

Problem 1: Linear algebra refresher

a) i Let $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ $AA^T = I$

$$\Rightarrow \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} a & c \\ b & d \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \begin{array}{l} a^2 + b^2 = 1 \\ ac + bd = 0 \\ b^2 + d^2 = 1 \end{array}$$

By inspection:

If $a^2 = \frac{1}{2}, b^2 = \frac{1}{2}$, then $a, b = \pm\sqrt{\frac{1}{2}}$ and $a^2 + b^2 = 1$

Some for $c^2, d^2 = \frac{1}{2}$, but if $c = -\sqrt{\frac{1}{2}}$, then $ac + bd = -\frac{1}{2} + \frac{1}{2} = 0 \checkmark$

$$a, b, d = \pm\sqrt{\frac{1}{2}} \quad c = -\sqrt{\frac{1}{2}}$$

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

$$Av = \lambda v$$

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

For non-trivial solutions solve $|A - \lambda I| = 0$

$$\Rightarrow \left| \begin{bmatrix} \sqrt{\frac{1}{2}} - \lambda & \sqrt{\frac{1}{2}} \\ \sqrt{\frac{1}{2}} & -\sqrt{\frac{1}{2}} - \lambda \end{bmatrix} \right| = 0 \quad \left(\sqrt{\frac{1}{2}} - \lambda \right) \left(-\sqrt{\frac{1}{2}} - \lambda \right) - \frac{1}{2} = 0$$

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

$$\lambda^2 = 1$$

$$\lambda_1 = 1, \lambda_2 = -1$$

$$(A - \lambda_1 I)v_1 = 0 \quad \text{Let } v_{11} = 1 \Rightarrow \sqrt{\frac{1}{2}} - 1 + \sqrt{\frac{1}{2}}v_{12} = 0 \quad v_{12} = \sqrt{\frac{1}{2}}(1 - \sqrt{\frac{1}{2}}) = \sqrt{2} - 1$$

\hookrightarrow scale invariance of eigenvectors

$$v_1 = \begin{bmatrix} 1 \\ \sqrt{2} - 1 \end{bmatrix}$$

$$(A - \lambda_2 I)v_2 = 0 \quad \text{Let } v_{21} = 1 \Rightarrow \sqrt{\frac{1}{2}} + 1 + \sqrt{\frac{1}{2}}v_{22} = 0 \quad v_{22} = \sqrt{\frac{1}{2}}(-1 - \sqrt{\frac{1}{2}}) = -\sqrt{2} - 1$$

$$v_2 = \begin{bmatrix} 1 \\ -\sqrt{2} - 1 \end{bmatrix}$$

$v_1 \perp v_2$

λ 's are 1, -1 and eigenvectors are orthogonal

ii

Eig definition

$$Av = \lambda v$$

Norm 2 of both sides

$$\|Av\|_2 = \|\lambda v\|_2$$

$$\begin{aligned} &\Rightarrow \sqrt{(Av)^T Av} = \|\lambda\| \|v\| \\ &\quad \text{scalar} \\ &\quad = v^T A^T Av = v^T I v = v^T v \\ &\Rightarrow \sqrt{v^T v} = \|\lambda\| \sqrt{v^T v} \Rightarrow \|\lambda\| = 1 \end{aligned}$$

Thus A has λ 's where $\|\lambda\|=1$

iii

Goal: show $v_1^T v_2 = 0$ if $\lambda_1 \neq \lambda_2$

$$A v_1 = \lambda_1 v_1 \quad A v_2 = \lambda_2 v_2$$

$$\begin{aligned} v_1^T v_2 &= v_1^T A^T A v_2 = \{(Av_1)^T A v_2 = (\lambda_1 v_1)^T \lambda_2 v_2\} \\ &\Rightarrow v_1^T A^T A v_2 = \lambda_1 \lambda_2 v_1^T v_2 \end{aligned}$$

$$(1 - \lambda_1 \lambda_2) v_1^T v_2 = 0$$

Since $\|\lambda\|=1$, and $\lambda_1 \neq \lambda_2$, then $v_1^T v_2 = 0$,
Thus $v_1 \perp v_2$

iv

x can be reflected and rotated, but it cannot be scaled based on the properties derived above

