# CS 383: Machine Learning

Prof Adam Poliak

Fall 2024

09/19/2024
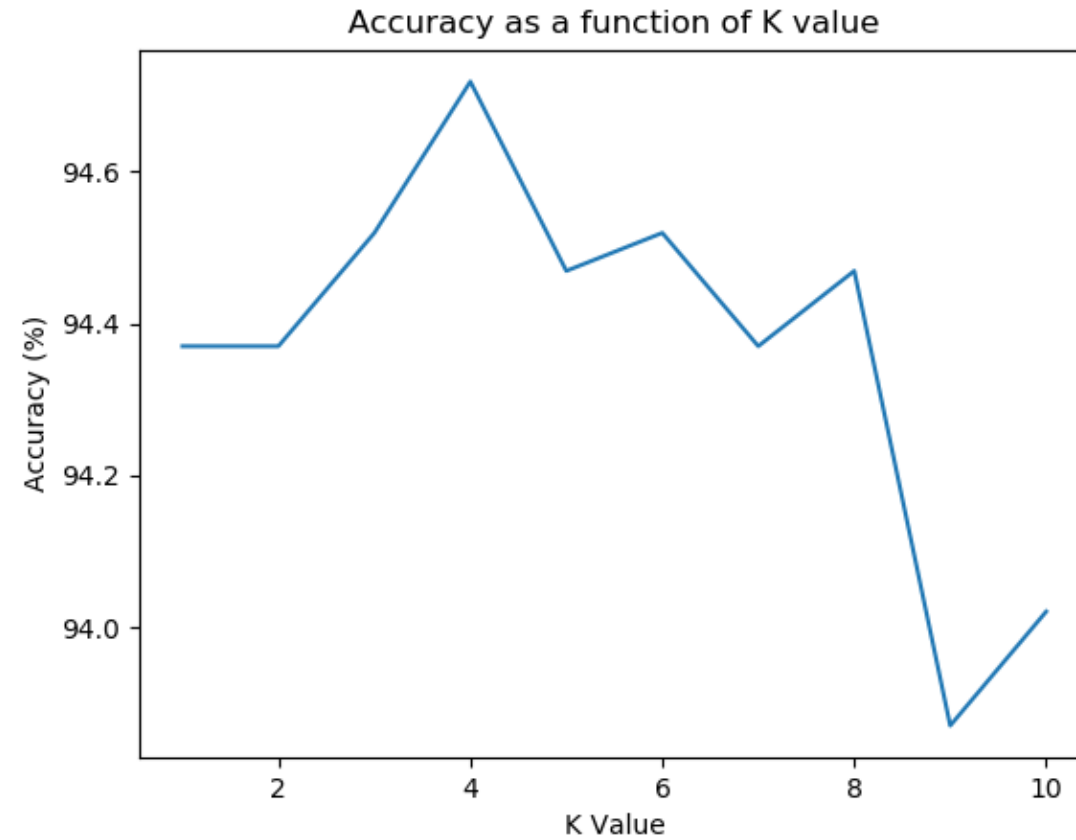
Lecture 07

# Announcements

HW02 is due Sunday night

- **Reading quiz: Thursday**
  - Duame 7.6 (2+ pages)
  - ISL 59-63 (4+ pages)

- Midterm 1: Thursday October 3rd

# My accuracy



Accuracy as a function of K value

# Speeding up K-NN

- Runtime: exercise!

- Don't need to sort all distances – for small $K$, we can find the top $K$ neighbors in linear time

- Save matrix of pair-wise distances across $K$

- Use less of the training data

- Put each training example in a "zone" or "cluster". For each test example, identify cluster and only consider neighbors within that cluster

# Outline

Reading quiz #3

Simple linear regression

SGD (Stochastic Gradient Descent)

Normal equations solution

# Outline

Reading quiz #3

**Simple linear regression**

SGD (Stochastic Gradient Descent)

Normal equations solution

# Goals of Inference

1) Which of the features/explanatory variables/predictors (x) are associated with the response variable (y)?

2) What is the relationship between x and y?

3) Is a linear model enough?

4) Can we predict y given a new x?
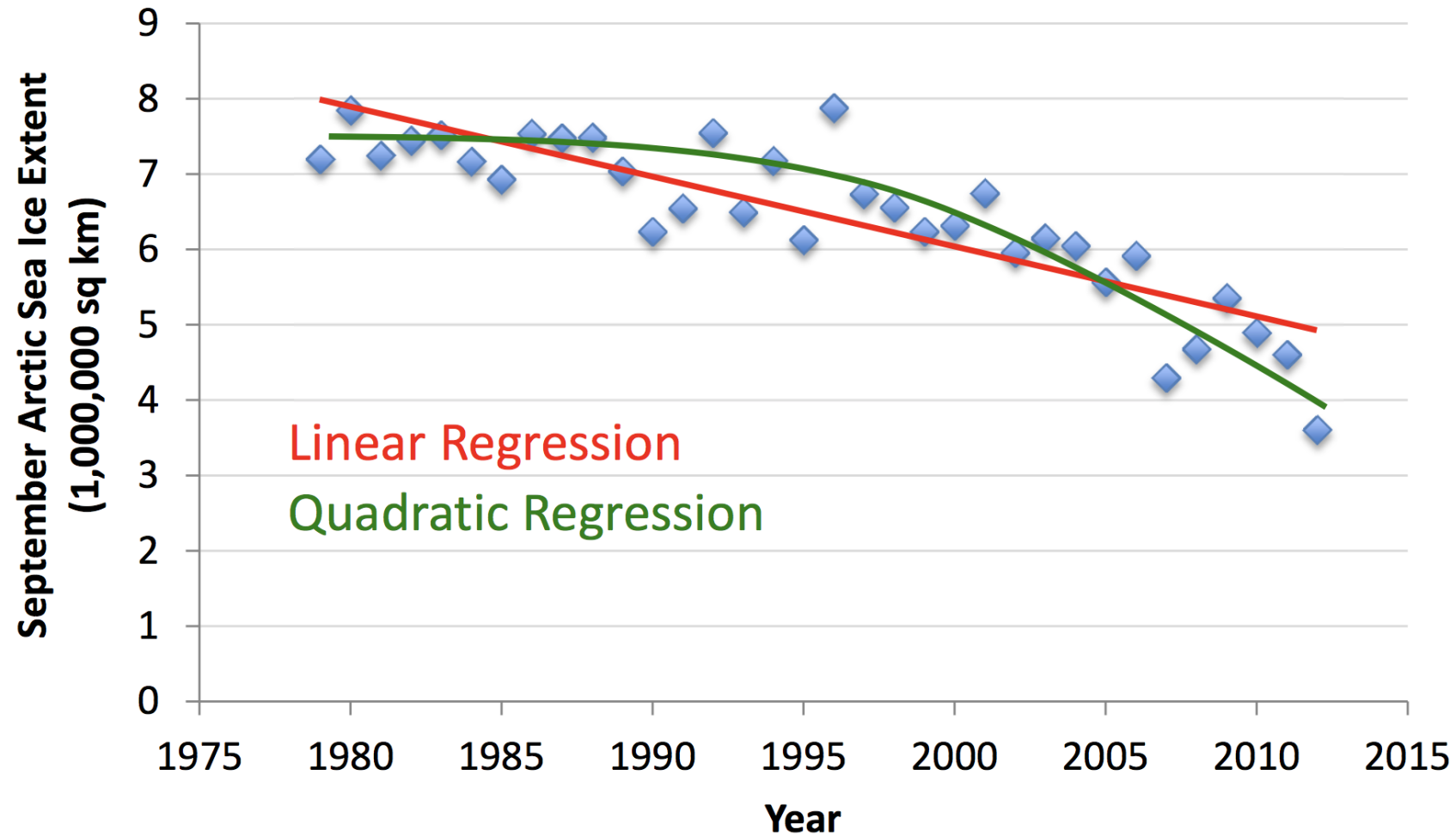
# Linear Regression so far

