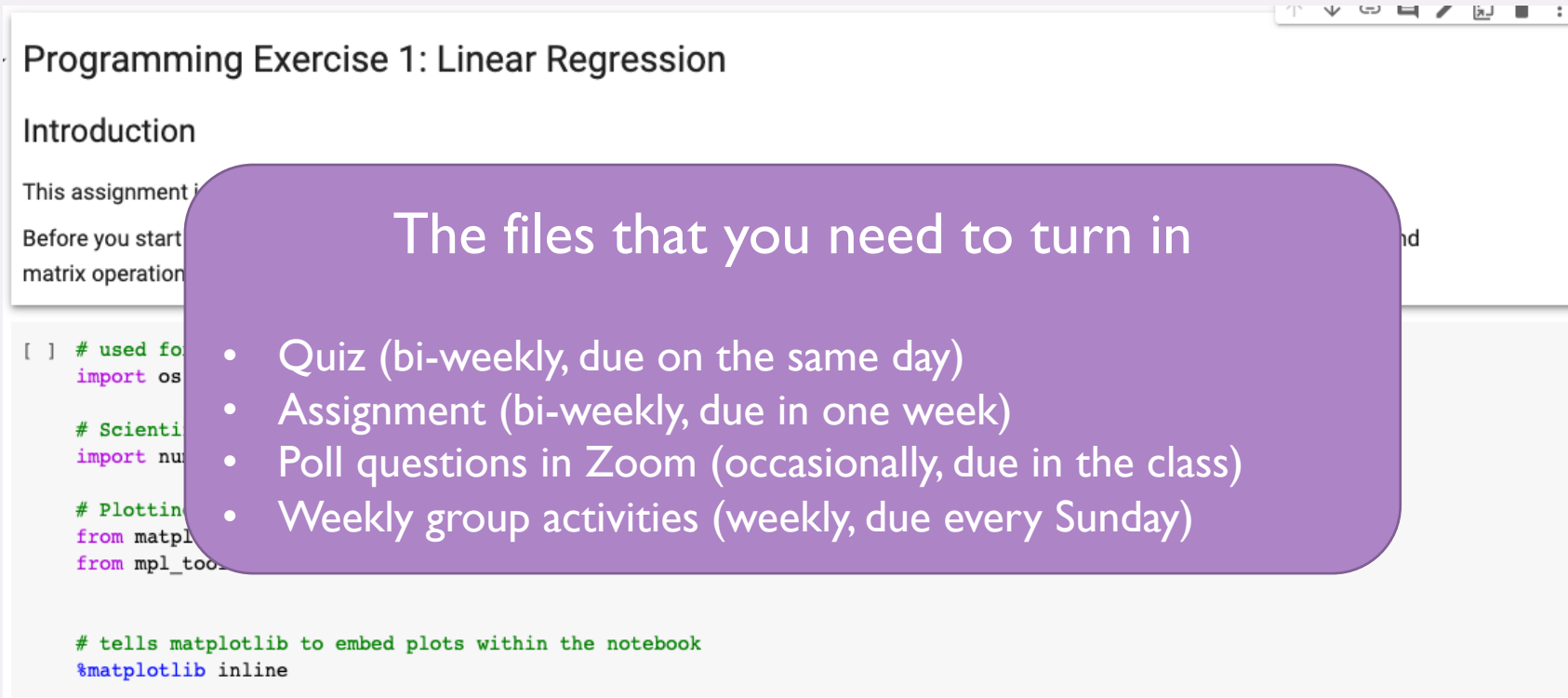# LECTURE 5

SPRING 2021
APPLIED MACHINE LEARNING
CIHANG XIE

1

# EXERCISES

- Google Colab Exercises (distributed during the lecture)