b) i

$$\text{SVD } A : A = U \Sigma V^T \quad \text{where } U \text{ and } V \text{ are orthogonal} \\ (U^T U, V^T V = I)$$
$$A A^T = U \Sigma V^T (U \Sigma V^T)^T \\ = U \Sigma V^T V \Sigma^T U^T \\ = U \Sigma \Sigma^T U^T \quad \begin{matrix} \text{left sv of } A \\ \text{is eigv of } A A^T \end{matrix}$$
$$A^T A = (U \Sigma V^T)^T U \Sigma V^T \\ = V \Sigma^T U^T U \Sigma V^T \\ = V \Sigma^T \Sigma V^T \quad \begin{matrix} \text{right sv of } A \\ \text{is eigv of } A^T A \end{matrix}$$

Thus the left singular vectors of A are eigenvectors of AA' ,
and right singular vectors of A are eigenvectors of $A'A$

ii

$$\text{From above } A = U \underline{\Sigma} V^T$$

$$A A^T = U \underline{\Sigma} \Sigma^T U^T \quad \Sigma \Sigma^T = \Sigma^T \Sigma \\ A^T A = V \underline{\Sigma^T \Sigma} V^T \quad = [\Sigma_1^2, \Sigma_n^2]$$

The eigenvalues of $A'A$ and AA' are the squares of the singular values of A .

c) i

False They need not be distinct:

Counter example? $\text{eig}(I_2)$, λ 's = 1, 1

ii False Counter: for $\begin{bmatrix} 3 & 0 \\ 0 & 1 \end{bmatrix}$ $\begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ are eigenvectors for $\lambda = 1, 3$ respectively

$\begin{bmatrix} 0 \\ 1 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$, but $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ is not an eigenvector of I .

iii

True

$$A\vec{v} = \lambda\vec{v} \text{ eig definition}$$

$$\vec{x}^T A \vec{x} \geq 0 \quad \forall \vec{x} \quad \text{PSD definition}$$

$$\vec{v}^T A \vec{v} = \vec{v}^T \lambda \vec{v} = \vec{v}^T \vec{v} \lambda$$

If $\vec{v}^T A \vec{v} \geq 0$, and $\vec{v}^T \vec{v} \geq 0$ by definition,

$$\underline{\text{then } \lambda \geq 0}$$

iv

True

Rank-Nullity Theorem:

$$\text{rank}(A) + \text{nullity}(A) = n \quad \text{for}$$

$\text{rank}(A) = \# \text{ of eigenvalues}$

But they need not be distinct

For example : Rank $\left(\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \right) = 2$ $\lambda_1, \lambda_2 = 1, \lambda_3 = 0$
 Rank 2, 1 distinct, non-zero λ

v

True

The sum of two vectors is a linear combination.

If $a \in \lambda$ has two eigenvectors, then $\dim(E_\lambda) = 2$
reiginspace
Dimension of eigenspace

$$\text{Sp}(E_\lambda) = \{V_{1\lambda}, V_{2\lambda}\} \quad \text{if } V_1, V_2 \text{ as eigenvectors for } \lambda$$

$$1V_{1\lambda} + 1V_{2\lambda} \in \text{Sp}(E_\lambda)$$

A sum of two eigenvectors for λ will be in the span of E_λ , so it is still an eigenvector.

Problem 2: Probability refresher

Notation: $t = T$, $h = H$, $h^3 = H, H, H$

a) i) $P(H_{50}|t)?? \quad ; \quad P(t|H_{50}) = 0.5 \quad P(H_{50}) = 0.5$

$$P(t) = P(H_{50}) + P(t|H_{50})P(H_{60}) = 0.5 \cdot 0.5 + 0.5 \cdot 0.4 = 0.45$$

$$P(H_{50}|t) = \frac{P(t|H_{50})P(H_{50})}{P(t)} = \frac{0.5 \cdot 0.5}{0.45} = \frac{5}{9}$$

$$P(H_{50}|t) = \frac{5}{9}$$

ii) $P(H_{50}|t, h, h, h)$:

$$P(t, h, h, h | H_{50}) = 0.5^4 \quad P(H_{50}) = 0.5$$

$$P(t, h, h, h) = 0.5 \cdot 0.5^4 + 0.5(0.4 \cdot 0.6^3)$$

$$P(H_{50}|t, h^3) = \frac{P(t, h^3 | H_{50})P(H_{50})}{P(t, h^3)} = \frac{0.5^5}{0.5^5 + 0.5(0.4)(0.6)^3} \approx 0.4197$$

$$P(H_{50}|t, h^3) \approx 0.4197$$

iii)

$$P(H_{50}), P(H_{55}), P(H_{60}) = \frac{1}{3}$$

$$P(h^a t) = \frac{1}{3} \cdot 0.5^{10} + \underbrace{\frac{1}{3} \cdot 0.55^9 \cdot 0.45}_{P(h^a t | H_{55})} + \underbrace{\frac{1}{3} \cdot 0.60^9 \cdot 0.4}_{P(h^a t | H_{60})} \approx 0.0024$$

$$P(H_{50}|h^a t) = \frac{\frac{1}{3} \cdot 0.5^{10}}{P(h^a t)} \approx 0.1379$$

$$P(H_{55}|h^a t) = \frac{\frac{1}{3} \cdot 0.55^9 \cdot 0.45}{P(h^a t)} \approx 0.2927$$

$$P(H_{60} | h^9 t) = \frac{V_3 \cdot 0.68^9 \cdot 0.3}{P(h^9 t)} \approx 0.5694$$

$$P(H_{50} | h^9 t) \approx 0.1379$$

$$P(H_{55} | h^9 t) \approx 0.2927$$

$$P(H_{60} | h^9 t) \approx 0.5694$$

b)

A = Pregnant T = Positive Test

$$P(T|A) = 0.99 \quad P(T|\bar{A}) = 0.10 \quad P(A) = 0.01$$

$$P(A|T) = ??$$

$$P(T) = \frac{P(T|A)P(A)}{P(T|A)P(A) + P(T|\bar{A})P(\bar{A})} = \frac{0.99 \cdot 0.01}{0.99 \cdot 0.01 + 0.1 \cdot 0.99} \approx 0.1089$$

$$P(A|T) = \frac{P(T|A)P(A)}{P(T)} = \frac{0.99 \cdot 0.01}{0.1089} \approx 0.0901$$

$$P(\text{Pregnant} | \text{Positive}) = 0.0901$$

Note the high false positive rate makes the probability so low...

$$P(T|\bar{A}) = 0.1 > P(A) = 0.01$$

Hence more false positives are generated than there are pregnant women.

$$\begin{aligned} c) \quad & \mathbb{E}(Ax+b) = \underbrace{\mathbb{E}[Ax]}_{\substack{\text{Linearity of } \mathbb{E} \\ \Rightarrow b \text{ is deterministic}}} + \mathbb{E}[b] \\ & = \mathbb{E}[Ax] = A \mathbb{E}[x] \quad \mathbb{E}[Ax] = \mathbb{E}\left[\sum_{i=1}^n A_{ij} x_j\right] = \sum_{j=1}^n A_{ij} \mathbb{E}[x_j] \\ & \qquad \qquad \qquad \text{By determinism of } A \text{ and iid } x; \\ & = A \mathbb{E}[x] \end{aligned}$$

$$\text{Thus } \mathbb{E}[Ax+b] = A \mathbb{E}[x] + b$$

$$\begin{aligned}
 d) \quad \text{cov}(Ax+b) &= \mathbb{E}((Ax+b - \mathbb{E}[Ax+b])(Ax+b - \mathbb{E}[Ax+b])^T) \\
 &= \mathbb{E}[(Ax+b - A\mathbb{E}[x] - b)(Ax+b - A\mathbb{E}[x] - b)^T] \quad \text{from c) above} \\
 &= \mathbb{E}[A(x - \mathbb{E}[x]) (x - \mathbb{E}[x])^T A^T] \\
 &= A \underbrace{\mathbb{E}[(x - \mathbb{E}[x])(x - \mathbb{E}[x])^T]}_{= \text{cov}(x)} A^T \quad \text{determinism of } A + \text{iid } x \\
 &= \boxed{A \text{ cov}(x) A^T}
 \end{aligned}$$

