only need 1 negative term

$$A = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}$$

$$Av = \lambda v$$

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

$$\begin{bmatrix} \cos \theta - \lambda & \sin \theta \\ \sin \theta & -\cos \theta - \lambda \end{bmatrix} = A - \lambda I$$

$$|A - \lambda I| = (\cos \theta - \lambda)(-\cos \theta - \lambda) - \sin^2 \theta$$

$$0 = -\cos^2 \theta + \lambda^2 + \lambda \cos \theta - \lambda \cos \theta - \sin^2 \theta$$

$$\lambda^2 = \sin^2 \theta + \cos^2 \theta$$

$$\lambda = \pm 1$$

$$\lambda_1 = 1$$

$$\lambda_2 = -1$$

$$\lambda_1: \begin{bmatrix} \cos \theta - 1 & \sin \theta \\ \sin \theta & -\cos \theta - 1 \end{bmatrix} \begin{bmatrix} v_{11} \\ v_{12} \end{bmatrix} = 0$$

$$v_{12} = 1$$

$$(\cos \theta - 1)v_{11} + \sin \theta = 0$$

$$v_{11} = \frac{-\sin \theta}{\cos \theta - 1}$$

$$(\cos \theta - 1)v_{11} + \sin \theta v_{12} = 0$$

$$\sin \theta v_{11} - (\cos \theta + 1)v_{12} = 0$$

$$(\sin \theta)v_{11} - \cos \theta - 1 = 0$$

$$v_{11} = \frac{\cos \theta}{\sin \theta} + \frac{1}{\sin \theta}$$

$$v_{11} = \cot(\theta) + \csc(\theta)$$

$$(\cos \theta - 1)\left(\frac{\cos \theta}{\sin \theta} + \frac{1}{\sin \theta}\right) + \sin \theta$$

$$\frac{\cos^2 \theta}{\sin \theta} - \frac{1}{\sin \theta} - \frac{\cos \theta}{\sin \theta} + \frac{\cos \theta}{\sin \theta} + \sin \theta$$

$$\frac{\cos^2 \theta - 1}{\sin \theta} + \sin \theta = \frac{-\sin^2 \theta}{\sin \theta} + \sin \theta$$

$$= 0 \checkmark$$

$$v_1 = \begin{bmatrix} \cot(\theta) + \csc(\theta) \\ 1 \end{bmatrix}$$

$$\lambda_2: \begin{bmatrix} \cos\theta + 1 & \sin\theta \\ \sin\theta & -\cos\theta + 1 \end{bmatrix} \begin{bmatrix} v_{21} \\ v_{22} \end{bmatrix} = 0$$

$$(\cos\theta + 1)v_{21} + (\sin\theta)v_{22} = 0$$

$$(\sin\theta)v_{21} + (-\cos\theta + 1)v_{22} = 0$$

for $v_{22} = 1$

$$v_{21} = \frac{\cos\theta}{\sin\theta} - \frac{1}{\sin\theta}$$

$$v_{21} = \cot(\theta) - \csc(\theta)$$

$$v_2 = \begin{bmatrix} \cot(\theta) - \csc(\theta) \\ 1 \end{bmatrix}$$

Eigenvalues are $\lambda_1 = 1$ and eigenvectors are:
 $\lambda_2 = -1$

$$v_1 = \begin{bmatrix} \cot(\theta) + \csc(\theta) \\ 1 \end{bmatrix}, \quad v_2 = \begin{bmatrix} \cot(\theta) - \csc(\theta) \\ 1 \end{bmatrix}$$

$$\begin{aligned} \vec{v}_1 \cdot \vec{v}_2 &= (\cot(\theta) + \csc(\theta))(\cot(\theta) - \csc(\theta)) + 1 \\ &= \cot^2(\theta) + \csc(\theta)\cot(\theta) - \csc^2(\theta) - \csc(\theta)\cot(\theta) + 1 \\ &= \frac{\cos^2\theta}{\sin^2\theta} - \frac{1}{\sin^2\theta} + 1 \\ &= \frac{\cos^2\theta - 1}{\sin^2\theta} + 1 \\ &= \frac{-\sin^2\theta}{\sin^2\theta} + 1 \\ &= -1 + 1 \\ &= 0 \end{aligned}$$

$$\boxed{\vec{v}_1 \cdot \vec{v}_2 = 0}$$

→ the dot product of the eigenvectors is 0, so they are orthogonal.