NO NEED TO SUBMIT! JUST TO HELP YOU UNDERSTAND LECTURE

Programming Exercise 1: Linear Regression

Introduction

This assignment

Before you start

matrix operation

```
[ ]  # used fo
     import os

     # Scienti
     import nu

     # Plottin
     from matpl
     from mpl_too
```
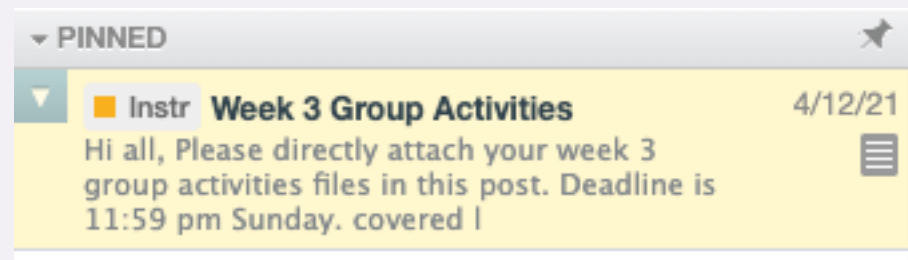
```
     # tells matplotlib to embed plots within the notebook
     %matplotlib inline
```

**The files that you need to turn in**

- Quiz (bi-weekly, due on the same day)
- Assignment (bi-weekly, due in one week)
- Poll questions in Zoom (occasionally, due in the class)
- Weekly group activities (weekly, due every Sunday)

# GROUP ACTIVITIES

- ALL groups have submitted the files

  https://piazza.com/class/kmmw9butjod4ay?cid=13

- Please read notes & exercises from other groups

- Week 3 (lecture 4-5)

# TODAY

- Review of Gradient Descent Algorithm

- Choosing Learning Rate $\alpha$

- Basis Functions

# LINEAR REGRESSION

# GRADIENT DESCENT

- Initialize $\theta$

- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial Cost(\theta)}{\partial \theta_j} \qquad \text{(simultaneous update for } \theta_0, \theta_1, \ldots, \theta_d)$$

- For linear regression:

$$\frac{\partial Cost(\theta)}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^{n} \left( h_\theta\left(x^{(i)}\right) - y^{(i)} \right)^2$$

With $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \ldots + \theta_d x_d = \sum_{j=0}^{d} \theta_j x_j$

# GRADIENT DESCENT

$$\frac{\partial Cost(\theta)}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^{n} \left(h_\theta\left(x^{(i)}\right) - y^{(i)}\right)^2$$

$$= \frac{1}{2n} \quad \frac{\partial}{\partial \theta_j} \sum_{i=1}^{n} \left(h_\theta\left(x^{(i)}\right) - y^{(i)}\right)^2$$

Scalar multiple rule

$$= \frac{1}{2n} \sum_{i=1}^{n} \quad \frac{\partial}{\partial \theta_j} \left(h_\theta\left(x^{(i)}\right) - y^{(i)}\right)^2$$

Sum rule

$$= \frac{1}{2n} \sum_{i=1}^{n} \quad 2\left(h_\theta\left(x^{(i)}\right) - y^{(i)}\right) \quad \frac{\partial}{\partial \theta_j} \left(h_\theta\left(x^{(i)}\right) - y^{(i)}\right)$$

Power rule

$$= \frac{1}{n} \sum_{i=1}^{n} \left(\sum_{k=0}^{d} \theta_k x_k^{(i)} - y^{(i)}\right) \quad \frac{\partial}{\partial \theta_j} \left(\sum_{k=0}^{d} \theta_k x_k^{(i)} - y^{(i)}\right)$$