Problem 3: Multivariate Derivatives

a) $\nabla_x x^T A y ? \quad x \in \mathbb{R}^n, \text{ so } \nabla_x x^T A y \in \mathbb{R}^n$

$$x^T A y = x_1 a_{11} y_1 + x_2 a_{12} y_1 + \dots + x_n a_{nn} y_n \quad \nabla_x x^T A y = a_{11} y_1 + a_{12} y_2 + \dots + a_{nn} y_n \\ = A y \in \mathbb{R}^n \checkmark$$

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

b) $\nabla_y x^T A y \quad \nabla_y x^T A y \in \mathbb{R}^m$

$$x^T A y = x_1 a_{11} y_1 + x_2 a_{12} y_1 + \dots + x_n a_{nn} y_n \quad \nabla_y x^T A y = x_1 a_{11} + x_2 a_{12} + \dots + x_n a_{nn} \\ = x^T A$$

\Rightarrow To match dimension ($n \times 1$) need $(x^T A)^T = A^T x \in \mathbb{R}^m \checkmark$

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

c) $\nabla_A x^T A y \quad \in \mathbb{R}^{n \times m}$

$$\Rightarrow$$

$$\left[\begin{array}{l} \frac{\partial x^T A y}{\partial a_{11}} = x_1 y_1, \dots, \frac{\partial x^T A y}{\partial a_{1n}} = x_1 y_n \\ \vdots \\ \frac{\partial x^T A y}{\partial a_{n1}} = x_n y_1, \dots, \frac{\partial x^T A y}{\partial a_{nn}} = x_n y_n \end{array} \right] \\ \Rightarrow x y^T \in \mathbb{R}^{n \times m} \checkmark$$

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

$$d) \nabla_x (x^T A x + b^T x) \in \mathbb{R}^n$$

$$= \nabla_x (x^T A x) + \nabla_x (b^T x) \quad \text{linearity of } \nabla$$

$$\nabla_x (x^T A x) = Ax + A^T x \quad \text{derived in class :}$$

$$A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \quad f(x) = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j$$

$$\sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j = \sum_{i=1}^n a_{ii} x_i^2 + 2 \sum_{i<j} a_{ij} x_i x_j$$

$$\Rightarrow \frac{\partial f(x)}{\partial x_i} = a_{ii} x_i + \sum_{j=1}^n a_{ij} x_j + \sum_{i>j} a_{ii} x_i$$

$$= \sum_{j=1}^n a_{ij} x_j + \sum_{i=1}^n a_{ii} x_i \quad \begin{matrix} \uparrow \\ (Ax)_i \end{matrix} \quad \begin{matrix} \uparrow \\ (A^T x) \end{matrix}$$

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

$$\in \mathbb{R}^n \checkmark$$

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

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

$$\text{tr}(AB) = \text{tr} \left(\begin{bmatrix} a_{11}b_{11} + a_{12}b_{21} + \dots + a_{1n}b_{n1} \\ \vdots \\ a_{m1}b_{11} + \dots + a_{mn}b_{n1} \end{bmatrix} \right)$$

$$= \sum_{i=1}^n a_{ii} b_{ii} + \sum_{i=1}^n a_{1i} b_{i1} + \dots + \sum_{i=1}^n a_{ni} b_{i1}$$

$$\frac{\partial \text{tr}(AB)}{\partial a_{11}} = b_{11} \quad \frac{\partial \text{tr}(AB)}{\partial a_{12}} = b_{21} \quad \frac{\partial \text{tr}(AB)}{\partial a_{1n}} = b_{1n}$$

$$\Rightarrow \frac{\partial \text{tr}(AB)}{\partial a_{ij}} = b_{ji}$$

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

Problem 4: Deriving LS

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

To minimize: $\frac{\partial \mathcal{L}}{\partial W} = 0 \quad \nabla_W \mathcal{L} = 0$