aii) Let A be an orthogonal matrix s.t. $A^T A = I$ is satisfied.
 Let λ_i be some eigenvalue associated with eigenvector v_i

$$A v_i = \lambda_i v_i$$

$$\|A v_i\|^2 = \|\lambda_i v_i\|^2 \quad \left. \begin{array}{l} \end{array} \right\} \text{since } \lambda_i \text{ is a scalar}$$

$$(A v_i)^T (A v_i) = |\lambda_i|^2 \|v_i\|^2$$

$$v_i^T A^T A v_i = |\lambda_i|^2 \|v_i\|^2$$

$$v_i^T I v_i = |\lambda_i|^2 \|v_i\|^2$$

$$\|v_i\|^2 = |\lambda_i|^2 \|v_i\|^2 \quad \left. \begin{array}{l} \end{array} \right\} \text{since } v_i^T v_i = \|v_i\|^2$$

$$\begin{array}{l} A, I \in \mathbb{R}^{m \times m} \\ v_i \in \mathbb{R}^{m \times 1} \\ (n \times m)(m \times m)(m \times n) \rightarrow n \times n \end{array}$$

$$|\lambda_i|^2 = 1$$

$|\lambda_i| = \pm 1 \rightarrow$ length cannot be negative, so

$$\boxed{|\lambda_i| = 1}$$

a.iii) $A \in \mathbb{R}^{n \times n}$

$$\left. \begin{array}{l} Av_1 = \lambda_1 v_1 \\ Av_2 = \lambda_2 v_2 \end{array} \right\} \lambda_1 \neq \lambda_2 \quad (\text{distinct eigenvalues})$$

$$\begin{aligned} v_1^T v_2 &= v_1^T I v_2 \quad \rightarrow A^T A = I \\ &= v_1^T A^T A v_2 \\ &= (A v_1)^T (A v_2) \\ &= \lambda_1 v_1^T \lambda_2 v_2 \\ &= \lambda_1 \lambda_2 v_1^T v_2 \end{aligned}$$

Since $\lambda_i \in \{-1, 1\}$, if $\lambda_1 \neq \lambda_2$ then $\lambda_1 \lambda_2 \neq 1$ and $v_1^T v_2 = 0$, meaning $v_1 \perp v_2$.

a.iv) The transformation Ax maps some coefficient $\wedge a_{ij} \in A$ to $x_i \in \bar{x}$ $\forall i, j$. The coefficients in A are associated with some variable in x . This forms a system of equations when a b is given s.t. $Ax = b$.

b.i) The eigenvectors of AA^T are the left-singular vectors of A .
The eigenvectors of $A^T A$ are the right-singular vectors of A .

b.ii) The square roots of the eigenvalues of $A^T A$ and AA^T are the singular (non-zero) values of A .

c i) \boxed{F} Can have an operator w/ no eigenvalues

c ii) \boxed{F}

c iii) \boxed{T} - IF A is PSD, $x^T A x \geq 0 \forall x$, then its eigenvalues are all positive or zero-valued

c iv) \boxed{F} By rank-nullity theorem, $\text{rank}(A) + \text{nullity}(A) = n$ for $A \in \mathbb{R}^{n \times n}$,
thus, $\text{rank}(A)$ cannot exceed # of non-zero eigenvalues - it is in fact equal to the # of non-zero eigenvalues.

c v) \boxed{T}
 $A v_1 = \lambda v_1$
 $A v_2 = \lambda v_2 \quad \Rightarrow \quad u = v_1 + v_2 \quad \text{Prove: } Au = \lambda u$

In general, let $v_1, \dots, v_n \in \mathbb{R}^m$ be the eigenvectors of A sharing eigenvalue λ .

$$A v_i = \lambda v_i \quad \forall v_i \in \{v_1, \dots, v_n\}$$

$$\left. \begin{aligned} A u &= A \sum_{i=1}^m v_i \\ &= \sum_{i=1}^m A v_i \\ &= \sum_{i=1}^m \lambda v_i \\ &= \lambda \sum_{i=1}^m v_i \end{aligned} \right\} \text{ where } u = \sum_{i=1}^m v_i$$