$$= \frac{1}{n} \sum_{i=1}^{n} \left(\sum_{k=0}^{d} \theta_k x_k^{(i)} - y^{(i)}\right) x_j^{(i)}$$
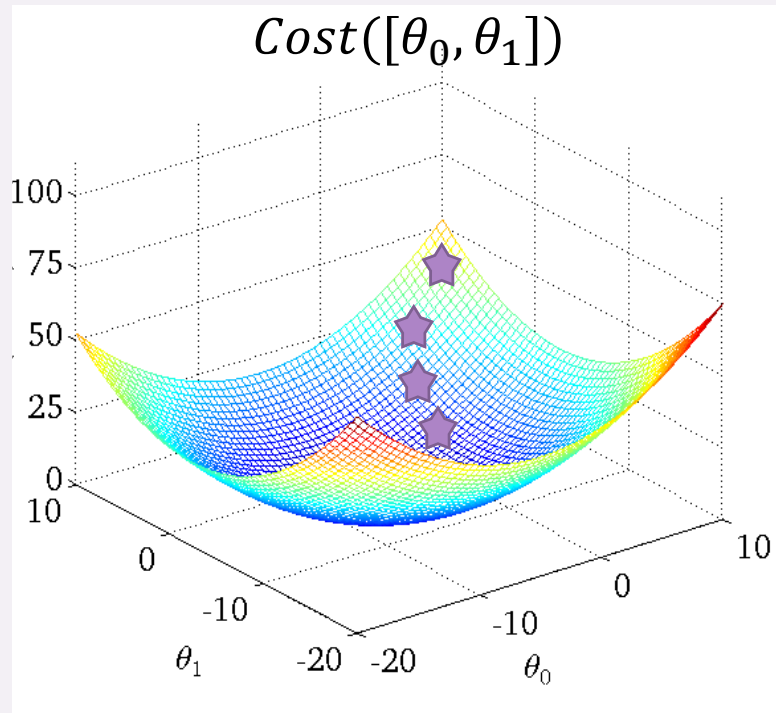
# GRADIENT DESCENT

- Initialize $\theta$

- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial Cost(\theta)}{\partial \theta_j} \qquad \text{(simultaneous update for } \theta_0, \theta_1, \ldots, \theta_d)$$

$$Cost([\theta_0, \theta_1])$$

**Left:**

$$Temp0 \leftarrow \theta_0 - \alpha \frac{\partial Cost\,(\theta_0, \theta_1)}{\partial \theta_0}$$

$$Temp1 \leftarrow \theta_1 - \alpha \frac{\partial Cost\,(\theta_0, \theta_1)}{\partial \theta_1}$$

$$\theta_0 \leftarrow Temp0$$

$$\theta_1 \leftarrow Temp1$$

**Right:**

$$Temp0 \leftarrow \theta_0 - \alpha \frac{\partial Cost\,(\theta_0, \theta_1)}{\partial \theta_0}$$
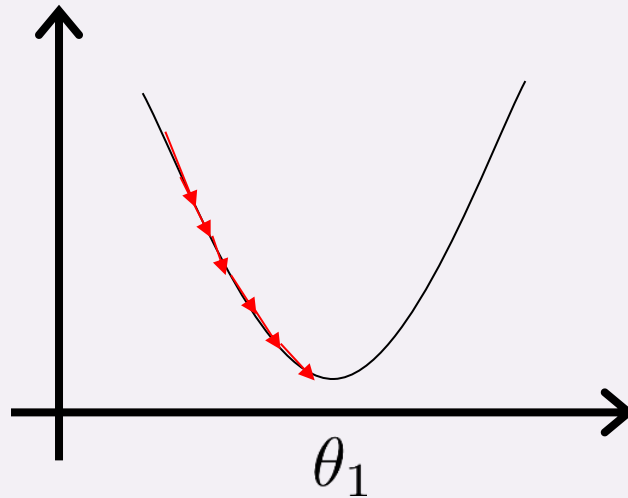
$$\theta_0 \leftarrow Temp0$$

$$Temp1 \leftarrow \theta_1 - \alpha \frac{\partial Cost\,(\theta_0, \theta_1)}{\partial \theta_1}$$

$$\theta_1 \leftarrow Temp1$$

# CHOOSE $\alpha$

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial Cost(\theta)}{\partial \theta_j}$$