$$\begin{aligned}\mathcal{L} &= \frac{1}{2} \sum_{i=1}^n (y^{(i)} - Wx^{(i)})^T (y^{(i)} - Wx^{(i)}) \\ &= \frac{1}{2} \sum_{i=1}^n y^{(i)T} - 2y^{(i)T}Wx^{(i)} + x^{(i)T}W^TWx^{(i)} \quad \text{only need } W \text{ terms} \\ &= \sum_{i=1}^n -y^{(i)T}Wx^{(i)} + \frac{1}{2} x^{(i)T}W^TWx^{(i)} \quad \left. \right\} \in \mathbb{R} \text{ so } \text{tr}(.) = . \\ &= \sum_{i=1}^n -\text{tr}(y^{(i)T}Wx^{(i)}) + \frac{1}{2}\text{tr}(x^{(i)T}W^TWx^{(i)}) \quad \text{Notation: Dropping } '^{(i)} \text{ for readability...} \\ &= \underbrace{-\text{tr}(Wxy^T)}_{A} + \frac{1}{2}\text{tr}(Wxx^TW^T) \quad \text{Tr}(AB) = \text{Tr}(BA) \\ &= -\text{tr}(Wxy^T) + \frac{1}{2}\text{tr}(Wxx^TW^T) \quad \left. \begin{array}{l} \sum x^{(i)}y^{(i)T} = XY^T \\ \sum x^{(i)x^{(i)T}} = XX^T \end{array} \right\}\end{aligned}$$

Now $\nabla_W \mathcal{L}$,

$$\begin{aligned}\nabla_W \mathcal{L} &= -YX^T + \frac{1}{2}WX^T + \frac{1}{2}WXX^T \\ &= -YX^T + WXX^T = 0\end{aligned}$$

$$\left. \begin{array}{l} \frac{\partial \text{tr}(WA)}{\partial W} = A^T, \quad A = XY^T \\ \frac{\partial \text{tr}(WAW^T)}{\partial W} = WA^T + WA, \quad A = XX^T \end{array} \right\}$$

$$\Rightarrow W = (XY^T)(XX^T)^{-1}$$

Optimal W (to $\min \mathcal{L}(W)$) is $(XY^T)(XX^T)^{-1}$

linear_regression

January 11, 2022

0.1 Linear regression workbook: Sunay Bhat - W2022

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 2022, Prof. J.C. Kao, TAs Y. Li, P. Lu, T. Monsoor, T. wang

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

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

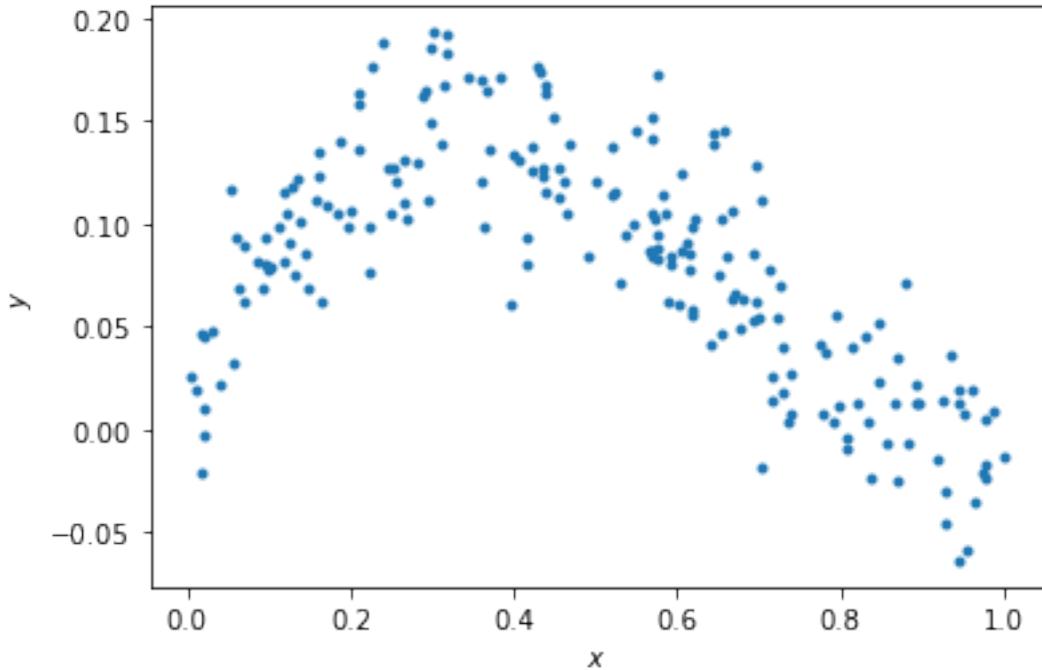
0.1.1 Data generation

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

```
[2]: 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('$y$')
```

```
[2]: 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) $U(0, 1)$ - Uniform distribution with parameters 0 and 1.
- (2) $N(0, 0.3)$ - Normal distribution with parameters mean $\mu = 0$ and standard deviation $\sigma = 0.3$.

0.1.4 Fitting data to the model (5 points)

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