Output (*y*) is continuous, not a discrete label

Learned model: *linear function* mapping input to output (a *weight* for each feature + *bias*)
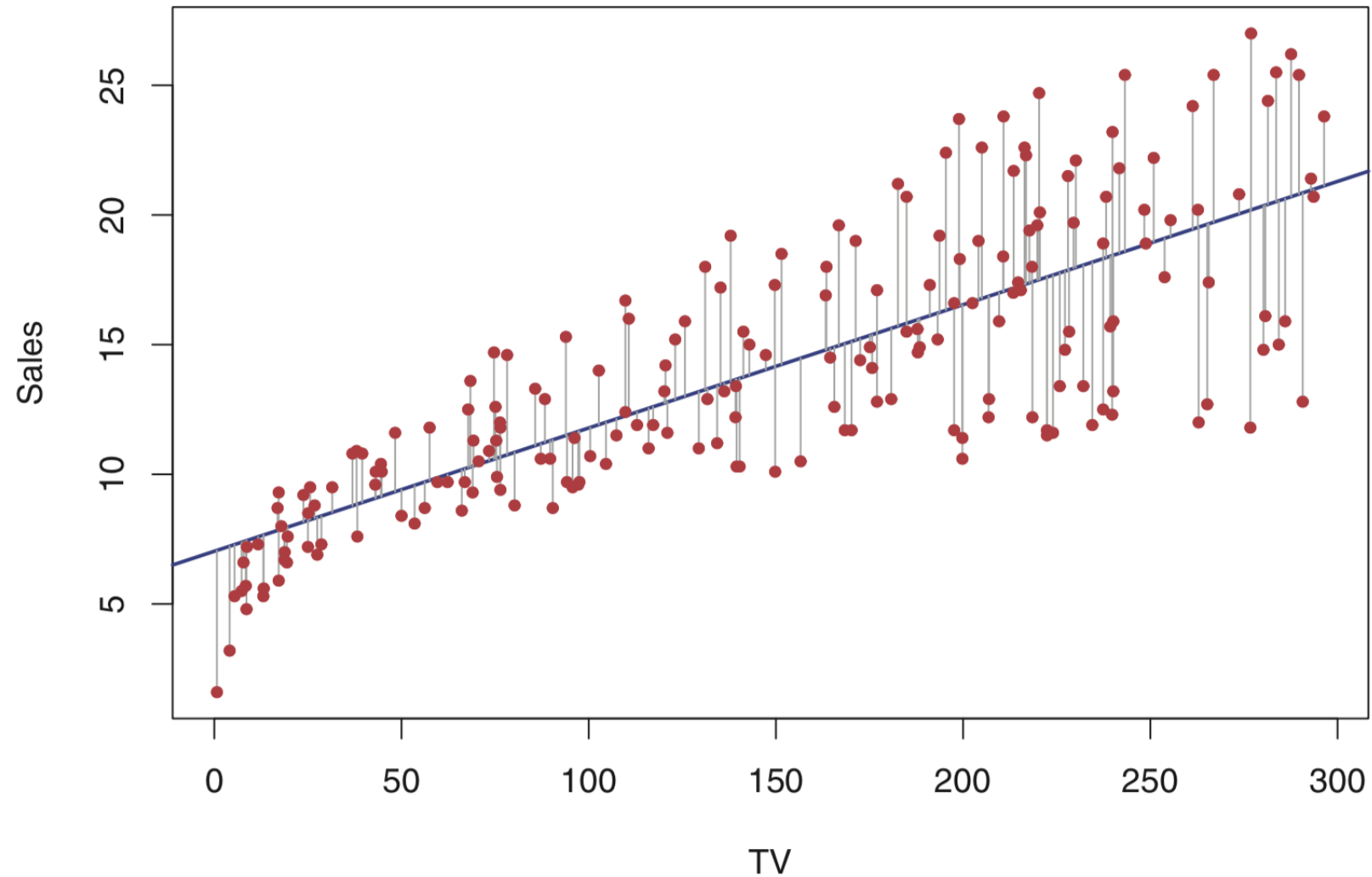
Goal: minimize the *RSS* (residual sum of square) or *SSE* (sum of squared errors)

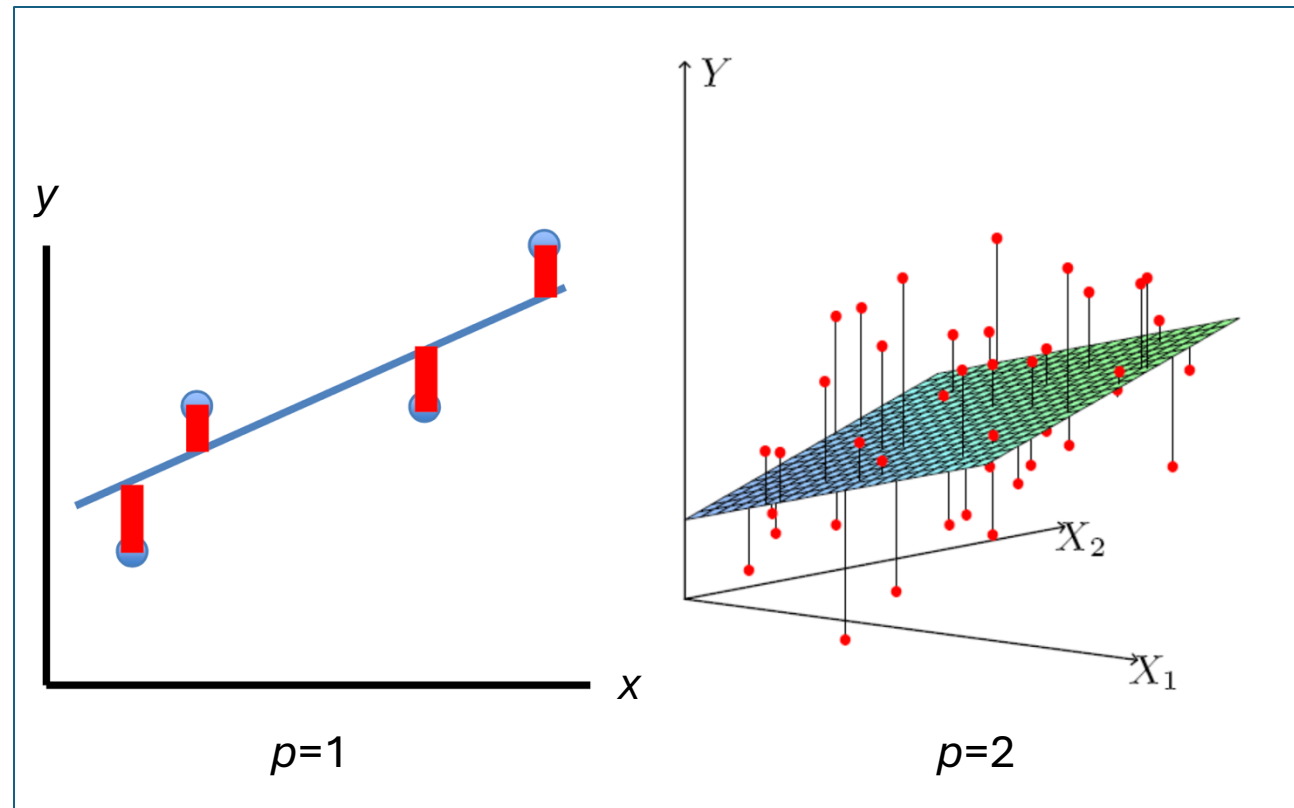$$RSS = \sum_{i}^{n} (y_i - \widehat{y_i})^2$$