\longrightarrow for $m=2$ case:

$$A u = \lambda u$$

$$\begin{aligned} A u &= A(v_1 + v_2) \\ &= A v_1 + A v_2 \\ &= \lambda v_1 + \lambda v_2 \\ &= \lambda(v_1 + v_2) \end{aligned}$$

$$A u = \lambda u$$

$$\textcircled{2} \text{ ai) } P(c = H50 | x = \text{tails}) = \frac{P(x = \text{tails} | c = H50) P(c = H50)}{P(x = \text{tails})}$$

$$= \frac{(0.5)(0.5)}{P(x = \text{tails})}$$

$$\begin{aligned} \rightarrow P(x = \text{tails}) &= P(x = \text{tails} | c = H50) P(c = H50) \\ &\quad + P(x = \text{tails} | c = H60) P(c = H60) \\ &= (0.5)(0.5) + (0.4)(0.5) \\ &= 0.25 + 0.2 \end{aligned}$$

$$P(x = \text{tails}) = 0.45$$

$$= \frac{(0.5)(0.5)}{0.45}$$

$$\boxed{P(c = H50 | x = \text{tails}) = 0.55}$$

$$\text{aii) } P(c = H50 | x_1 = T, x_2 = H, x_3 = H, x_4 = H) = P(c = H50 | T, H, H, H)$$

$$\begin{aligned} P(c = H50 | T, H, H, H) &= P(c = H50) P(x = T | c = H50) P(x = H | c = H50) P(x = H | c = H50) P(x = H | c = H50) \\ &= (0.5)(0.5)^4 \end{aligned}$$

$$\boxed{P(c = H50 | T, H, H, H) = 0.03125}$$

$$\text{aiii) } \left. \begin{array}{l} H50 = 0.5 \\ H55 = 0.55 \\ H60 = 0.6 \end{array} \right\} \text{ p of heads}$$

$$\begin{aligned} P(c = H50 | 9H, 1T) &= P(c = H50) P(x = H | c = H50)^9 P(x = T | c = H50) \\ &= \left(\frac{1}{3}\right) \left(\frac{1}{2}\right)^9 \left(\frac{1}{2}\right) \\ &= \left(\frac{1}{3}\right) \left(\frac{1}{2^{10}}\right) \end{aligned}$$

$$\boxed{P(c = H50 | 9H, 1T) = 0.00033}$$

$$\begin{aligned} P(c = H55 | 9H, 1T) &= P(c = H55) P(x = H | c = H55)^9 P(x = T | c = H55) \\ &= \left(\frac{1}{3}\right) (0.55)^9 (0.45) \end{aligned}$$

$$\boxed{P(c = H55 | 9H, 1T) = 0.00669}$$

$$\begin{aligned} P(c = H60 | 9H, 1T) &= P(c = H60) P(x = H | c = H60)^9 P(x = T | c = H60) \\ &= \left(\frac{1}{3}\right) (0.6)^9 (0.4) \end{aligned}$$

$$\boxed{P(c = H60 | 9H, 1T) = 0.00134}$$

$$\begin{aligned}
 b) \quad & p(t=P | w=preg) = 0.99 \\
 & p(t=P | w=\neg preg) = 0.1 \\
 & \Rightarrow \\
 & p(w=preg) = 0.01 \\
 & p(w=\neg preg) = 0.99
 \end{aligned}$$

$$\begin{aligned}
 p(w=preg | t=P) &= \frac{p(t=P | w=preg) p(w=preg)}{p(t=P)} \\
 &= \frac{(0.99)(0.01)}{0.1089}
 \end{aligned}$$

$$\boxed{p(w=preg | t=P) = 0.091}$$

$$\begin{aligned}
 p(t=P) &= p(t=P | w=preg) p(w=preg) \\
 &\quad + p(t=P | w=\neg preg) p(w=\neg preg) \\
 &= (0.99)(0.01) + (0.1)(0.99) \\
 p(t=P) &= 0.1089
 \end{aligned}$$

Makes sense because most of the female population is not pregnant yet test returns positive for 10% of this majority population.