```
[3]: # 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]
```

```
# Normal Equation:  $(X^T * X)^{-1} * X^T * y$ 
theta = np.linalg.inv(xhat.dot(xhat.T)).dot(xhat.dot(y))

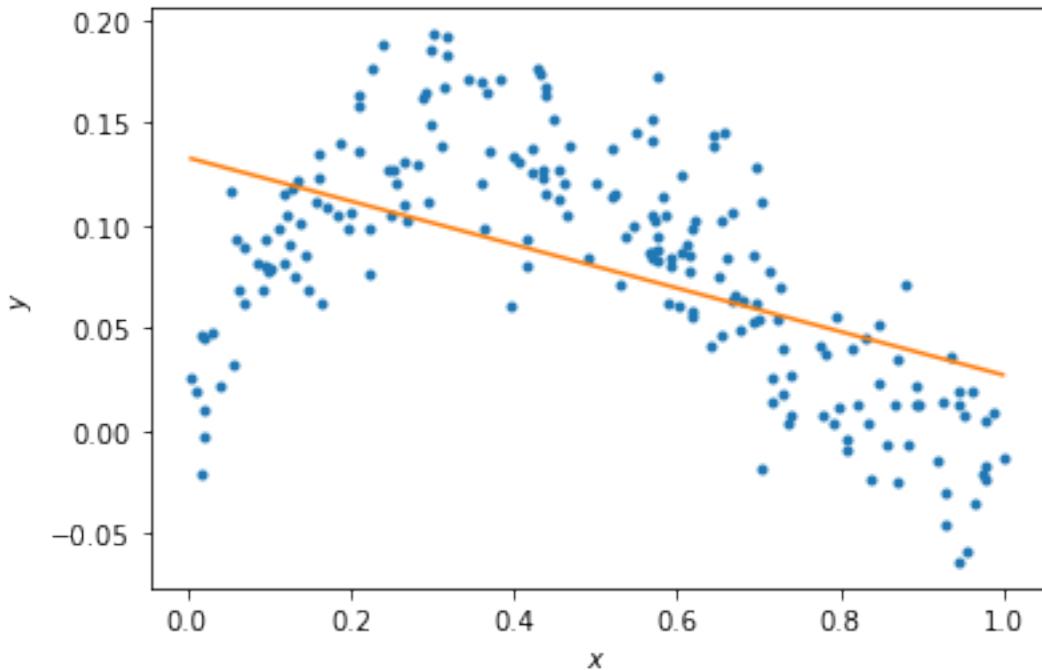
# ===== #
# END YOUR CODE HERE #
# ===== #
```

[4]: # 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))
```

[4]: [`<matplotlib.lines.Line2D at 0x7fb9242b5eb0>`]



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** the data.

(2) We can use a **higher order polynomial** like a 2nd or 3rd order polynomial or many other models with **higher complexity**.