# Regression Example



September Arctic Sea Ice Extent (1,000,000 sq km) vs Year

Linear Regression
Quadratic Regression

# Example: predict sales from TV advertising budget

ISL: Figure 3.1

# Cost Function: sum of squared errors



p=1

p=2

CS383 - Lecture 07 - ML

# Outline

Reading quiz #3
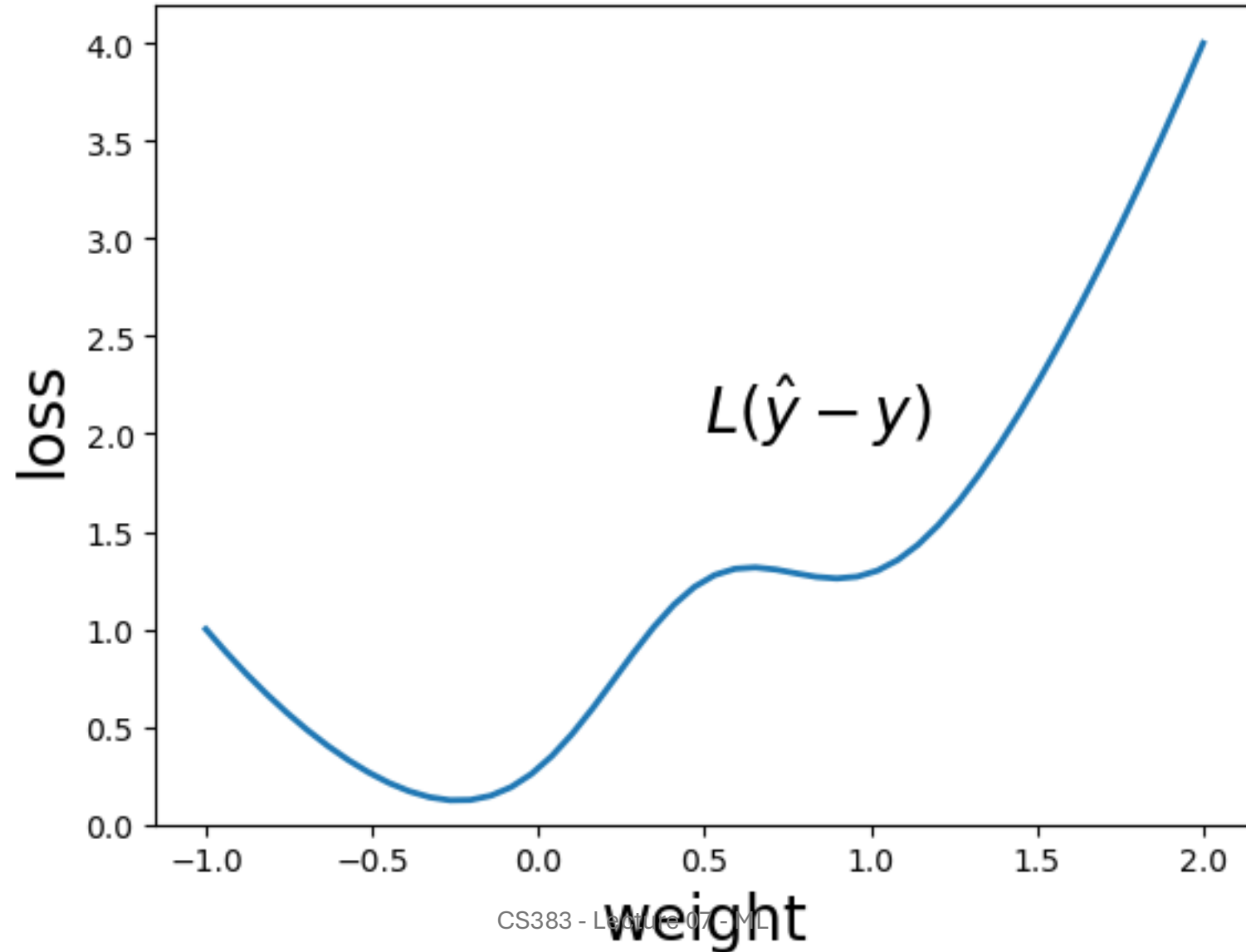
Simple linear regression