$$c) \quad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad E[\alpha f(x) + \beta g(x)] = \alpha E[f(x)] + \beta E[g(x)]$$

$$E(Ax + b) = E(Ax) + E(b)$$

$$\boxed{E(Ax + b) = AE[x] + b}$$

$$d) \quad cov(x) = E((x - E[x])(x - E[x])^T)$$

$$\begin{aligned}
 cov(Ax + b) &= E[(Ax + b - E[Ax + b])(Ax + b - E[Ax + b])^T] \\
 &= E[(Ax + b - AE[x] - b)(Ax + b - AE[x] - b)^T] \\
 &= E[(Ax - AE[x])(Ax - AE[x])^T] \\
 &= E[A^2(x - E[x])(x - E[x])^T] \\
 &= A^2 E[(x - E[x])(x - E[x])^T]
 \end{aligned}$$

$$\boxed{cov(Ax + b) = A^2 cov(x)}$$

$$③ \quad a) \quad x \in \mathbb{R}^{n \times 1}, y \in \mathbb{R}^{m \times 1}, A \in \mathbb{R}^{n \times m}$$

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

$$\rightarrow \nabla_x x^T A y \quad \begin{matrix} (1 \times n) & (n \times m) & (m \times 1) \end{matrix} \rightarrow (1 \times 1)$$

$$\begin{bmatrix} a_1 \end{bmatrix} y_1 + \dots + \begin{bmatrix} a_m \end{bmatrix} y_m$$

$$\frac{\partial (x^T a)}{\partial x} = \frac{\partial (a^T x)}{\partial x} = a^T \quad [x_1 \dots x_n] \begin{bmatrix} s_1 \\ \vdots \\ s_n \end{bmatrix} = s_1 x_1 + \dots + s_n x_n$$

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

$$\boxed{\nabla_y x^T A y = x^T A}$$

$$\begin{matrix} x^T A y \\ (1 \times n) (n \times m) (m \times 1) \end{matrix}$$

$$[\longleftrightarrow]_{1 \times m} [y]$$

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

$$\nabla_A x^T A y = \nabla_A x^T \sum_{i=1}^n \sum_{j=1}^m A_{ij} y_j$$

$$\begin{matrix} x^T A y \\ (1 \times n) \quad (n \times 1) \end{matrix}$$

$$= \nabla_A x^T \sum_{j=1}^m \sum_{i=1}^n A_{ij} y_j$$

$$= x^T \sum_{j=1}^m \sum_{i=1}^n \nabla_A A_{ij} y_j$$

$$= x^T \sum_{j=1}^m \sum_{i=1}^n y_j$$

$$= \sum_{j=1}^m \sum_{i=1}^n x_j y_j$$

$$\boxed{\nabla_A x^T A y = x y^T}$$

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

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

$$= (A x)^T + x^T A + b^T$$

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

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

$$\frac{\partial (x^T)}{\partial x} = \frac{\partial (a^T x)}{\partial x} = a^T$$

$$\frac{\partial ((a x^T)^T)}{\partial x} = a^T$$

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

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

$$\text{Assume } A \in \mathbb{R}^{m \times n} \\ B \in \mathbb{R}^{n \times m}$$

$$\text{tr}(AB) = \sum_{i=1}^m \left(\sum_{j=1}^n A_{ij} B_{ji} \right)$$

$$= \sum_{i=1}^m \sum_{j=1}^n A_{ij} B_{ji} = \sum_{j=1}^n \sum_{i=1}^m B_{ji} A_{ij}$$

$$\nabla_A \left(\sum_{i=1}^m \sum_{j=1}^n A_{ij} B_{ji} \right) = \sum_{i=1}^m \sum_{j=1}^n \frac{\partial (A_{ij} B_{ji})}{\partial A_{ij}}$$

$$\frac{\partial A_{ij} B_{ji}}{\partial A_{ij}} = B_{ji}$$

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

④

$$\hat{y} = Wx$$

$$\min_w \frac{1}{2} \sum_{i=1}^n \|y^{(i)} - Wx^{(i)}\|^2 = \min_w \frac{1}{2} \|\bar{y} - XW\|^2 \quad \text{drop } \frac{1}{2}$$

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

$$\|y - XW\|^2 = (y - XW)^T (y - XW)$$

$$= (y^T - X^T W^T) (y - XW)$$

$$= y^T y + W^T X^T X W - W^T X^T y - y^T X W$$

$$\begin{matrix} (1 \times n)(n \times 1) & (1 \times m)(m \times n)(n \times m)(m \times 1) & (1 \times m)(m \times n)(n \times 1) & (1 \times n)(n \times m)(m \times 1) \end{matrix}$$

$$\|y - XW\|^2 = y^T y + W^T X^T X W - 2W^T X^T y$$

$$\nabla_w \|y - XW\|^2 = \nabla_w (y^T y) + \nabla_w (W^T X^T X W) - \nabla_w (2W^T X^T y)$$

$$\nabla_w \|y - XW\|^2 = (X^T X)W - 2X^T y$$

$$\nabla_w = 0 \rightarrow (X^T X)W - 2X^T y = 0$$

$$(X^T X)W = X^T y$$

$$W = (X^T X)^{-1} X^T y$$

$$\begin{aligned} y &\in \mathbb{R}^{n \times 1} \\ X &\in \mathbb{R}^{n \times m} \\ W &\in \mathbb{R}^{m \times 1} \end{aligned}$$

$$\begin{aligned} y^T &\in \mathbb{R}^{1 \times n} \\ X^T &\in \mathbb{R}^{m \times n} \\ W^T &\in \mathbb{R}^{1 \times m} \end{aligned}$$

set to 0 to find min

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

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

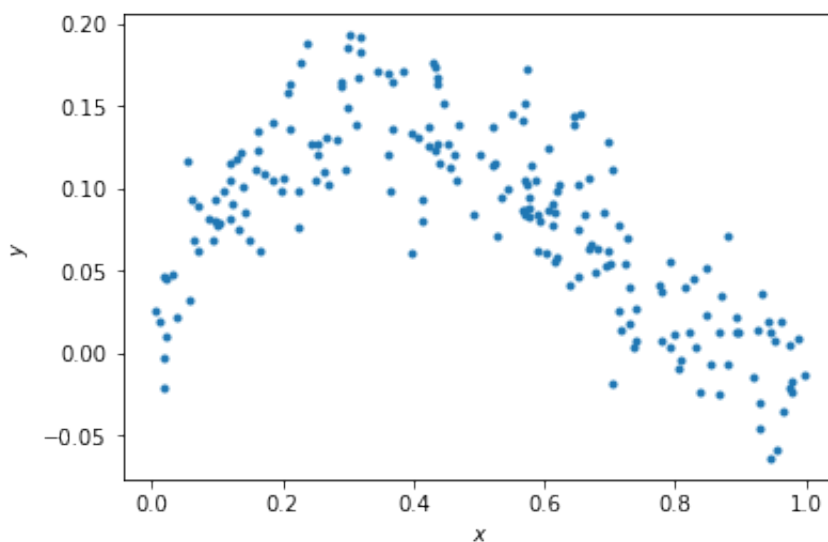
Data generation

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 [479]: 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[479]: 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) The generating distribution of x is uniform, since we're creating it via `np.random.uniform`.
- (2) The distribution of the additive noise is Gaussian, since we're creating it via `np.random.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 [480]: # xhat = (x, 1)
xhat = np.vstack((x, np.ones_like(x)))
#print(xhat)

# ===== #
# START YOUR CODE HERE #
# ===== #
# GOAL: create a variable theta; theta is a numpy array whose elements
# are [a, b]

#using normal eq
theta = np.linalg.inv(xhat.dot(xhat.T)).dot(xhat.dot(y))

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

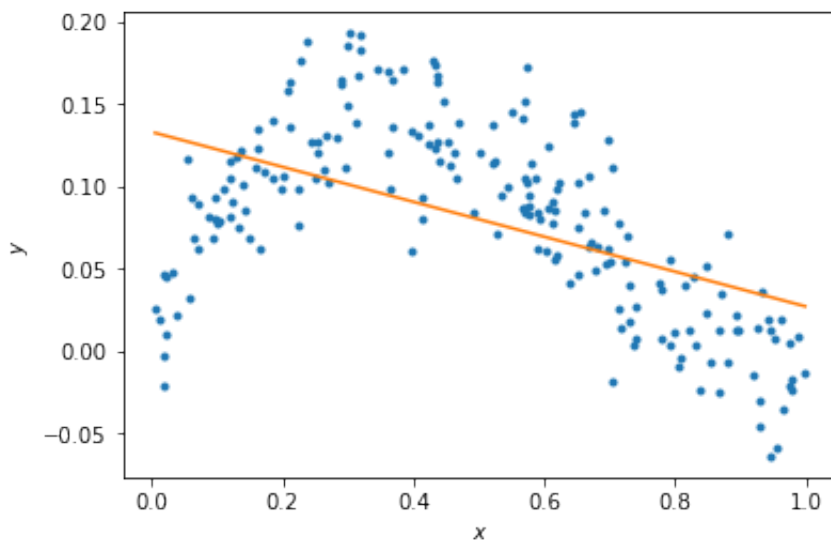


```
In [481]: # Plot the data and your model fit.
print(theta)
f = plt.figure()
ax = f.gca()
ax.plot(x, y, '.')
ax.set_xlabel('$x$')
ax.set_ylabel('$y$')

# 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))
```

```
[-0.10599633  0.13315817]
```

```
Out[481]: [<matplotlib.lines.Line2D at 0x114e09ac8>]
```



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 model underfits the data, since it is too simple. This data appears to resemble some higher order polynomial function (such as a parabola), but our linear model has no higher order terms.

(2) We can introduce higher order terms in order to improve the fit.

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

# ===== #
# START YOUR CODE HERE #
# ===== #
#create xhats
xhats.append(xhat)
#print(xhats[0])
for i in range(1,N):
    arr = np.vstack((np.array(xhat[0]**(i+1)), xhats[i-1]) )
    xhats.append(arr)

#print(xhats[3])

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

for i in range(0, N):
    cur_theta = np.linalg.inv(xhats[i].dot(xhats[i].T)).dot(xhats[i].dot(y))
    thetas.append(cur_theta)

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

```

```

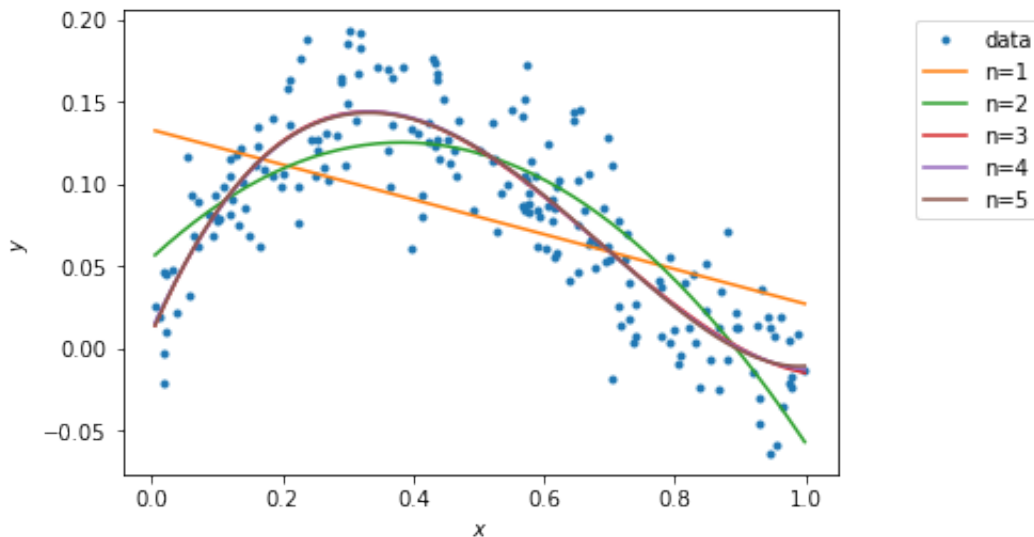
In [483]: # Plot the data
f = plt.figure()
ax = f.gca()
ax.plot(x, y, '.')
ax.set_xlabel('$x$')
ax.set_ylabel('$y$')

# 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), 200), np.ones(200)))
    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 [484]: 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 range(0, N):
    mse = ((y - thetas[i].dot(xhats[i]))**2).mean(axis=None)
    training_errors.append(mse)

min_error_index = np.argmin(training_errors) + 1
print("Best polynomial is order: " + str(min_error_index))

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

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

Best polynomial is order: 5
Training errors are:
[0.0023799610883627007, 0.001092492220926853, 0.0008169603801105369
4, 0.00081653537352969791, 0.00081614791955252942]
```

QUESTIONS

(1) What polynomial has the best training error?

(2) Why is this expected?

ANSWERS

(1) The 5th order polynomial has best fit.

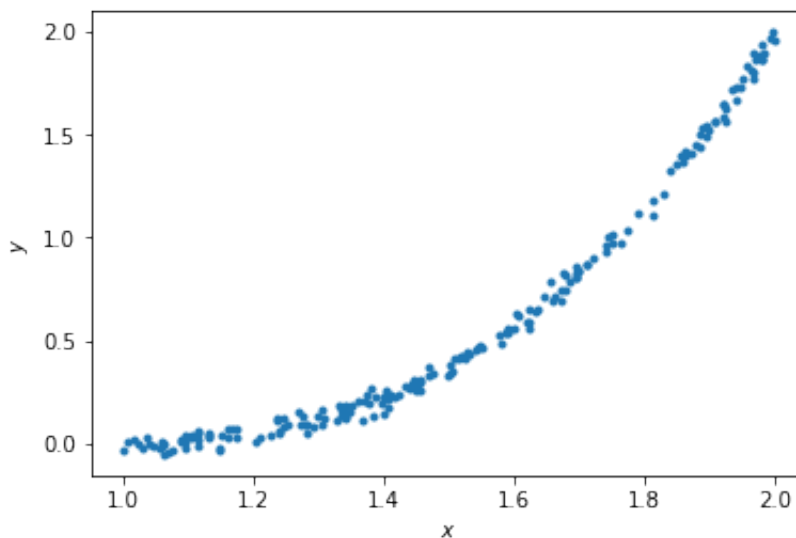
(2) This is expected since it is our most complicated model (it has the highest possible degree term). Since the goal is to minimize the cost function, weights will be adjusted such that the distribution can be best approximated. Given enough complexity, models will overfit and "memorize" the data, as shown here. It's probable that while the 5th order polynomial performed best in terms of training error, it will do poorly in terms of test error for a newly generated set. Models that have memorized the training set do not generalize well.

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 [485]: 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[485]: Text(0,0.5,'\$y\$')




```
In [486]: 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))

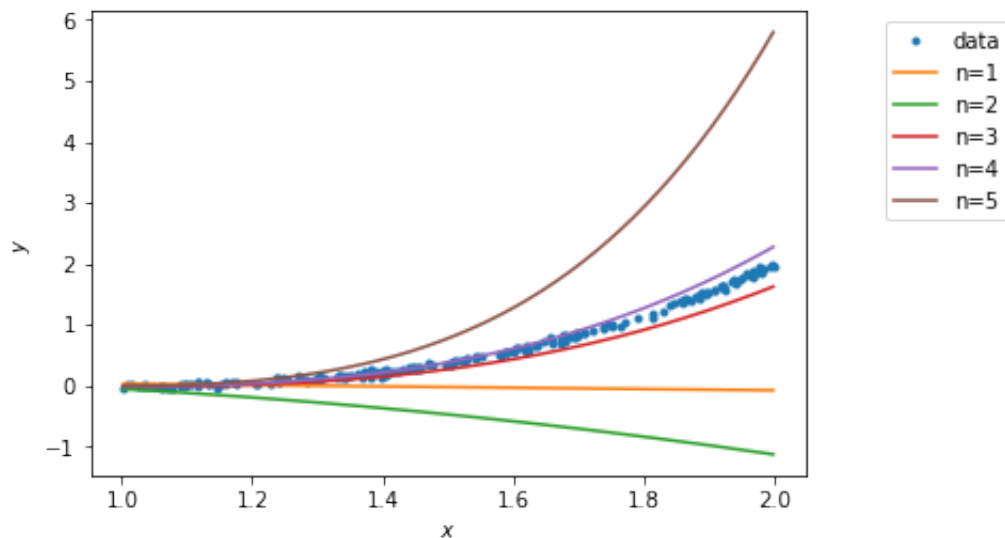
          xhats.append(xhat)
```

```
In [487]: # 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), 200), np.ones(200)))
              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 [488]: 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 range(0, N):
    mse = ((y - thetas[i].dot(xhats[i]))**2).mean(axis=None)
    testing_errors.append(mse)

index_min_error = np.argmin(testing_errors)
print("Lowest test error for polynomial: " + str(index_min_error + 1))
# ===== #
# END YOUR CODE HERE #
# ===== #

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

Lowest test error for polynomial: 4
Testing errors are:
[0.80861651845505866, 2.1319192445057897, 0.031256971083289675, 0.0
11870765196608337, 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 4th order polynomial has the best test error.
- (2) Order 5 polynomial models do not generalize well because they are overfitting. This means that the model is too complex and is able to effectively memorize the test data. Since linear regression "learning" is based on adjust weights the minimize a cost function, a complex higher order function that can memorize the test data will do very well in training (since the criteria is loss w/ respect to training data) but will do poorly on data that has not been incorporated into the generalization. Order 4 seems to be sufficiently complex such that it can reasonably output values from the target distribution without overfitting.