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.

```
[5]: 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.

xhats.append(np.vstack((x, np.ones_like(x)))) # Initial xhat = (x, 1)
for order in np.arange(N):
    if order > 0:
        xhats.append(np.vstack((x***(order+1), xhats[order-1]))) #add next order
#features to last xhat and append

    # Normal Equation: (X'*X)^-1 * X'*y
    thetas.append(np.linalg.inv(xhats[order].dot(xhats[order].T)).dot(xhats[order].dot(y)))

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

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

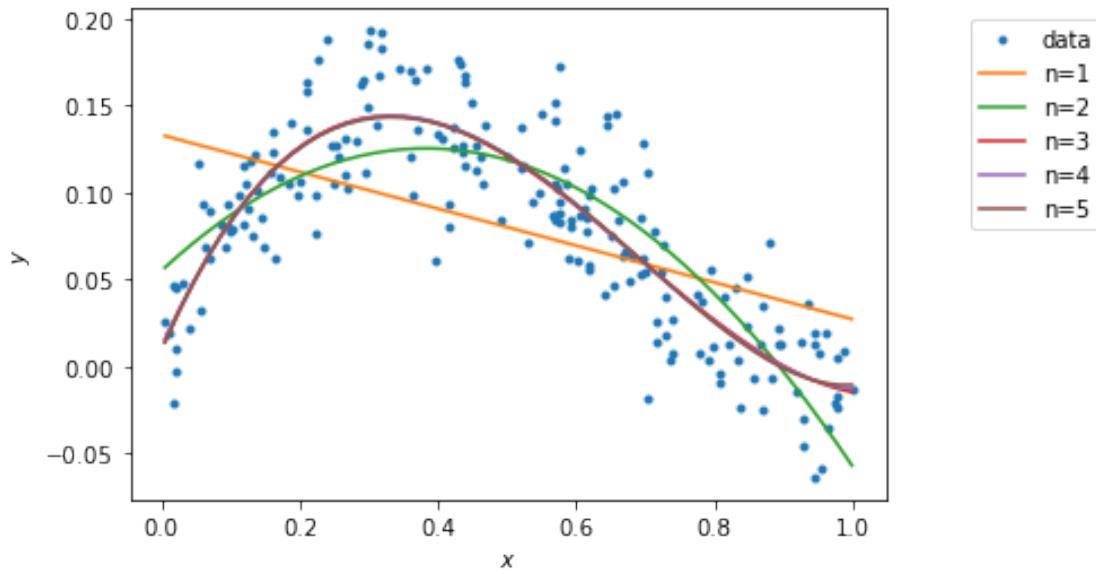
```

# 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_xs[-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:

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

[7]:

```

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 order in np.arange(N):
    #  $\frac{1}{2} * (\theta \cdot \hat{x} - y)' * (\theta \cdot \hat{x} - y)'$ 
    errors = thetas[order] @ xhats[order] - y
    training_errors.append(1/2 * (errors.T @ errors))

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

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

```

Training errors are:
[0.2379961088362701, 0.10924922209268528, 0.08169603801105374,
0.08165353735296982, 0.08161479195525295]

0.1.9 QUESTIONS

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

0.1.10 ANSWERS

(1) The **5th order** polynomial has the lowest training error.

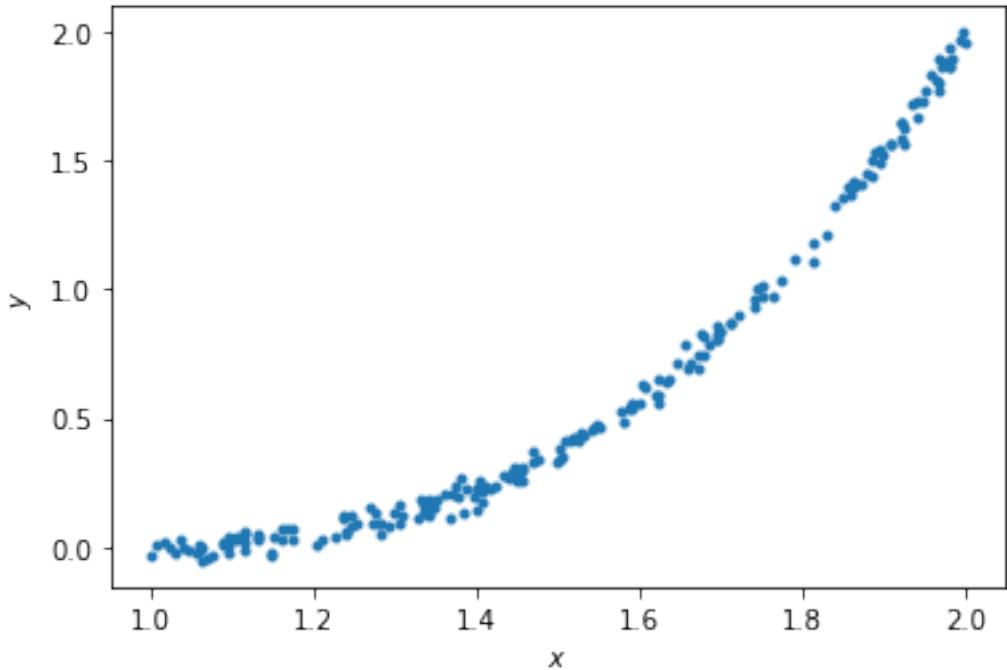
(2) Yes, higher order polynomials will have \leq training error than lower order. You can do as well as lower order polynomials by setting higher coefficients to zero (and others equal to each other), but you can generally utilize the added degrees of freedom from higher order terms to better fit the training data.