**SGD (Stochastic Gradient Descent)**

Normal equations solution

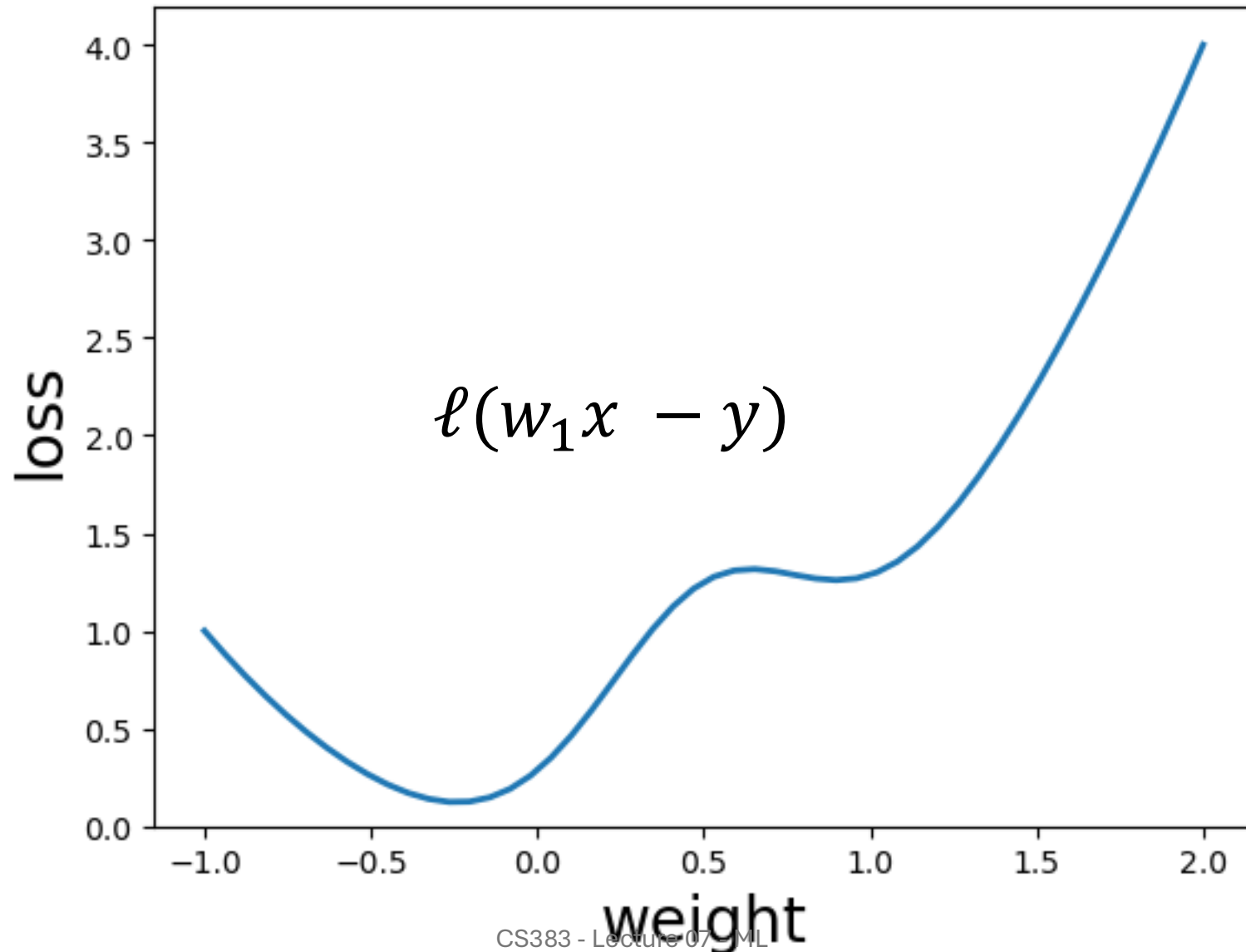# Process Learning Weights

1. Randomly initialize weights

2. Make predictions $\hat{y}$

3. Compute loss function - quantify how close $\hat{y}$ *and* $y$ are

4. Update weights accordingly based on the loss function
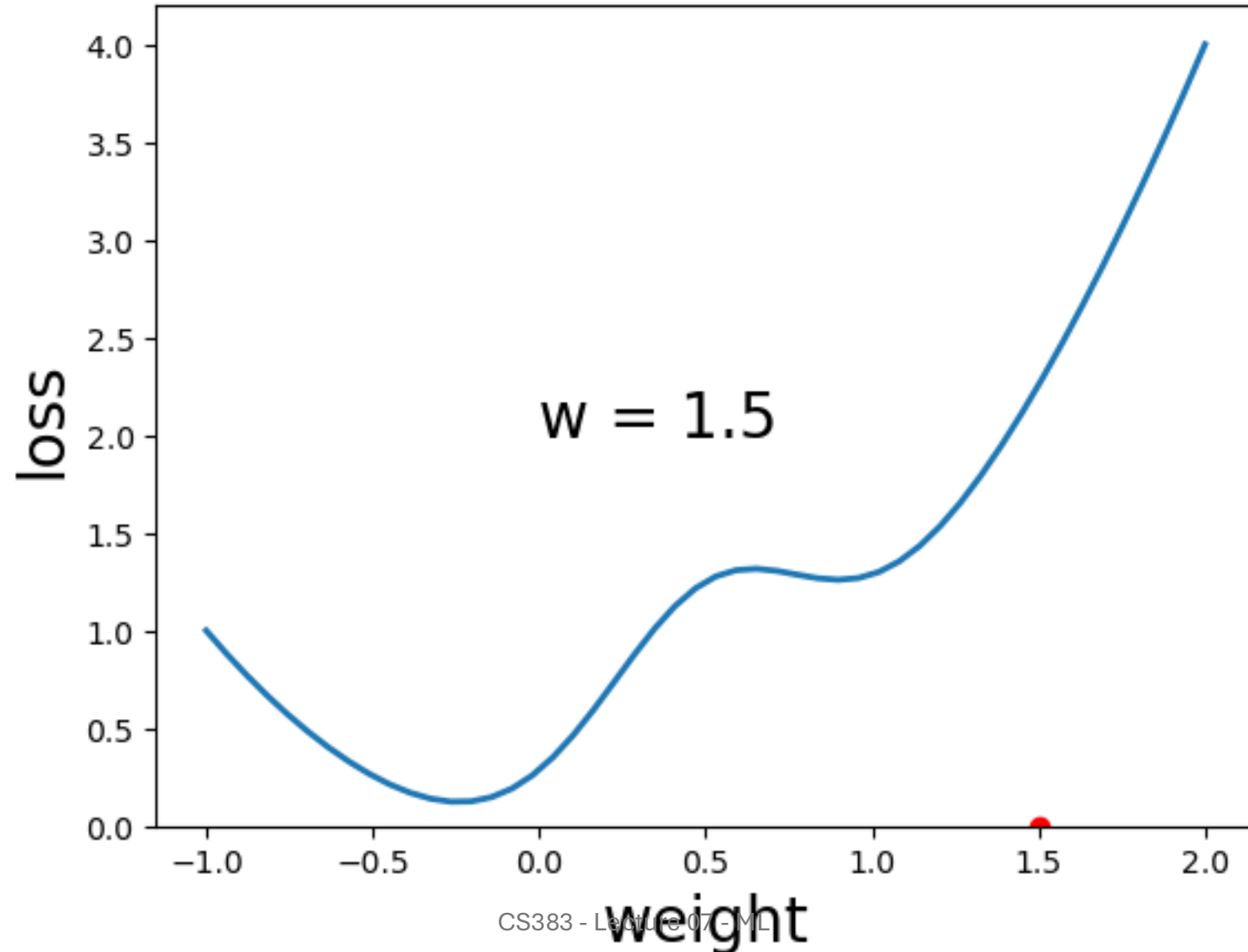                                              aka Optimization

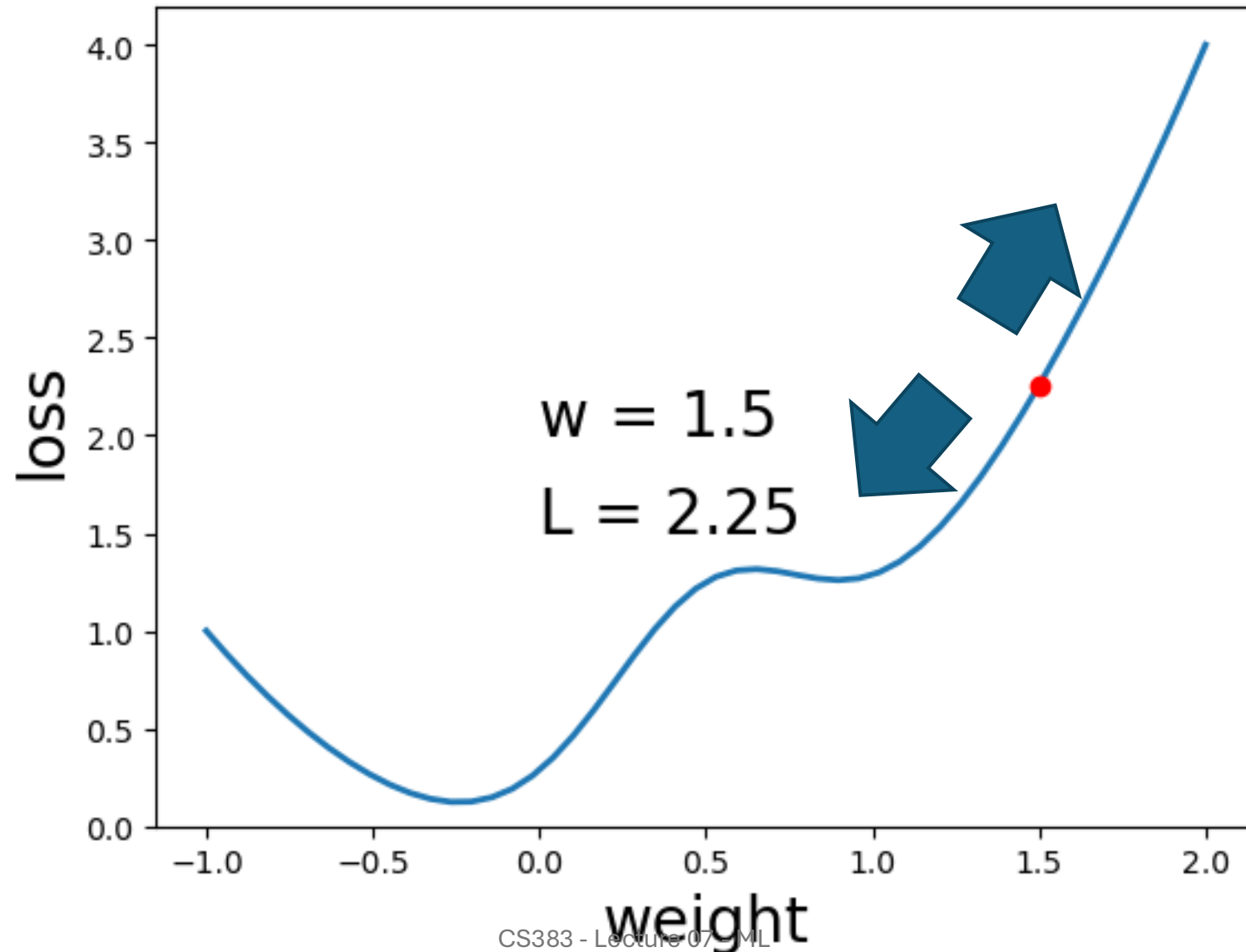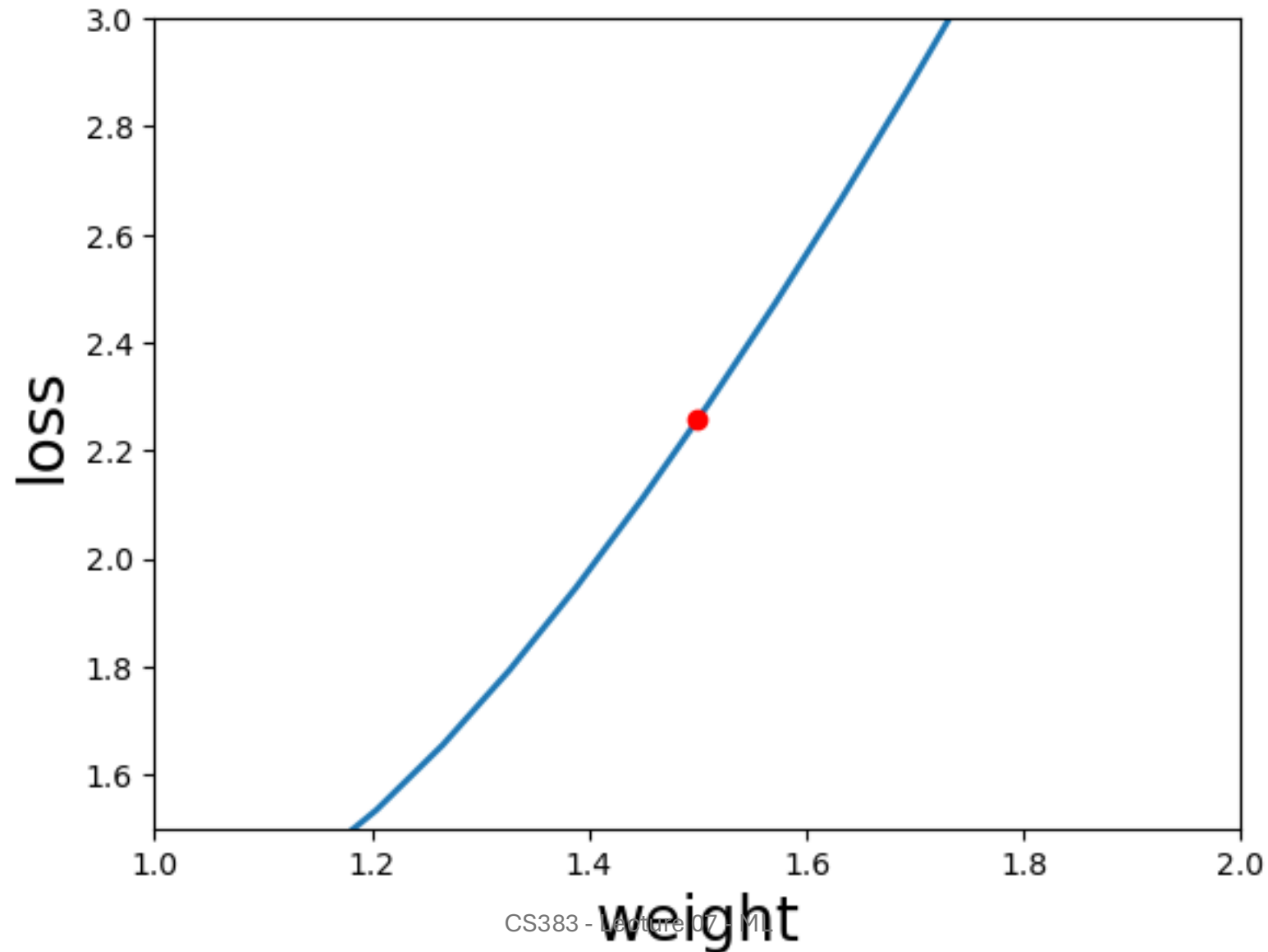5. Repeat 2-4

# Find weights that minimize the loss



CS383 - Lecture 07 - ML

# Find weights that minimize the loss



$$\ell(w_1 x - y)$$

# Find weights that minimize the loss



w = 1.5

# Find weights that minimize the loss



w = 1.5

L = 2.25

# Find weights that minimize the loss

CS383 - Lecture 07 - ML

# Find weights that minimize the loss

How far down should we move the weight?
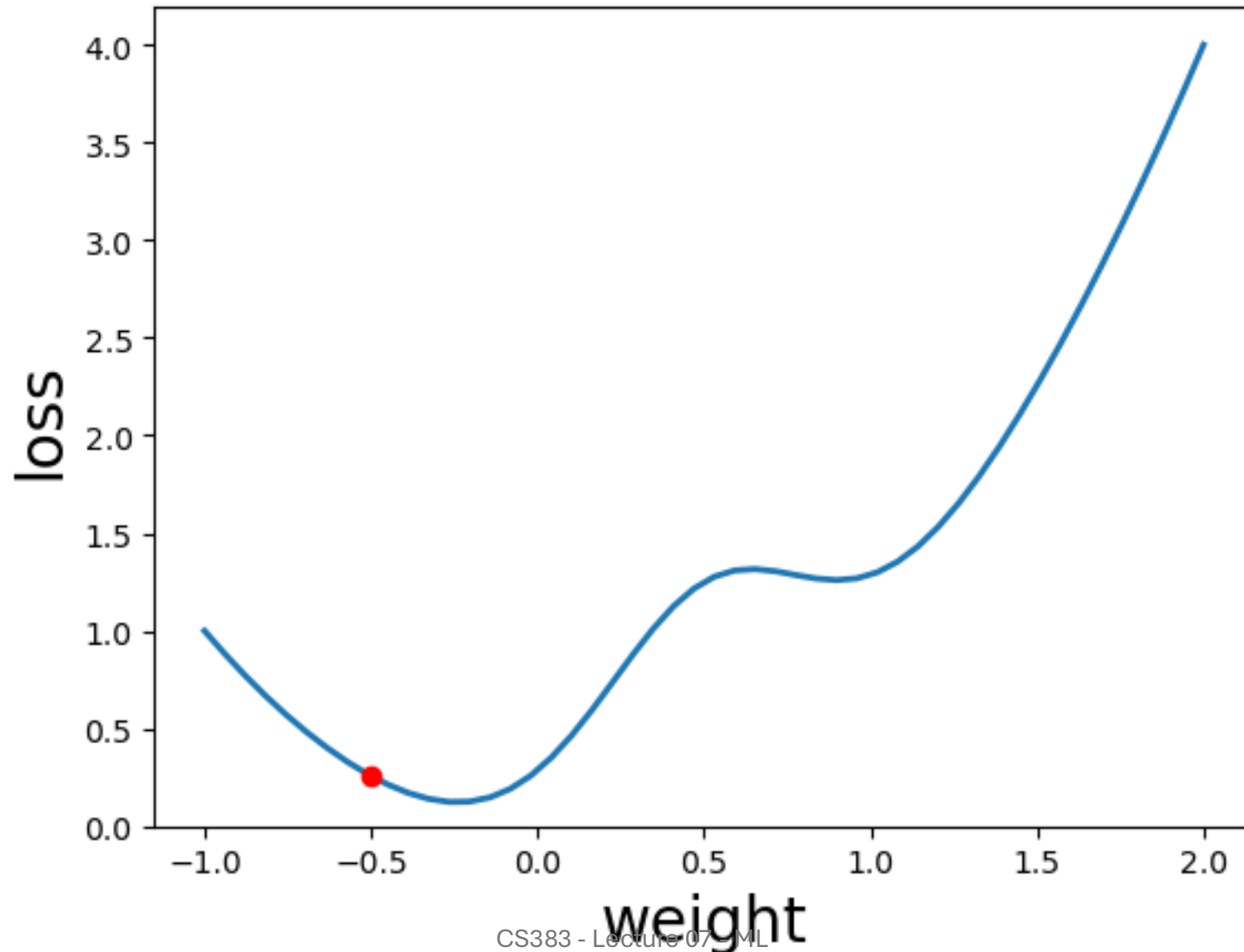
This is called the step-size or learning-rate

CS383 - Lecture 07 - ML

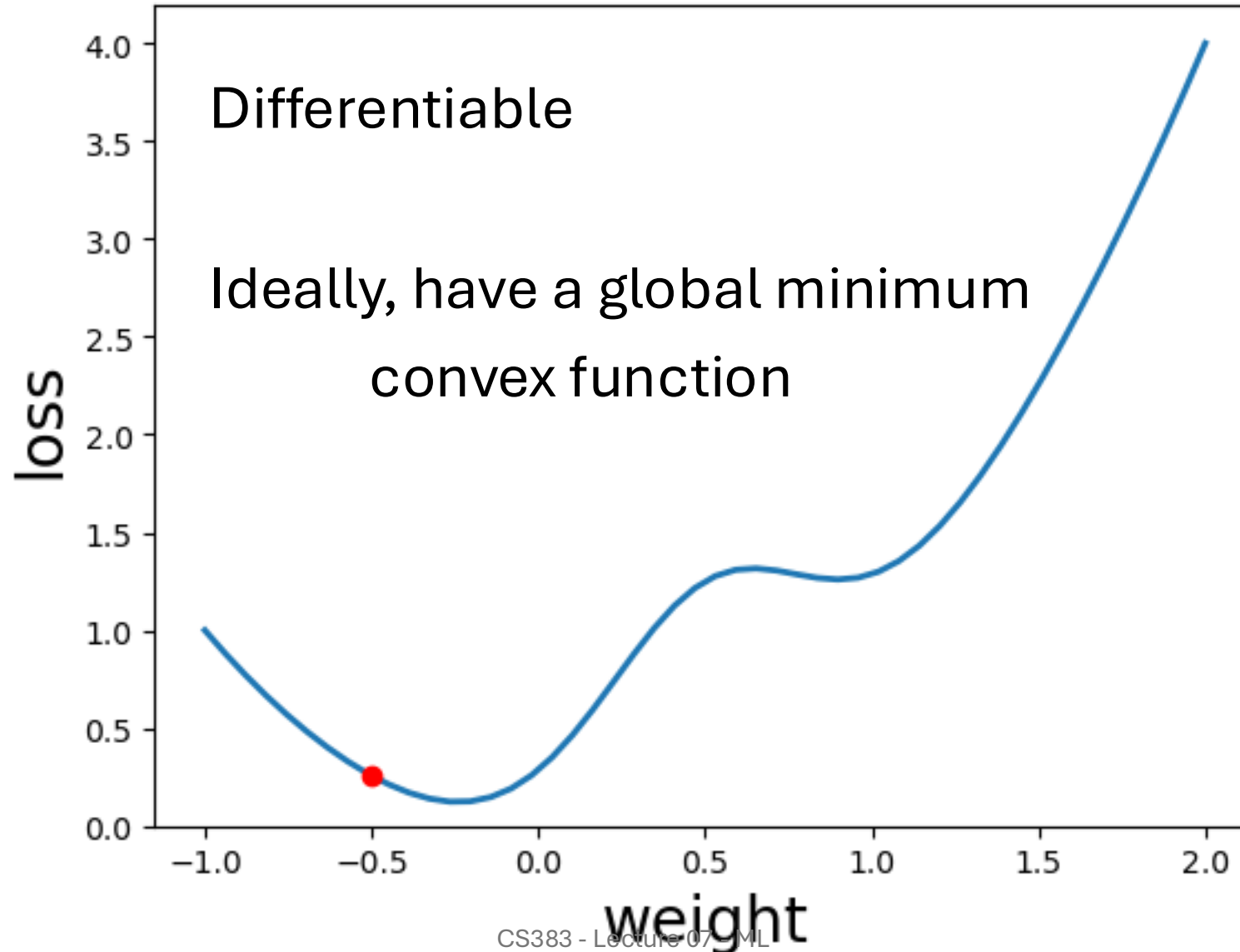# Find weights that minimize the loss
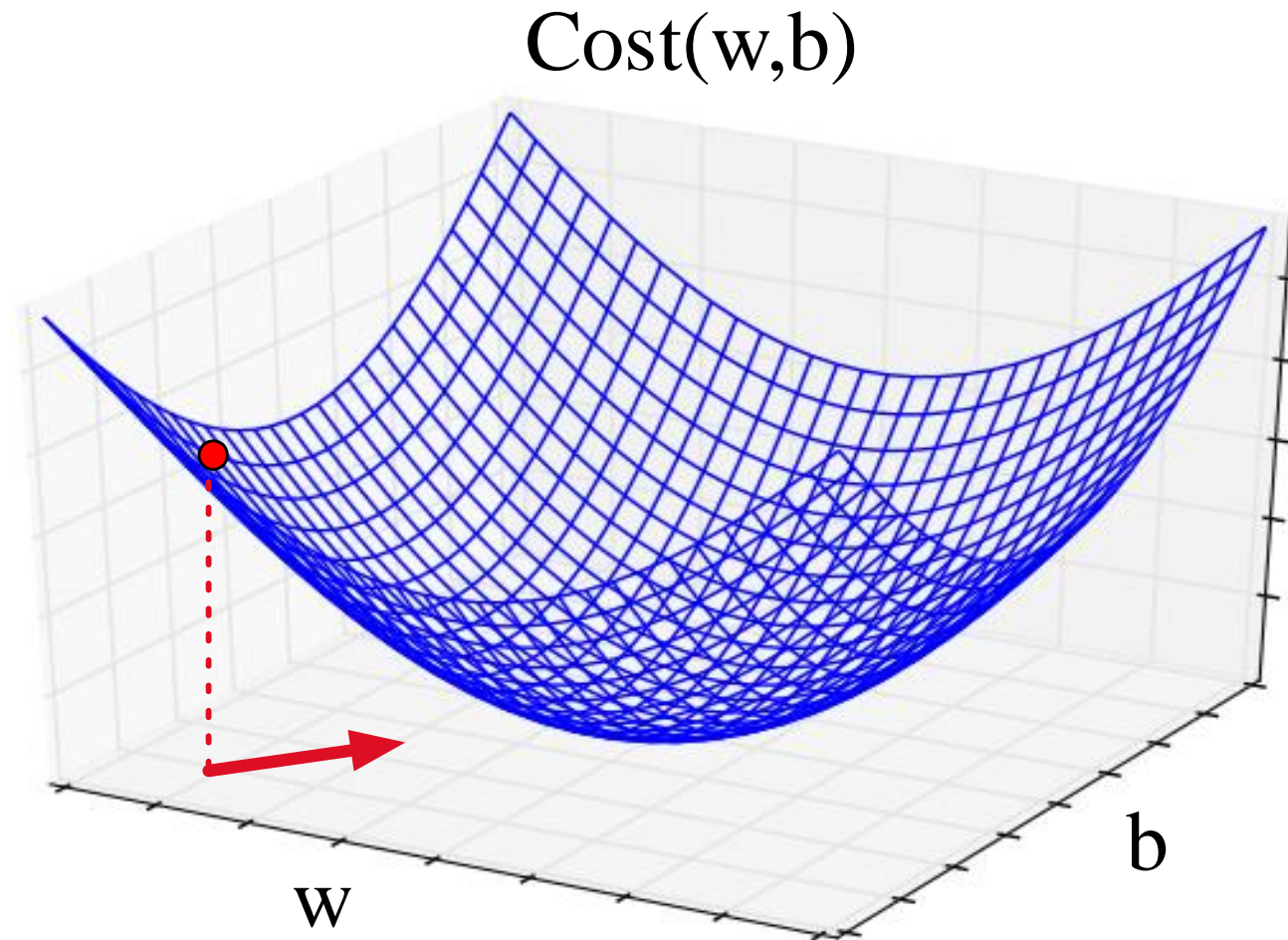
# Find weights that minimize the loss

# How to update the weights

1. Find the direction of the derivative of the loss function

   aka gradient of the loss

2. Move the weight in that direction

3. Then make a prediction on a new $x_{\{i\}}, y_{\{i\}}$ pair, and repeat 1 and 2

4. Repeat this for every example in our training set

# Loss function properties



Differentiable

Ideally, have a global minimum

convex function

# Moving to 2 weights



$Cost(w,b)$

w

b

# Computing the gradient of $\mathcal{L}$ partial derivatives

$$\nabla_{\vec{w}} \mathcal{L} = \begin{bmatrix} \dfrac{d\mathcal{L}}{dw_0} \\ \dfrac{d\mathcal{L}}{dw_1} \\ \dots \\ \dots \\ \dfrac{d\mathcal{L}}{dw_p} \end{bmatrix}$$

What can we do after we computed the gradients?

# Updating weights based on gradients

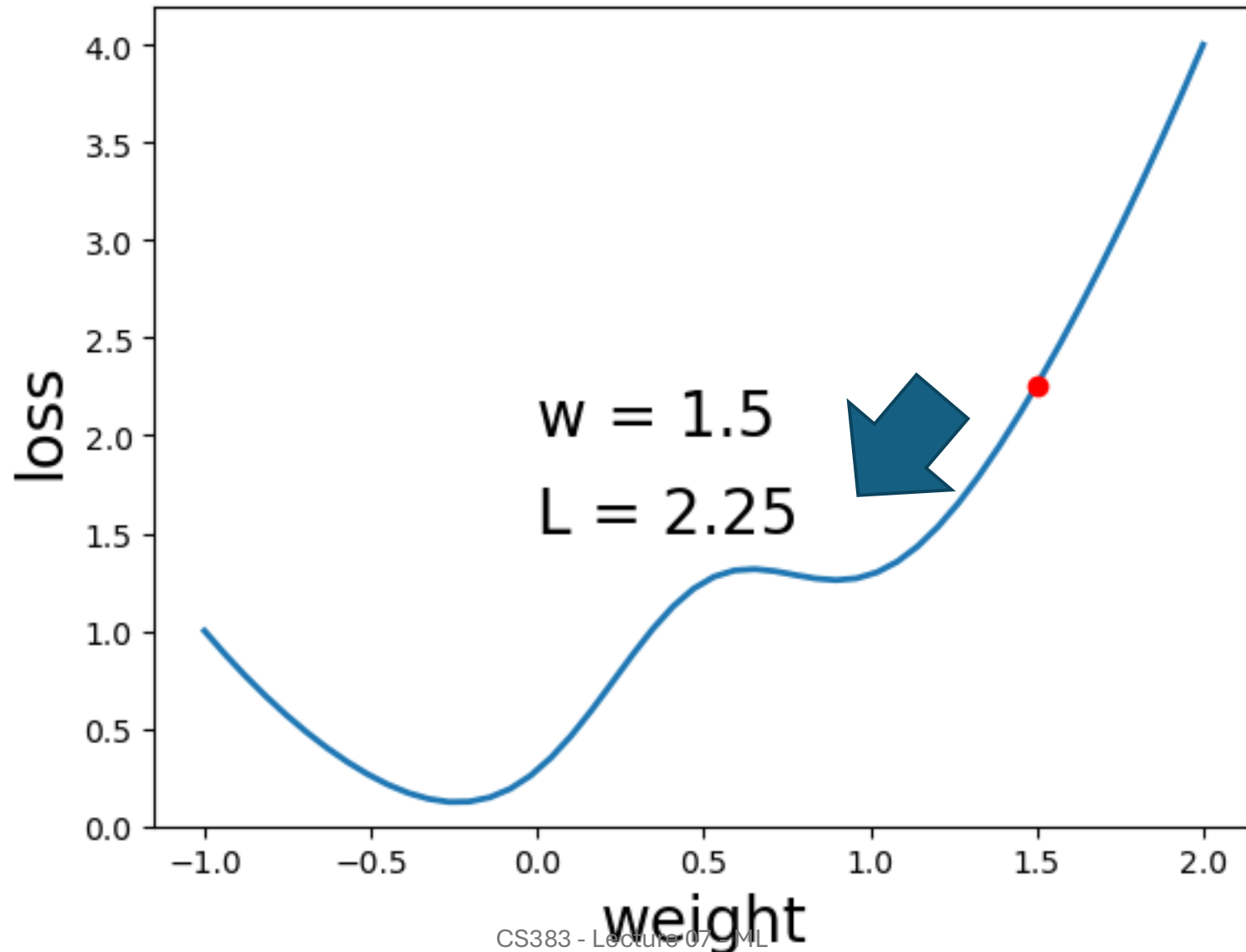$$\Delta_{\vec{w}}\mathcal{L} = \eta\nabla_{\vec{w}}\mathcal{L}\left(\vec{w}\right)$$

Update each individual weight:

$$w_i \leftarrow w_i - \eta\frac{d\mathcal{L}\left(\vec{w}\right)}{dw_i}$$

If we want to perform gradient ascent, we …

$$w_i \leftarrow w_i + \eta\frac{d\mathcal{L}\left(\vec{w}\right)}{dw_i}$$

# Find weights that minimize the loss



w = 1.5

L = 2.25

# Updating weights based on gradients

$$\Delta_{\vec{w}}\mathcal{L} = \eta \nabla_{\vec{w}}\mathcal{L}\left(\vec{w}\right)$$

Update each individual weight:

$$w_i \leftarrow w_i - \eta\frac{d\mathcal{L}\left(\vec{w}\right)}{dw_i}$$

If we want to perform gradient ascent, we …

$$w_i \leftarrow w_i + \eta\frac{d\mathcal{L}\left(\vec{w}\right)}{dw_i}$$

Step size

# Gradient Descent

1. Randomly initialize $\vec{w}$

2. For every $\{x_i, y_i\}$ pair in our training set:

      Compute the gradient of the loss

$$\frac{d\mathcal{L}(\vec{w})}{dw_i}$$

      Update each weights based on the gradients

$$w_i \leftarrow w_i - \eta \frac{d\mathcal{L}(\boldsymbol{w})}{dw_i}$$

3. Repeat 2 until convergance (or max epochs)
4. return $\beta_i$

# Stochastic Gradient Descent

1. Randomly initialize $\vec{w}$

2. Randomly choose a $\{x_i, y_i\}$ pair in our training set without replacement until all pairs are used:

      Compute the gradient of the loss

$$\frac{d\mathcal{L}\left(\vec{w}\right)}{dw_i}$$

      Update each weights based on the gradients

$$w_i \leftarrow w_i - \eta\frac{d\mathcal{L}(\boldsymbol{w})}{dw_i}$$

3. Repeat 2 until convergence (or max epochs)

4. return $\beta_i$

# Outline

Reading quiz #3

Simple linear regression

SGD (Stochastic Gradient Descent)

**Normal equations solution – see handout posted on slack with derivation**

# Pros and Cons

## **Gradient Descent**

- Requires multiple iterations

- Need to choose η

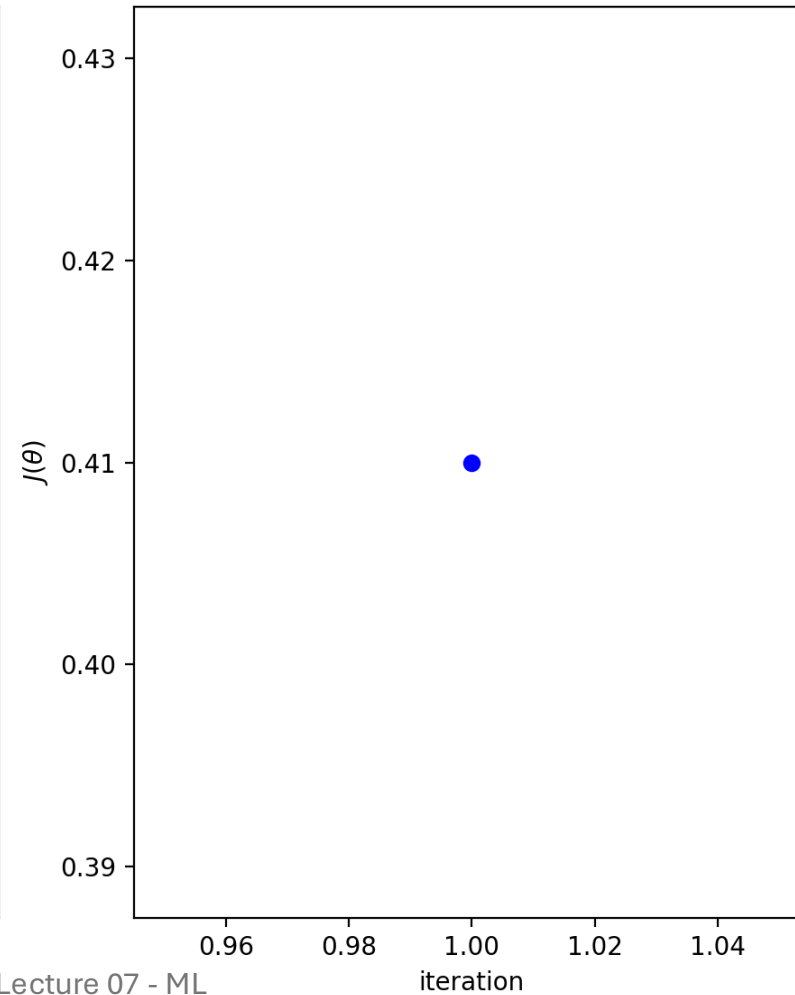- Works well when *n* is large
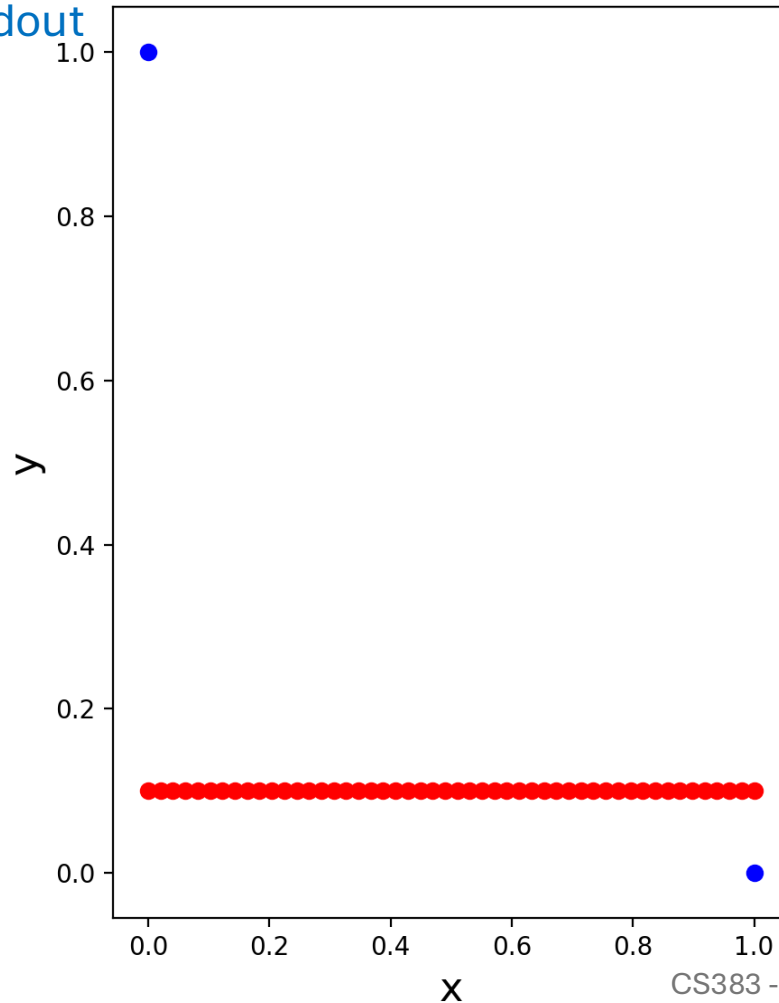
- Can support online learning

## **Normal Equation**

- Non-iterative

- No need to choose η

- Slow if *p* is large

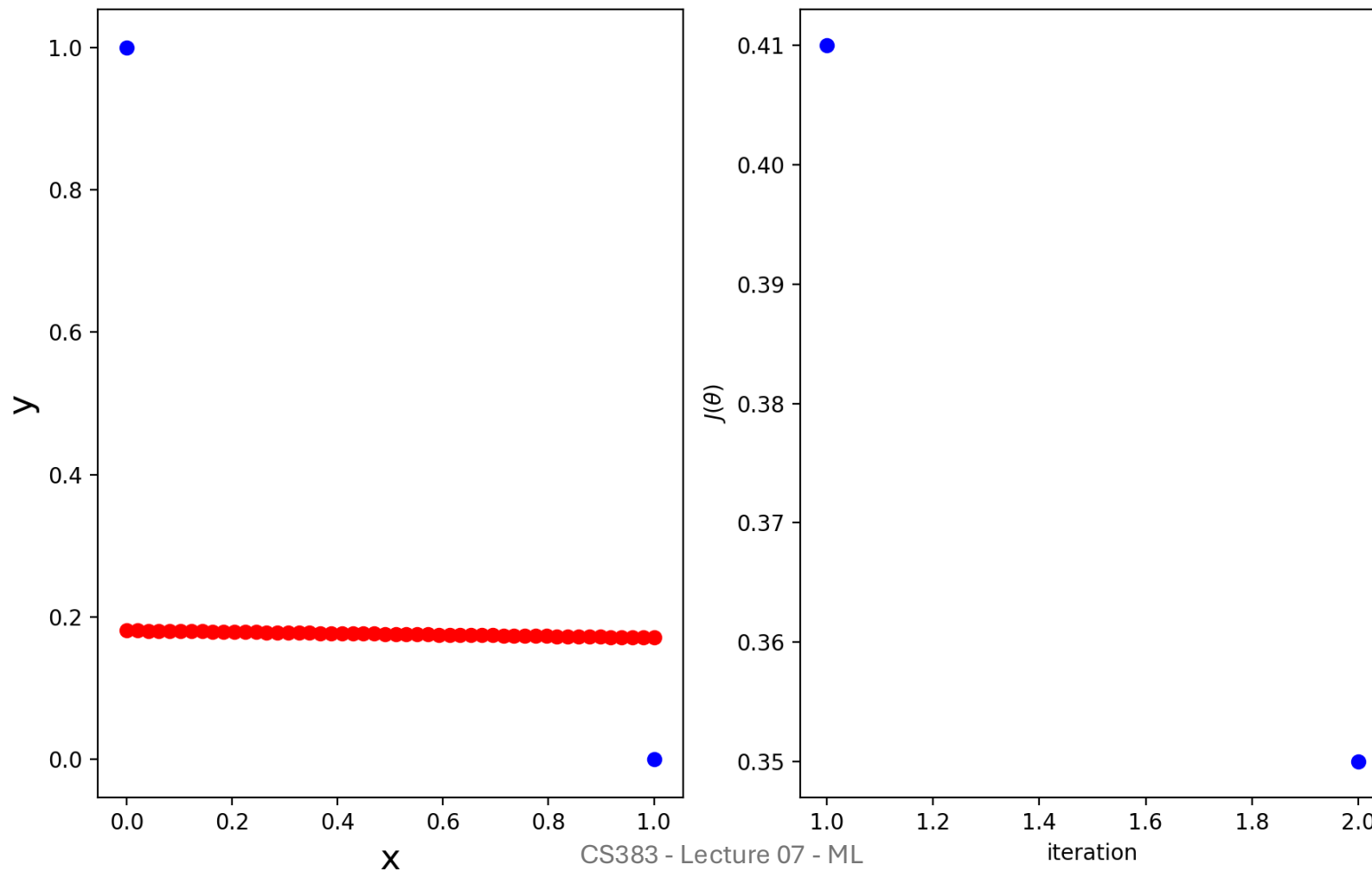  - Matrix inversion is $O(p^3)$

# Toy example, iteration 1

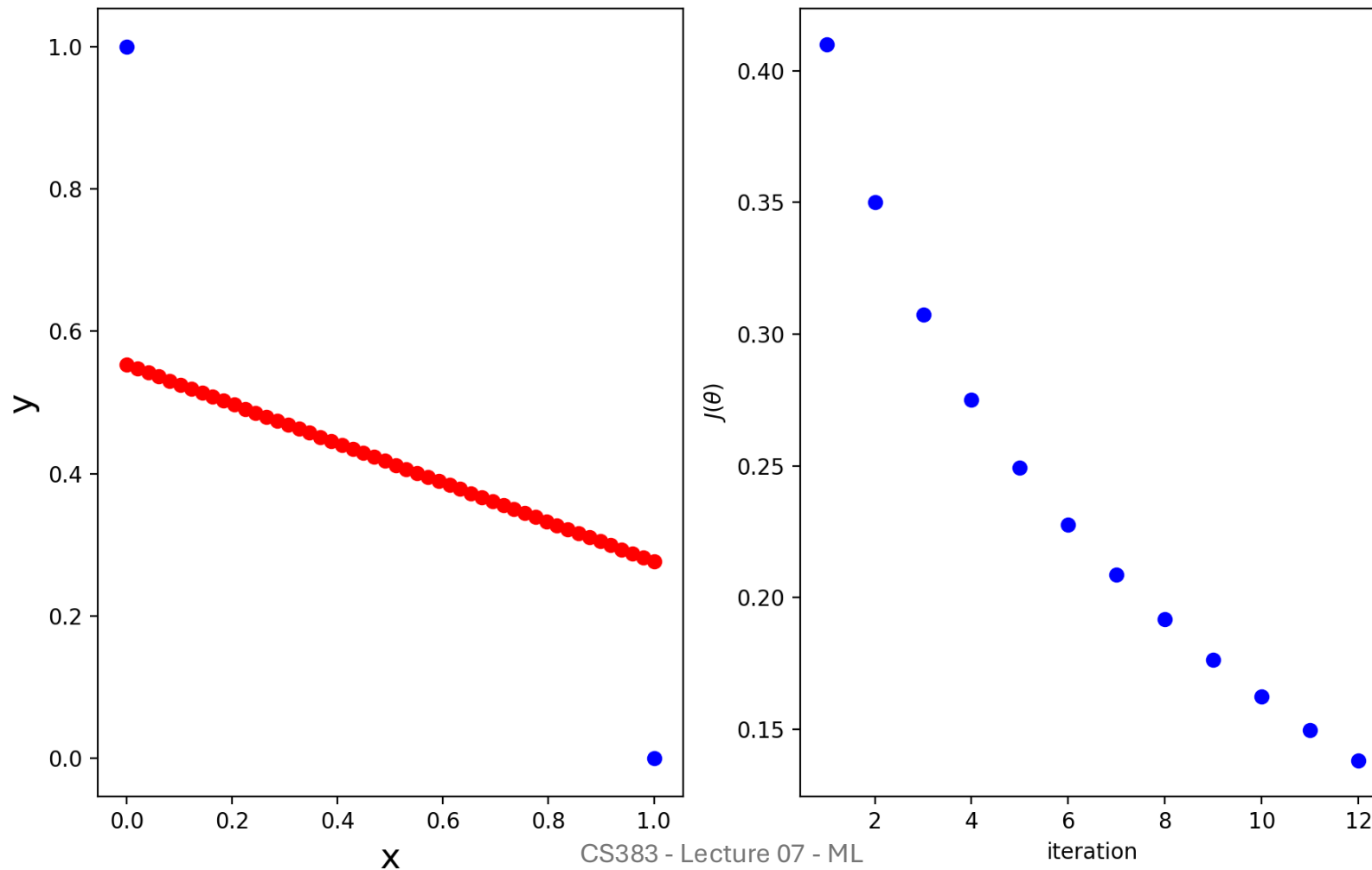This is what you should have obtained in Handout 7!

iteration: 1, cost: 0.410000

CS383 - Lecture 07 - ML

# Toy example, iteration 2

iteration: 2, cost: 0.350001

CS383 - Lecture 07 - ML

# Toy example, iteration 12

iteration: 12, cost: 0.138047

CS383 - Lecture 07 - ML

# Toy example, iteration 40

iteration: 40, cost: 0.014064

# Toy example, iteration 100

iteration: 100, cost: 0.000105