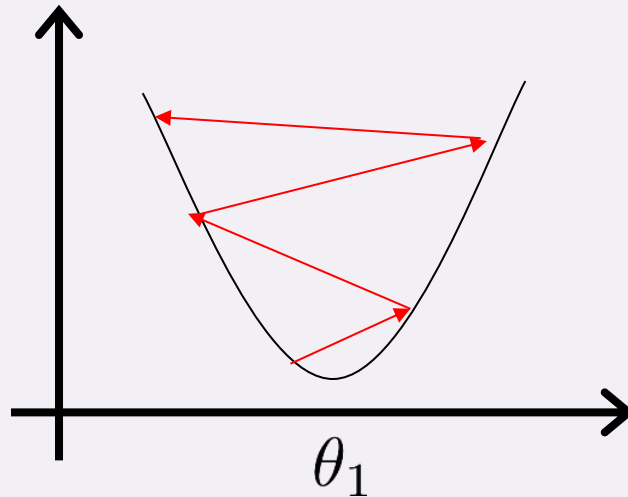(simultaneous update for $\theta_0, \theta_1, \ldots, \theta_d$)

Learning rate

**If α is too small, gradient descent can be very slow.**

$\theta_1$

# CHOOSE $\alpha$

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial Cost(\theta)}{\partial \theta_j}$$

(simultaneous update for $\theta_0, \theta_1, \dots, \theta_d$)

Learning rate



$\theta_1$

**If α is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.**

# CHOOSE $\alpha$

$$\theta_j \leftarrow \theta_j \; - \alpha \frac{\partial Cost(\theta)}{\partial \theta_j} \qquad \text{(simultaneous update for } \theta_0, \theta_1, \dots, \theta_d)$$

Learning rate

For certain functions $Cost(\theta),$ we can theoretically guarantee the convergence of gradient descent by choosing a appropriate $\alpha$

If interested, please read machine learning course at UBC: lecture 4, starting from page 9
https://www.cs.ubc.ca/~schmidtm/Courses/540-W18/L4.pdf

# EXTENDING LINEAR REGRESSION TO MORE COMPLEX MODELS

- The inputs $X$ for linear regression can be:

  - Original quantitative inputs

  - Transformation of quantitative inputs (log, exp, square, etc.)

  - Polynomial transformation (example: $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$)

  - Interactions between variables (example: $x_3 = x_1 \times x_2$)

- This allows use of linear regression techniques to fit non-linear datasets.

# LINEAR BASIS FUNCTION MODEL

- Generally,
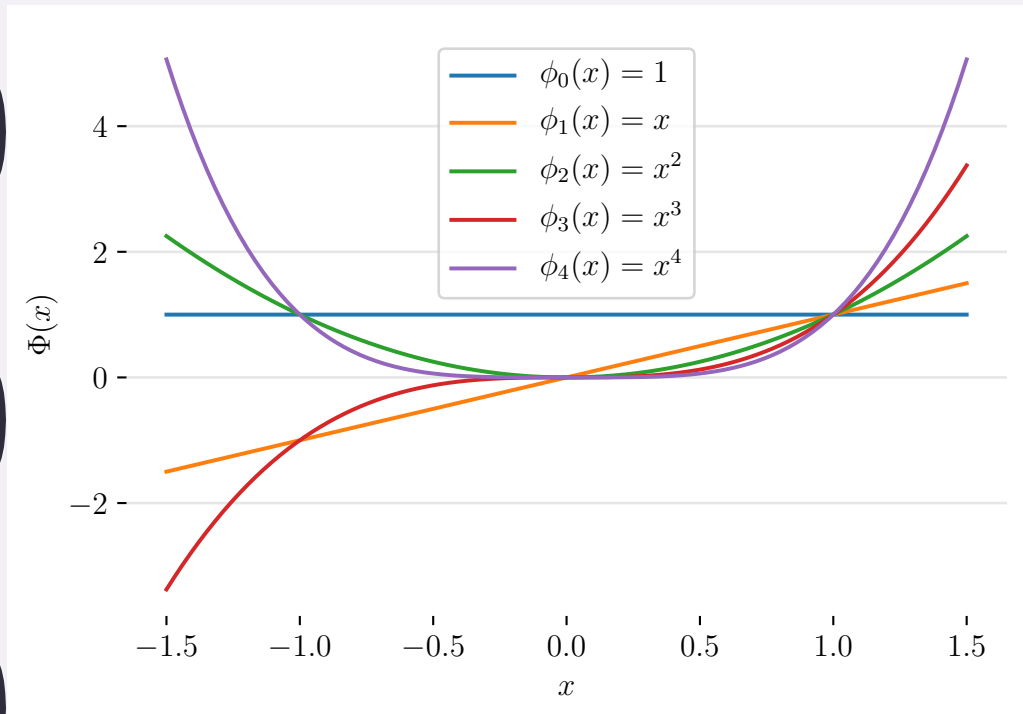
$$h_\theta(x) = \sum_{j=0}^{d} \theta_j \phi_j(x)$$

Basis Function

- Typically, $\phi_0(x) = 1$ so that $\theta_0$ acts as a bias.

- In the simplest case, we can use linear basis function:
$$\phi_j(x) = x_j$$

- Polynomial basis function: $\phi_j(x) = x^j$

- Gaussian basis function: $\phi_j(x) = e^{-\frac{(x-\mu_j)^2}{2s^2}}$
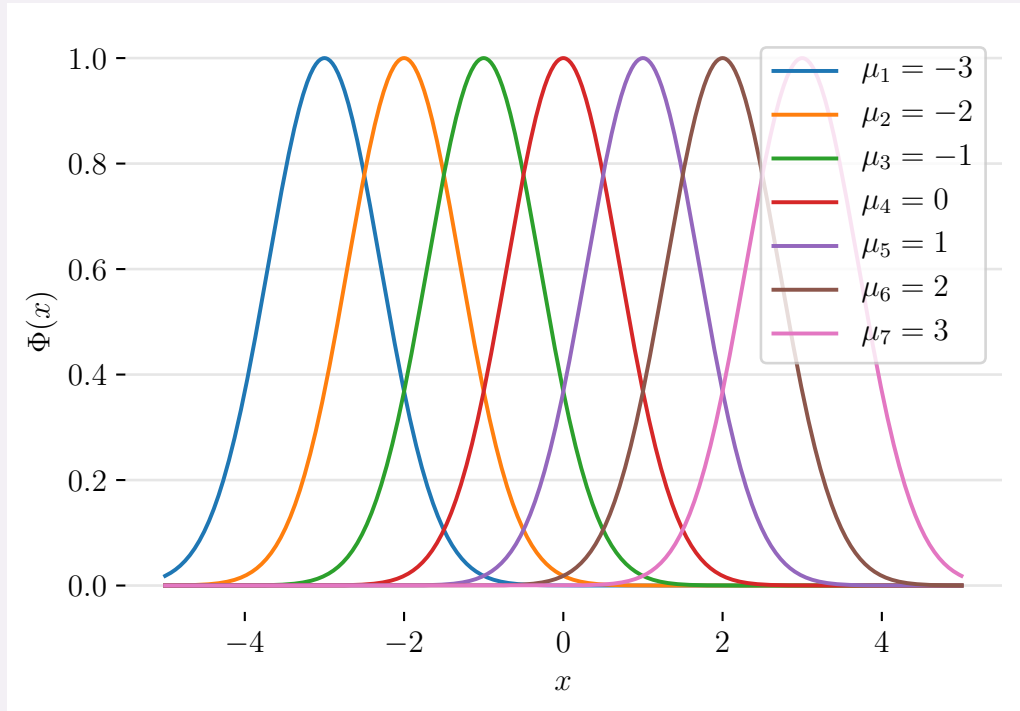
# EXAMPLE - POLYNOMIAL BASIS FUNCTION



(a) Polynomial basis out to degree 4.

$$y = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 = \sum_{j=0} \theta_j x^j$$

# EXAMPLE - GAUSSIAN BASIS FUNCTION



(a) Examples of Gaussian-type radial basis functions.

$$y = \theta_0 + \theta_1 e^{-\frac{(x-\mu_1)^2}{2s^2}} + \cdots + \theta_7 e^{-\frac{(x-\mu_7)^2}{2s^2}}$$

# EXERCISE

HTTPS://COLAB.RESEARCH.GOOGLE.COM/DRIVE/1V1FN_VBCFAXXAPHZW-WUGJRRGXGLFVFQ?USP=SHARING

# QUESTIONS?

# HW1 (DUE 4/20)

HW 1

Due Apr 20 at 11:59pm  |  60 pts