0.1.11 Generating new samples and validation error (5 points)

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

```
[8]: 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('$y$')
```

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



```
[9]: 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)
```

```
[10]: # 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:
```

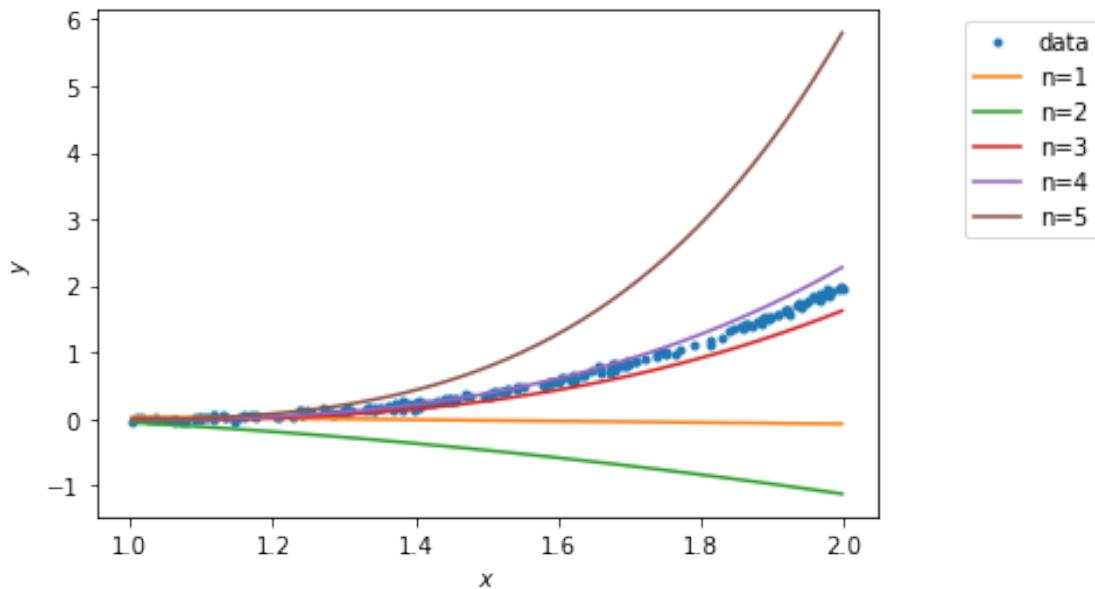
```

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)

```



```

[11]: validation_errors = []

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

# GOAL: create a variable validation_errors, a list of 5 elements,
# where validation_errors[i] are the validation loss for the polynomial fit of
# order i+1.

for order in np.arange(N):
    #  $\frac{1}{2} * (\theta^* \cdot \hat{x} - y)^T * (\theta^* \cdot \hat{x} - y)$ 
    errors = thetas[order] @ xhats[order] - y
    validation_errors.append(1/2 * (errors.T @ errors))

# ===== #

```

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

print ('Validation errors are: \n', validation_errors)
```

```
Validation errors are:
[80.86165184550586, 213.19192445057894, 3.125697108276393, 1.1870765189474703,
214.91021817652626]
```

0.1.12 QUESTIONS

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

0.1.13 ANSWERS

- (1) The **4th order** polynomial has the best, or lowest, validation error, with the 3rd order being close.
- (2) The 5th order model **overfits the data and does not generalize well when we validate our dataset in the previously unseen region of $x = [1 : 2]$** . In this case, the problem is made acute by validating in an extension of the domain ($x = [1 : 2]$) instead of new points from within the training data domain ($x = [0, 1]$). Specifically, the 4th order coefficient for the 5th order model is too large and rapidly diverges from the 3rd order generative model in the new domain space. The 4th order model's 4th order coefficient is much smaller though, allowing it to perform well even though the generative model is only 3rd order.

[]: