

# CS 260: Foundations of Data Science

Prof. Sara Mathieson

Fall 2023



**HAVERFORD**  
COLLEGE

# Admin

- **Lab 1** grades posted on Moodle
- **Lab 3** due Wednesday night
  - I'll check in with everyone today about Lab 3 during lab
- **Lab 4** posted today

# Outline for September 26

- Recap SGD (stochastic gradient descent)
- Introduction to classification
  - Decision tree models
  - Probabilistic interpretation
- Evaluation Metrics
  - Confusion matrices
  - Precision and recall
  - ROC curves

# Outline for September 26

- Recap SGD (stochastic gradient descent)
- Introduction to classification
  - Decision tree models
  - Probabilistic interpretation
- Evaluation Metrics
  - Confusion matrices
  - Precision and recall
  - ROC curves

# Stochastic Gradient Descent for Linear Regression

Key Idea: take the derivative of **one datapoint** at a time and use that to update  $w$

$$\nabla J = \begin{bmatrix} \frac{\partial J}{\partial w_0} \\ \vdots \\ \frac{\partial J}{\partial w_p} \end{bmatrix}$$

Handout 6  
142

derivative  
wrt  
example

$$J(\vec{w}) = \frac{1}{2} \sum_{i=1}^n (\underbrace{\vec{w} \cdot \vec{x}_i}_{\text{pred}} - \underbrace{y_i}_{\text{truth}})^2$$

derivative is very  
large + unstable

$$\nabla J(\vec{w})_{\vec{x}_i} = (\vec{w} \cdot \vec{x}_i - y_i) \vec{x}_i$$

datapoint      scalar      vector

$$\begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{bmatrix}$$

# Stochastic Gradient Descent for Linear Regression

SGD

Start  
while  
for

with  $\vec{w} = \vec{0}$  (vector of zeros)

(epoch)  
iteration  $t$ :

for  $i = 1, 2, \dots, n$ : (shuffle)

$\vec{w} \leftarrow \vec{w} - \alpha (\vec{w} \cdot \vec{x}_i - y_i) \vec{x}_i$

Step size derivative

check convergence  
if

$|J(\vec{w}^t) - J(\vec{w}^{t+1})| < \epsilon$

$\Rightarrow$  Stop!

$\epsilon$  is very small

not for Lab 3

# Handout 6 (#1, #2)

Handout 6

$$X = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}, Y = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

↑  
fake ones

①  $\vec{x}_1$

$$\vec{w} \leftarrow$$

$$\begin{bmatrix} 0 \\ 0 \end{bmatrix} - 0.1 \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \end{bmatrix} - 1 \right) \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

(0.1 + 0.0)  $y_i$

$$\begin{bmatrix} 0 \\ 0 \end{bmatrix} - 0.1(0 - 1) \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0 \end{bmatrix}$$

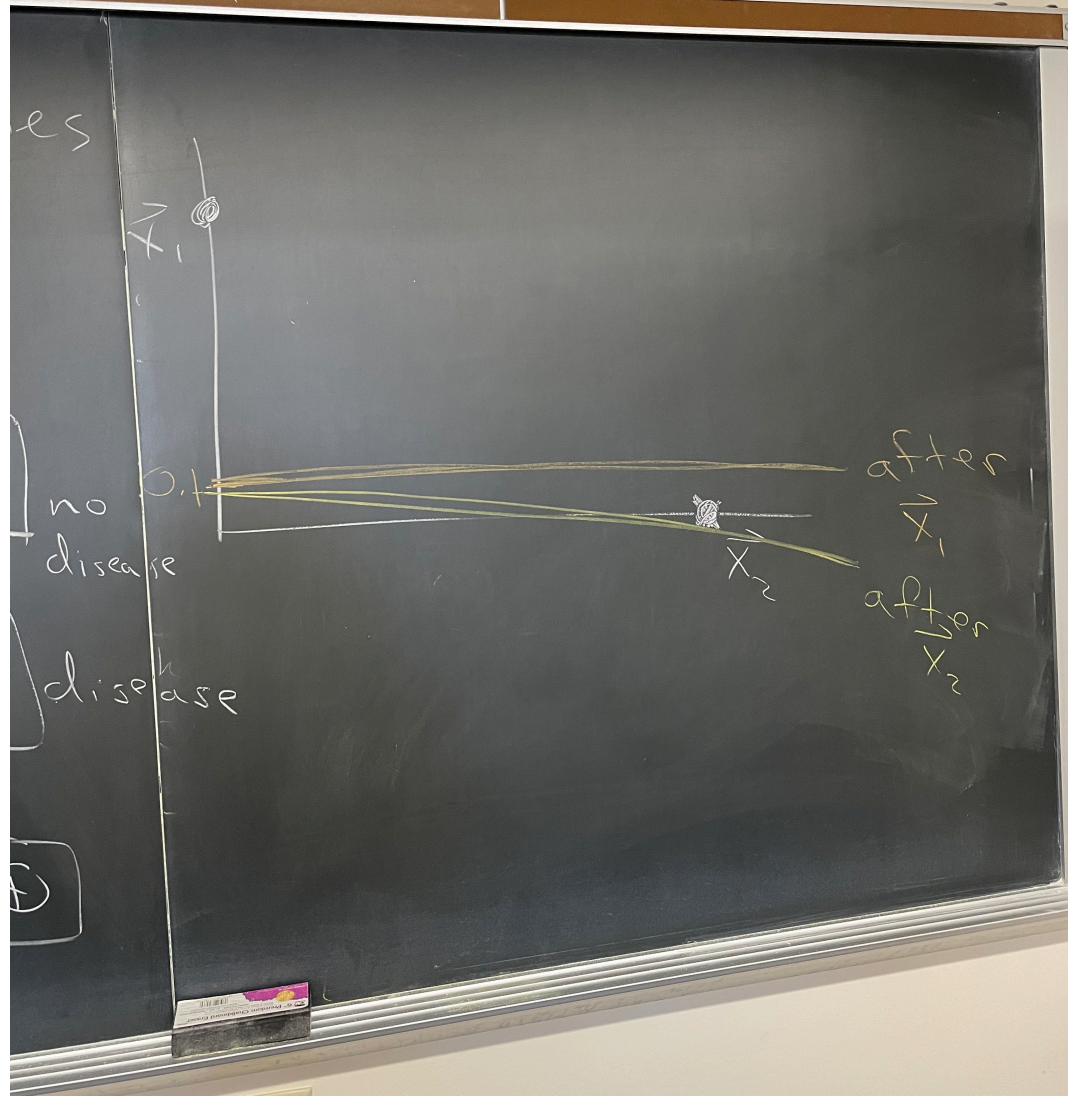
②

$$\vec{w} \leftarrow$$

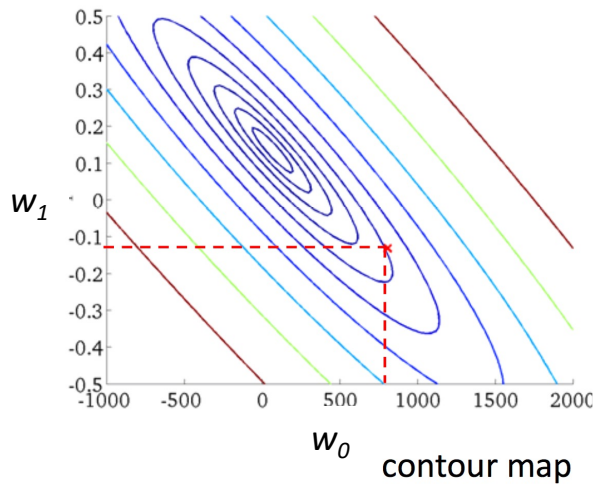
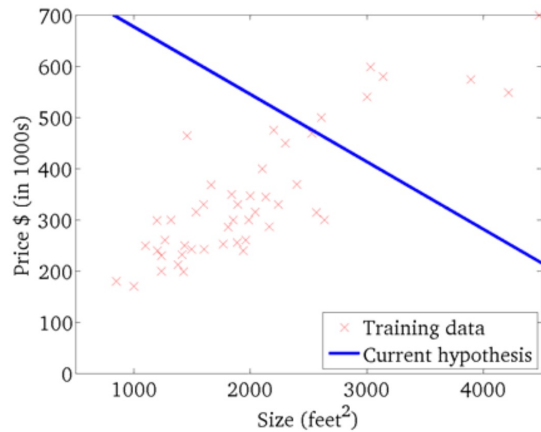
$$\vec{w} \leftarrow \begin{bmatrix} 0.09 \\ -0.01 \end{bmatrix}$$

$$\begin{bmatrix} 0.1 \\ 0 \end{bmatrix} - 0.1 \left( \begin{bmatrix} 0.1 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} - 0 \right) \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

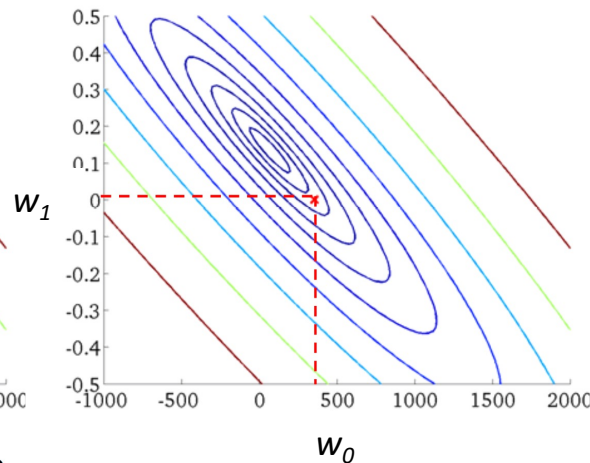
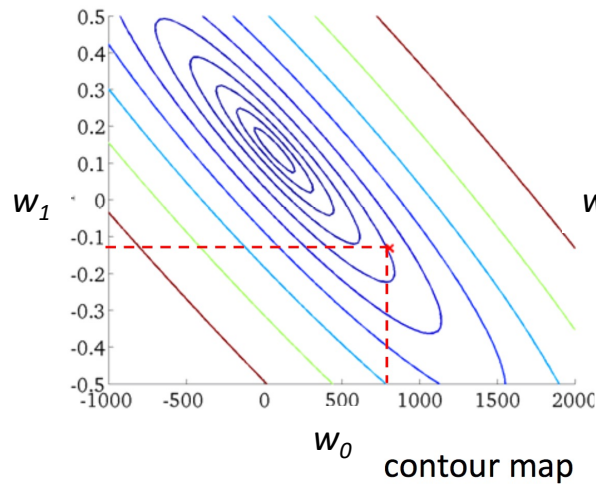
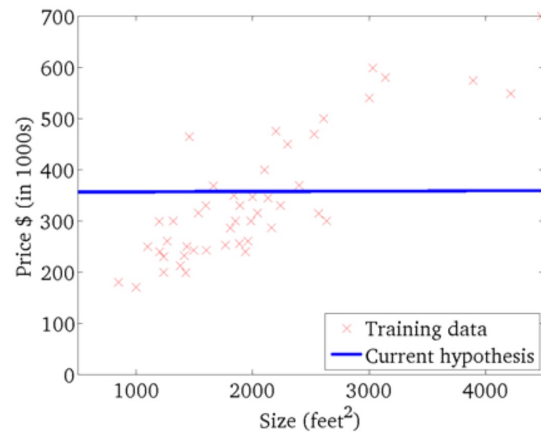
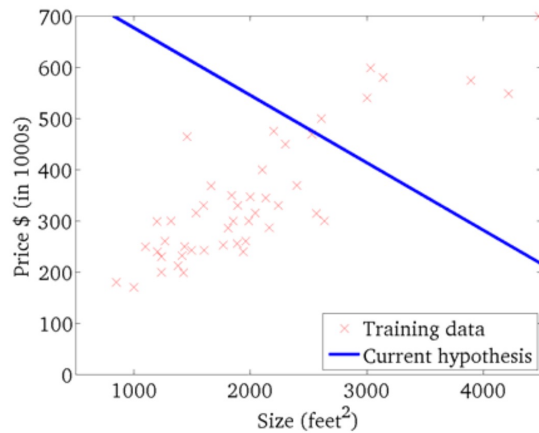
# Handout 6 (#4)



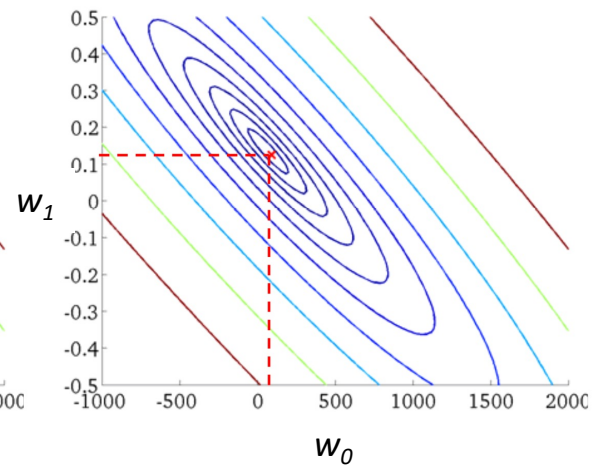
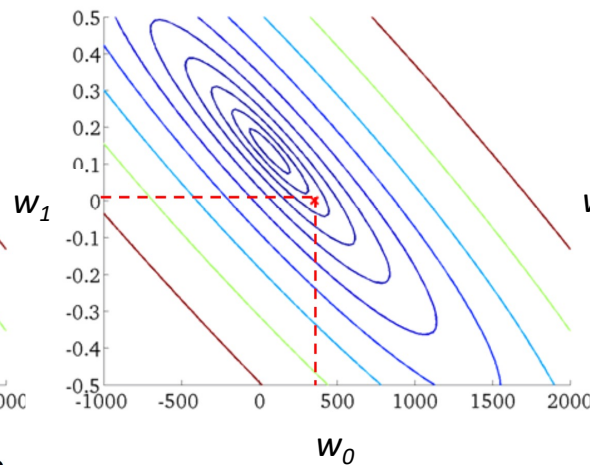
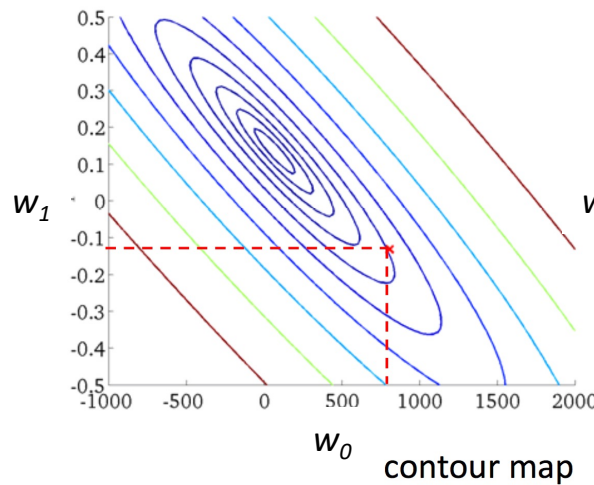
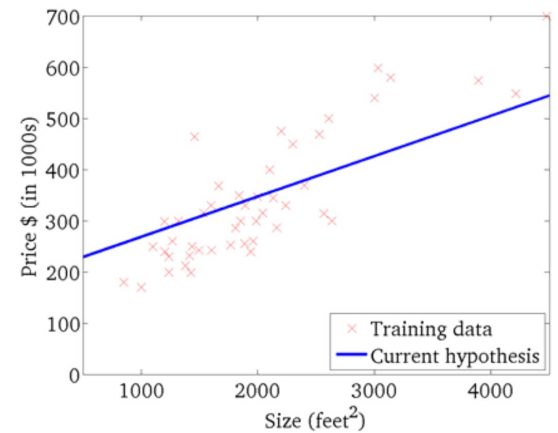
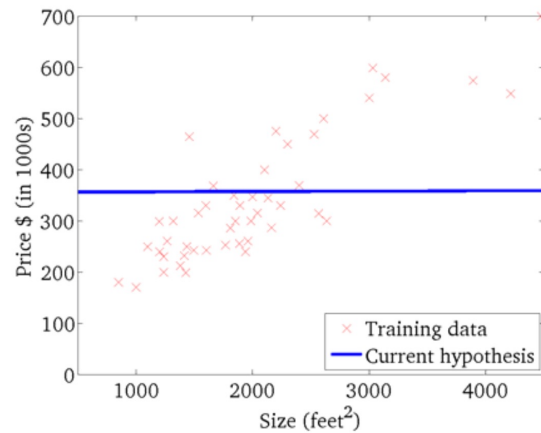
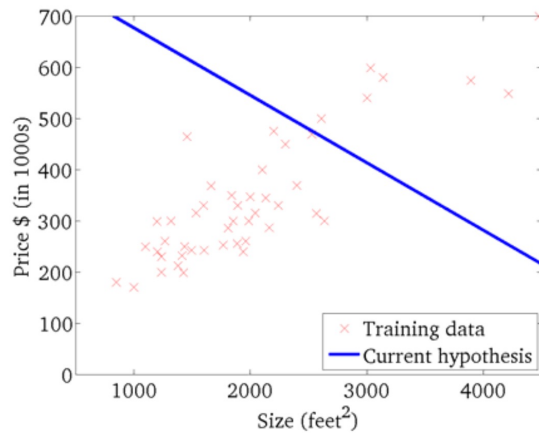
# Linear Model and Cost Function J



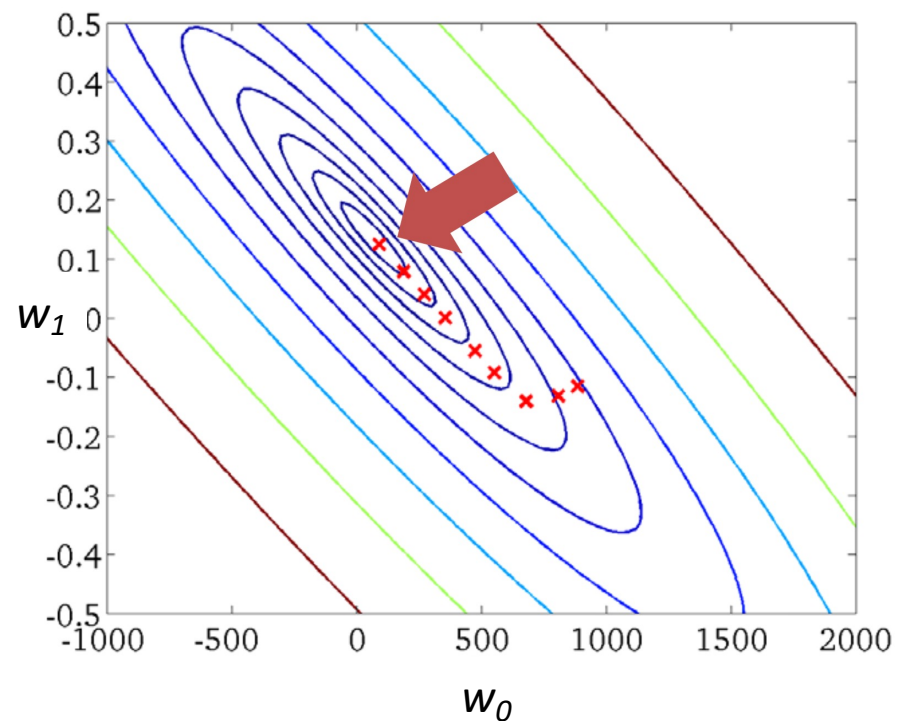
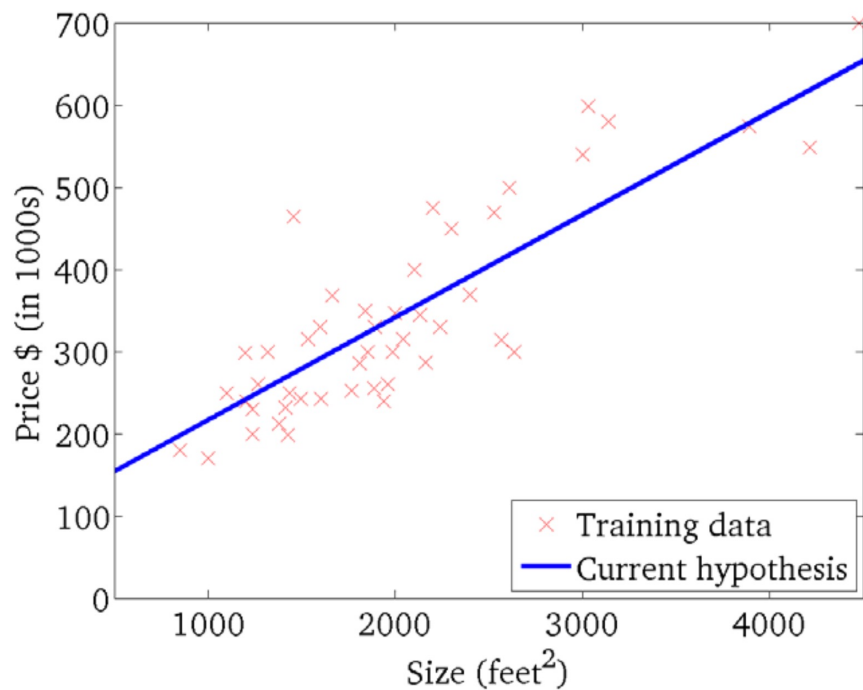
# Linear Model and Cost Function J



# Linear Model and Cost Function J



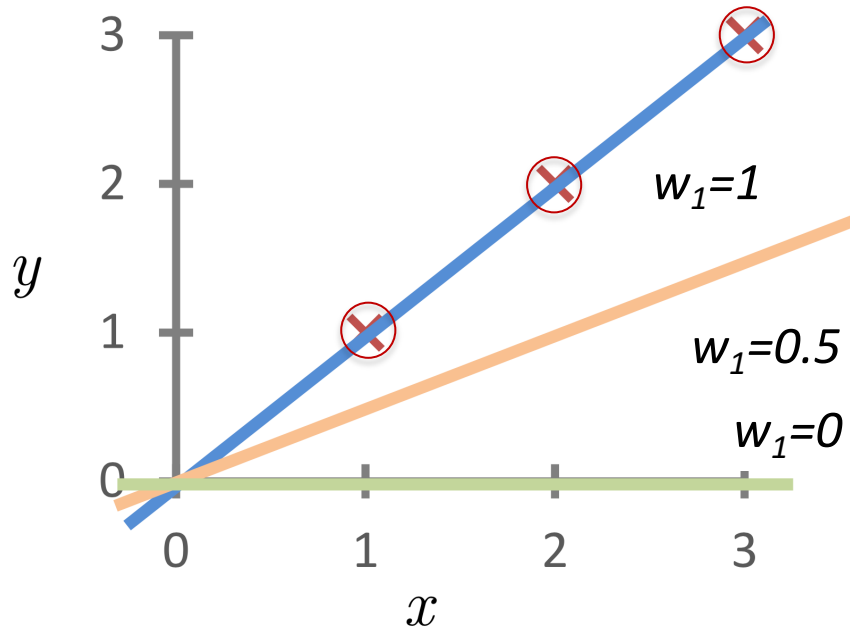
# Gradient Descent: walking toward the minimum



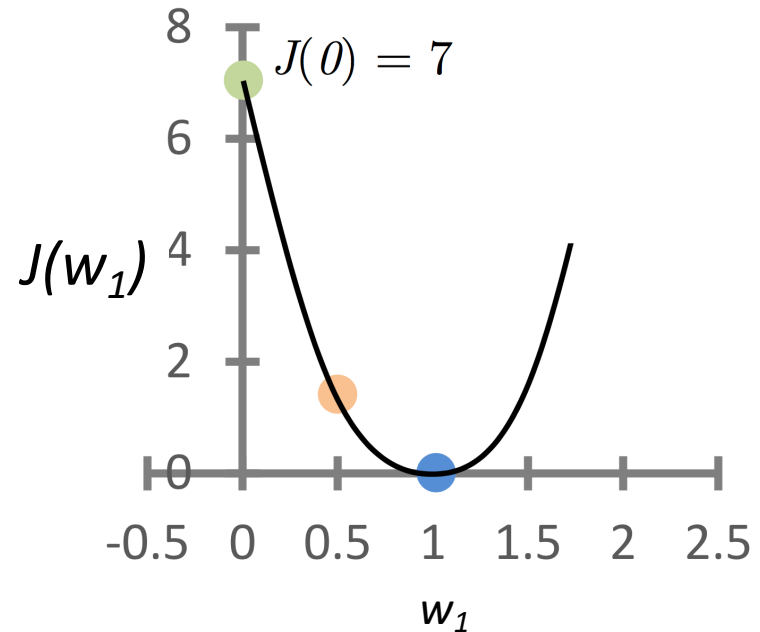
# Cost Function (extra practice)

$$h_w(x) = w_1 x$$

(assume  $w_0=0$  for this example)



$$J(w_1)$$

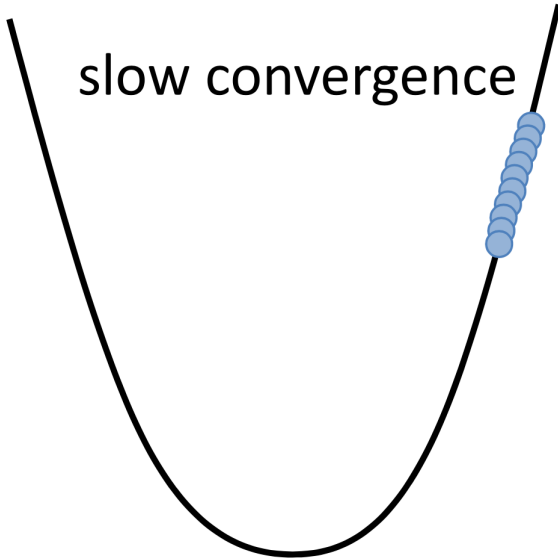


$$J(0.5) = \frac{1}{2} [(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2] = 1.75$$

# Choosing the step size alpha

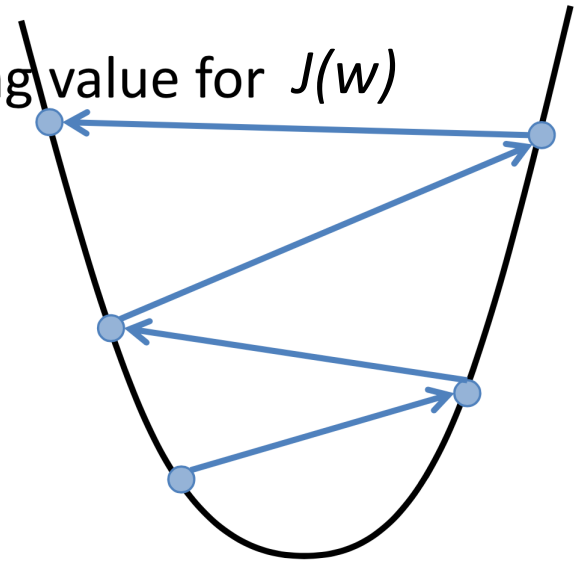
$\alpha$  too small

slow convergence



$\alpha$  too large

increasing value for  $J(w)$



- may overshoot minimum
- may fail to converge (may even diverge)

# Handout 6

## Linear Regression: SGD solution

*(find and work with a partner)*

In linear regression, we seek to minimize the sum of squared errors between the actual response and our prediction. We often call this RSS (residual sum of squares) or SSE (sum of squared errors). As an objective function, we often call it  $J$  and include a  $\frac{1}{2}$  in front to make the derivatives work out nicely.

$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^n (h_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2$$

For linear regression in general, one iteration of stochastic gradient descent includes the following updates (usually with the data points shuffled):

for  $i = 1, 2, \dots, n$ :

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha(\mathbf{w} \cdot \mathbf{x}_i - y_i)\mathbf{x}_i$$

We will begin with our same data from the previous two handouts:  $(x_1, y_1) = (0, 1)$  and  $(x_2, y_2) = (1, 0)$ , except we will reverse the order of the points to make the progress of gradient descent a bit clearer. So in this case our matrix/vector formulation is:

$$\mathbf{X} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

# Handout 6

$$X = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}, \quad y = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

1. Assuming  $\alpha = 0.1$  and our initial values are  $w_0 = 0$  and  $w_1 = 0$ , what are  $w_0$  and  $w_1$  after the just the first data point is used to update the gradient?

$$\vec{w} \leftarrow \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 0.1 \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \end{bmatrix} - 1 \right) \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\vec{w} \leftarrow \begin{bmatrix} 0 \\ 0 \end{bmatrix} + 0.1 \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\vec{w} \leftarrow \begin{bmatrix} 0.1 \\ 0 \end{bmatrix}$$

2. What are  $w_0$  and  $w_1$  after the second data point is used? Since we only have two examples here, your result would be the weight vector after the first iteration of SGD.

$$\vec{w} \leftarrow \begin{bmatrix} 0.1 \\ 0 \end{bmatrix} - 0.1 \left( \begin{bmatrix} 0.1 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} - 0 \right) \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\vec{w} \leftarrow \begin{bmatrix} 0.09 \\ -0.01 \end{bmatrix}$$

3. What is the value of the objective function (cost) after this initial iteration?

$$\hat{y} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 0.09 \\ -0.01 \end{bmatrix} = \begin{bmatrix} 0.09 \\ 0.08 \end{bmatrix}$$

$$J(\vec{w}) = \frac{1}{2} \begin{bmatrix} 0.09 \\ -0.08 \end{bmatrix} \cdot \begin{bmatrix} 0.09 \\ -0.08 \end{bmatrix}$$

$$\vec{y} - \hat{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} - \begin{bmatrix} 0.09 \\ 0.08 \end{bmatrix} = \begin{bmatrix} 0.91 \\ -0.08 \end{bmatrix}$$

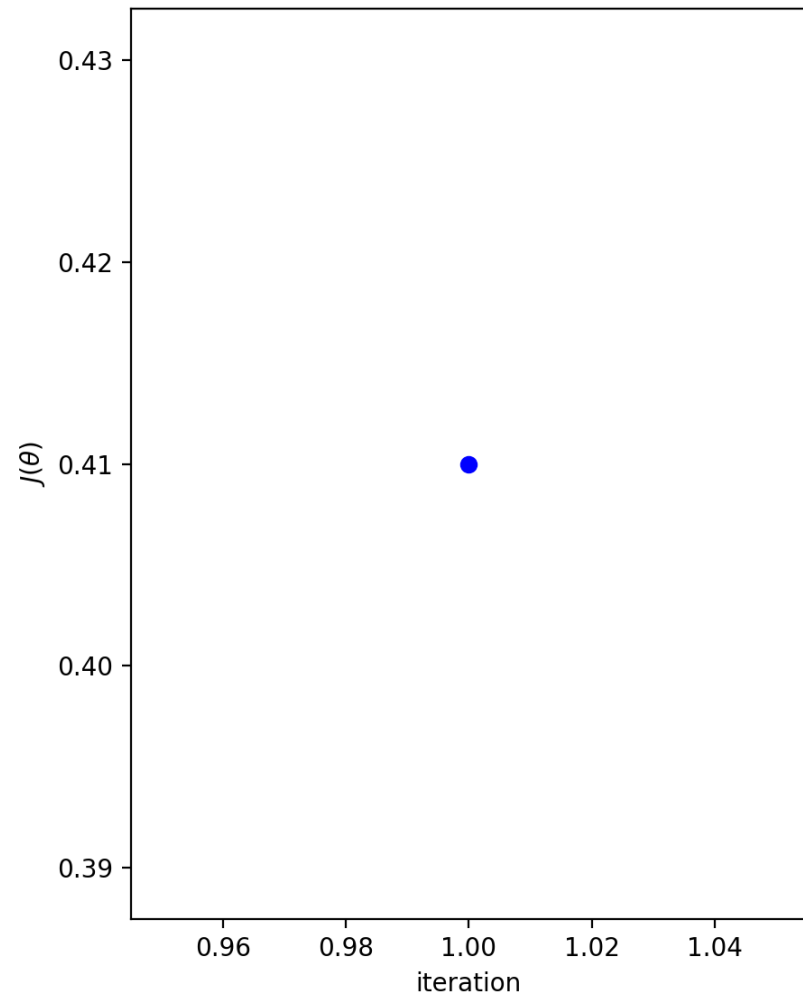
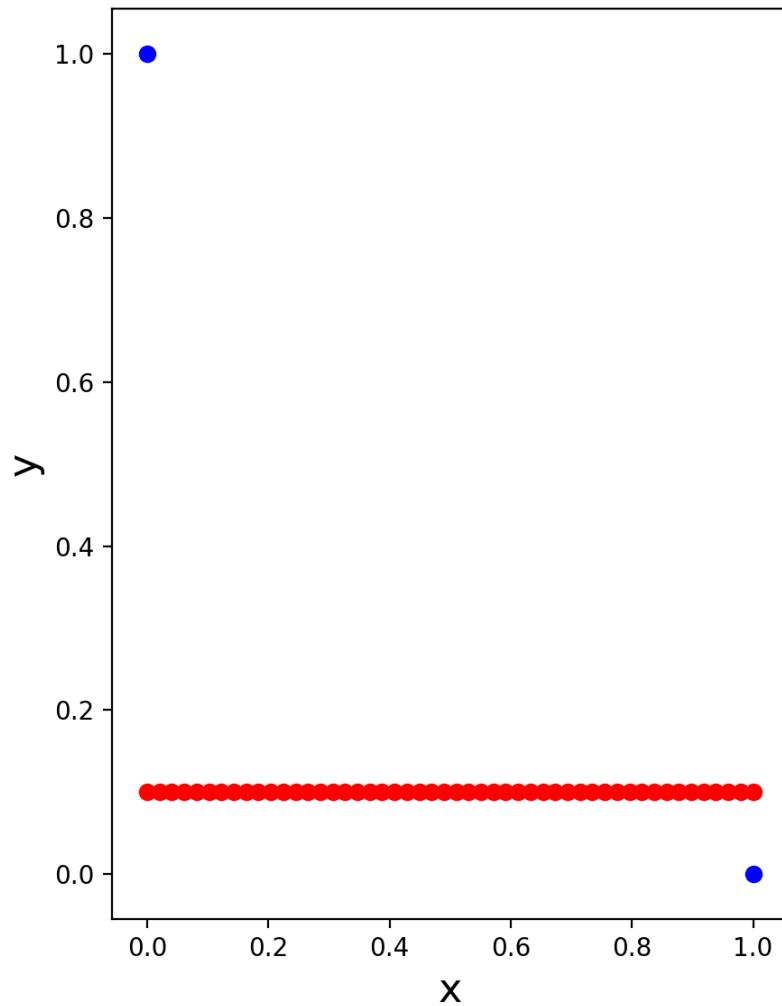
$$J(\vec{w}) = 0.417$$

# SGD with our small dataset from the handouts

Note: this is with the original order of the points

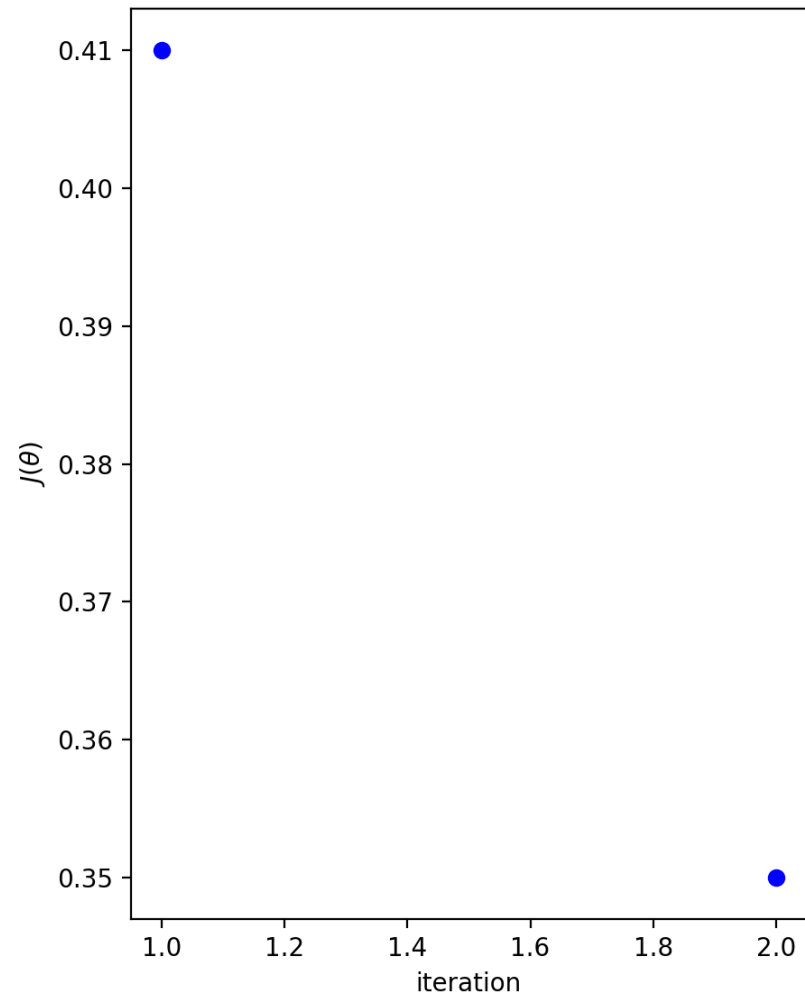
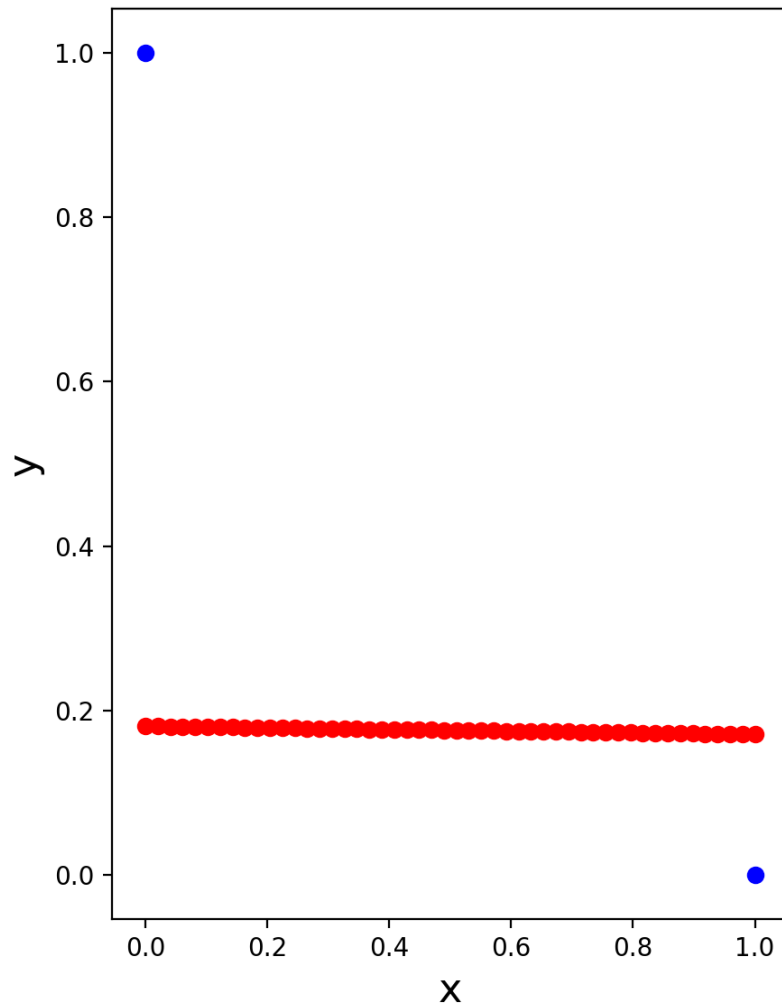
# Small example, iteration 1

iteration: 1, cost: 0.410000



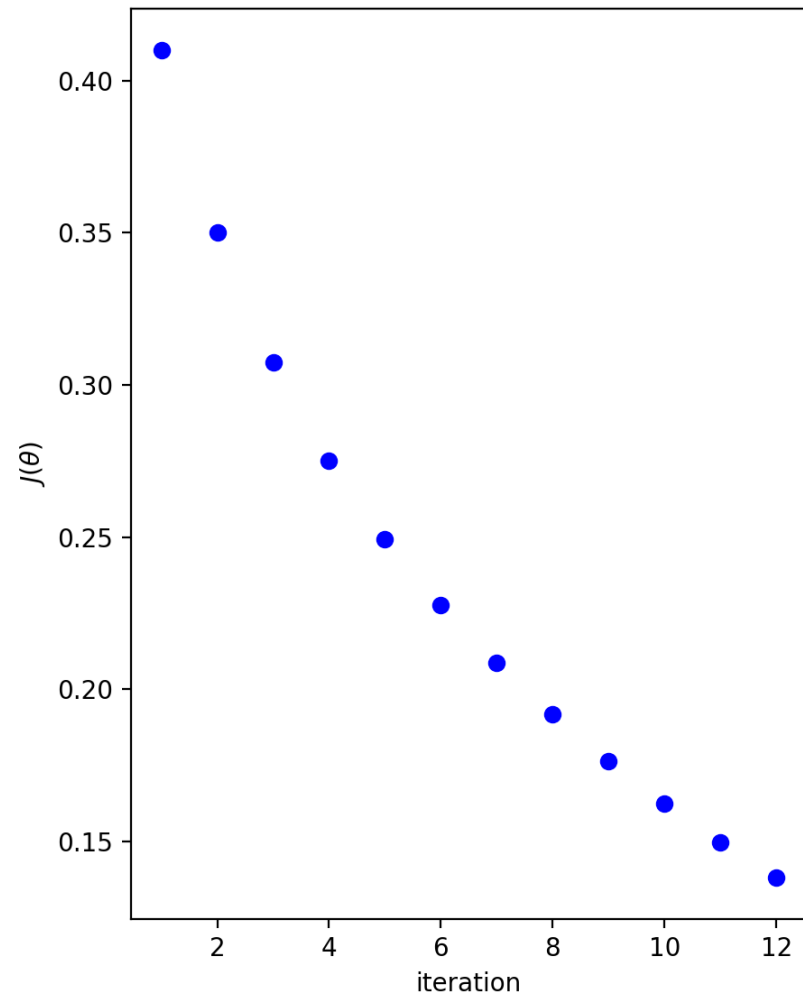
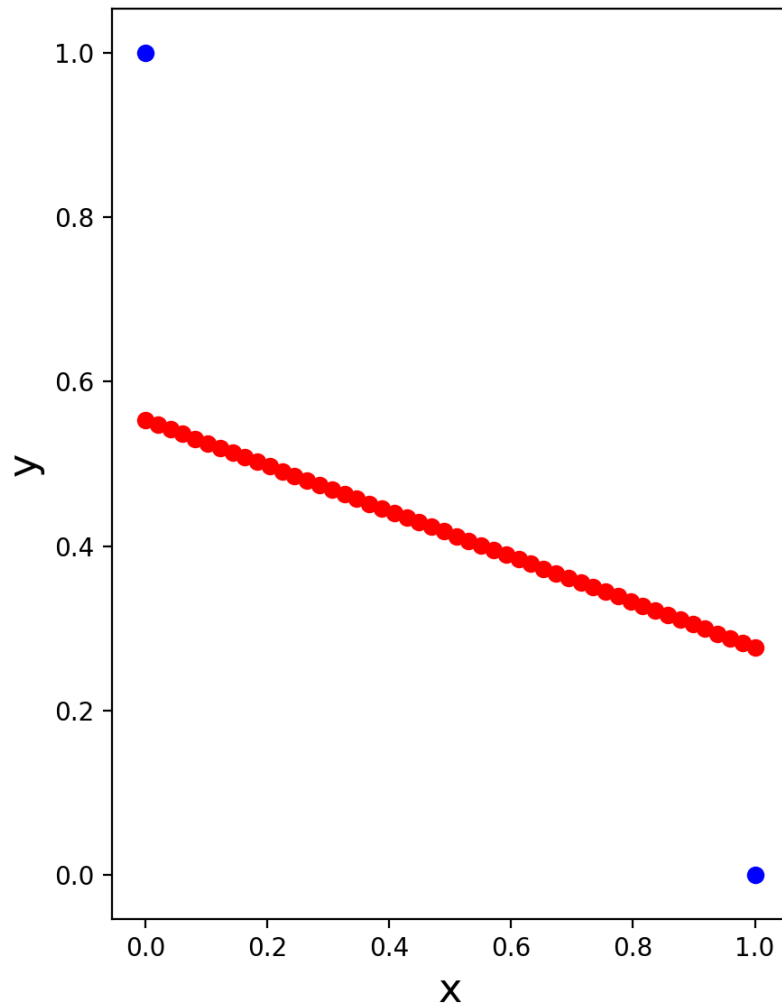
# Small example, iteration 2

iteration: 2, cost: 0.350001



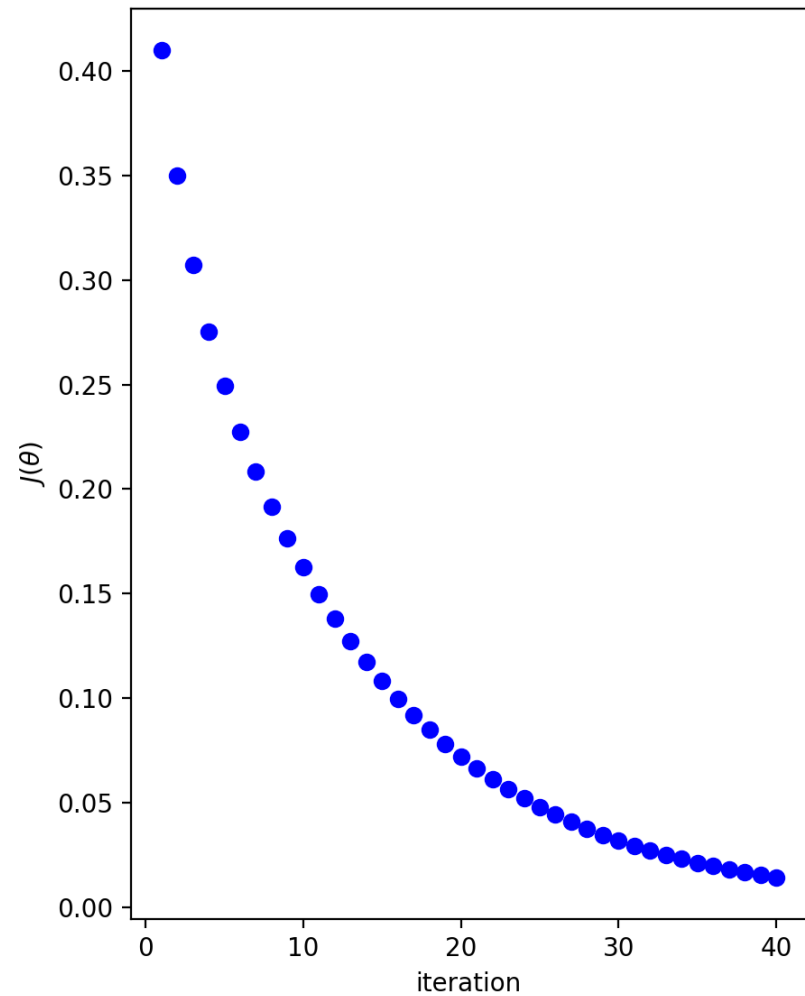
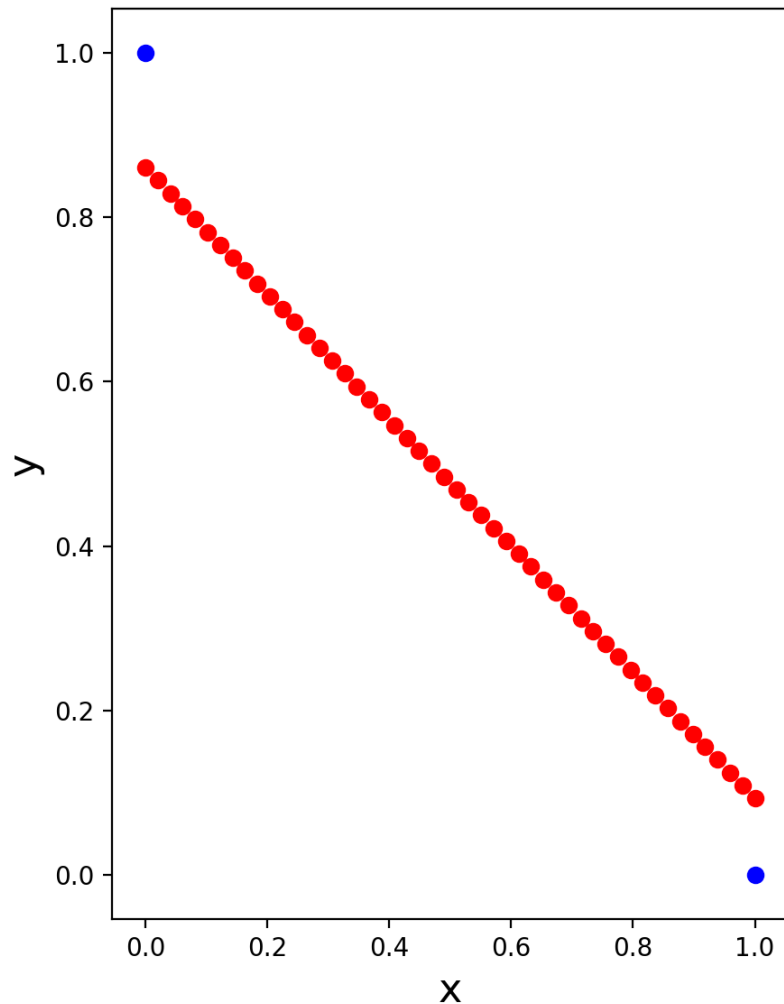
# Small example, iteration 12

iteration: 12, cost: 0.138047



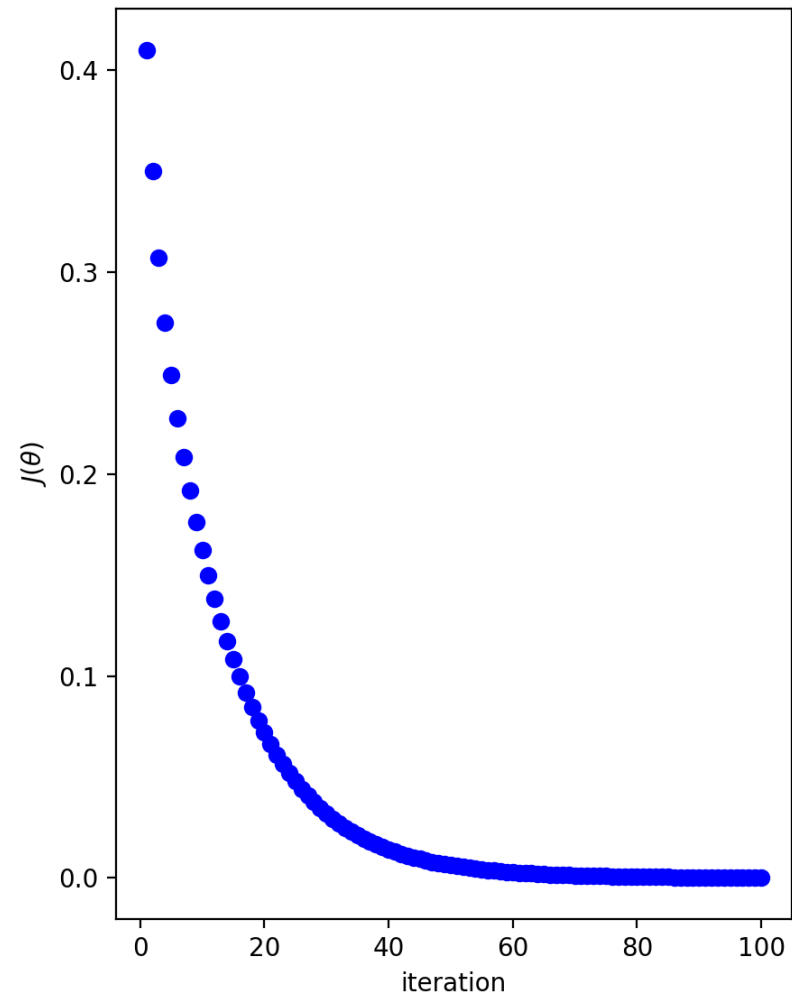
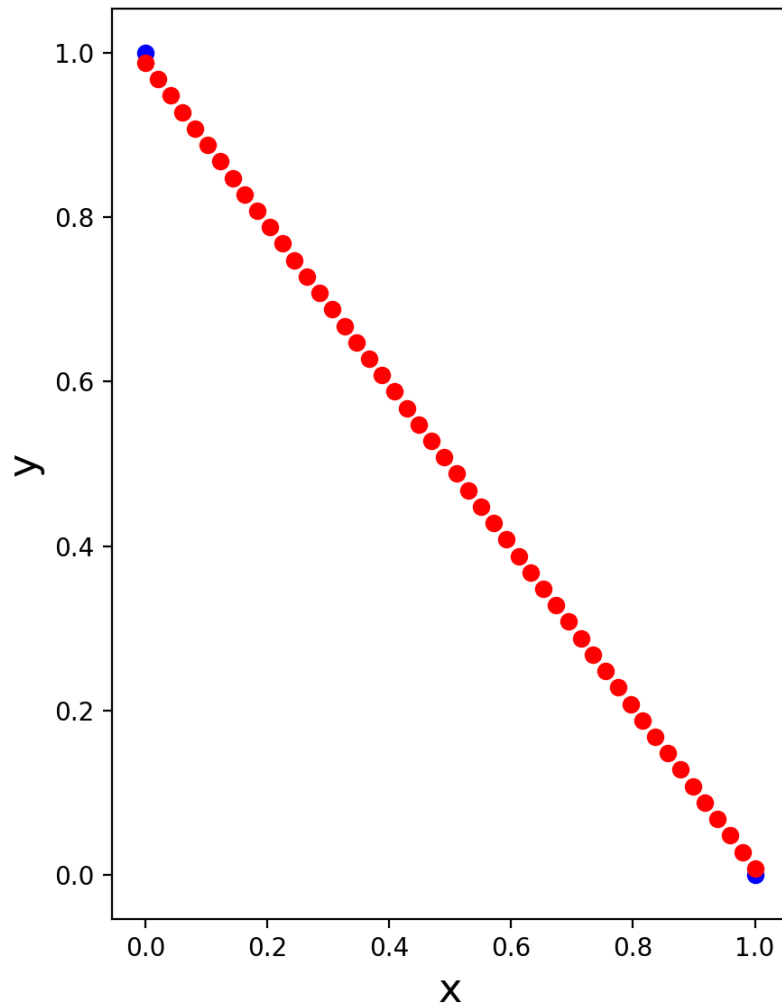
# Small example, iteration 40

iteration: 40, cost: 0.014064



# Small example, iteration 100

iteration: 100, cost: 0.000105



# Pros and Cons

(Analytic Solution)

## Gradient Descent

- requires multiple iterations
- need to choose  $\alpha$
- works well when  $p$  is large
- can support online learning

## Normal Equations

- non-iterative
- no need for  $\alpha$
- slow if  $p$  is large
  - matrix inversion is  $O(p^3)$

# Linear Regression Runtime

- $T$  = # iterations of SGD
- $n$  = # examples
- $p$  = # features

- 1) What is the runtime of SGD?
- 2) What is the runtime of the analytic solution?

# Outline for September 26

- Recap SGD (stochastic gradient descent)
- Introduction to classification
  - Decision tree models
  - Probabilistic interpretation
- Evaluation Metrics
  - Confusion matrices
  - Precision and recall
  - ROC curves

# Binary classification examples

- Transactions that indicate credit card fraud
- Accounts that are bots
- Detecting which scans show tumors
- Prenatal test for Down's Syndrome
- Finding genes under natural selection
- Finding regions of the genome with high recombination rate (“hotspots”)

In all these examples, we are trying to find unusual items (“needle in a haystack”) -- we call these *positives*

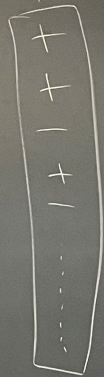
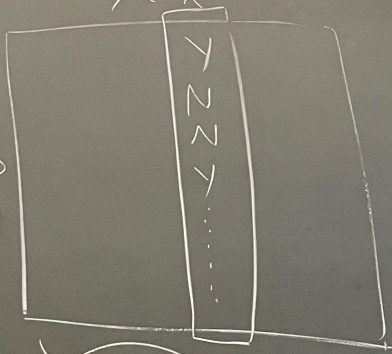
# Introduction to Classification

## Classification

X fever

Y (disease)

n examples



+  $\Rightarrow$  disease  
-  $\Rightarrow$  no disease

p features

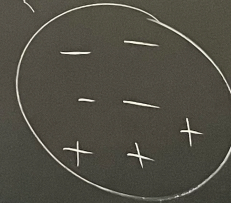
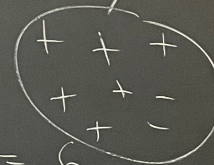
training data

model: decision tree with a single feature ("stump")

fever

Y

N



$$P_{\text{pos}} = \frac{6}{8}$$

$P_{\text{pos}}$  = prob of positive Y

$$P_{\text{pos}} = \frac{3}{7}$$

n=15

# Introduction to Classification

new idea : use probabilities  
to classify test examples

$$\vec{X}_{\text{test}} = \begin{bmatrix} \dots & \text{fever} & N & \dots \end{bmatrix}^T$$

threshold 0.5  $\Rightarrow$

$$\hat{y}_{\text{test}} = \ominus$$

no  
disease

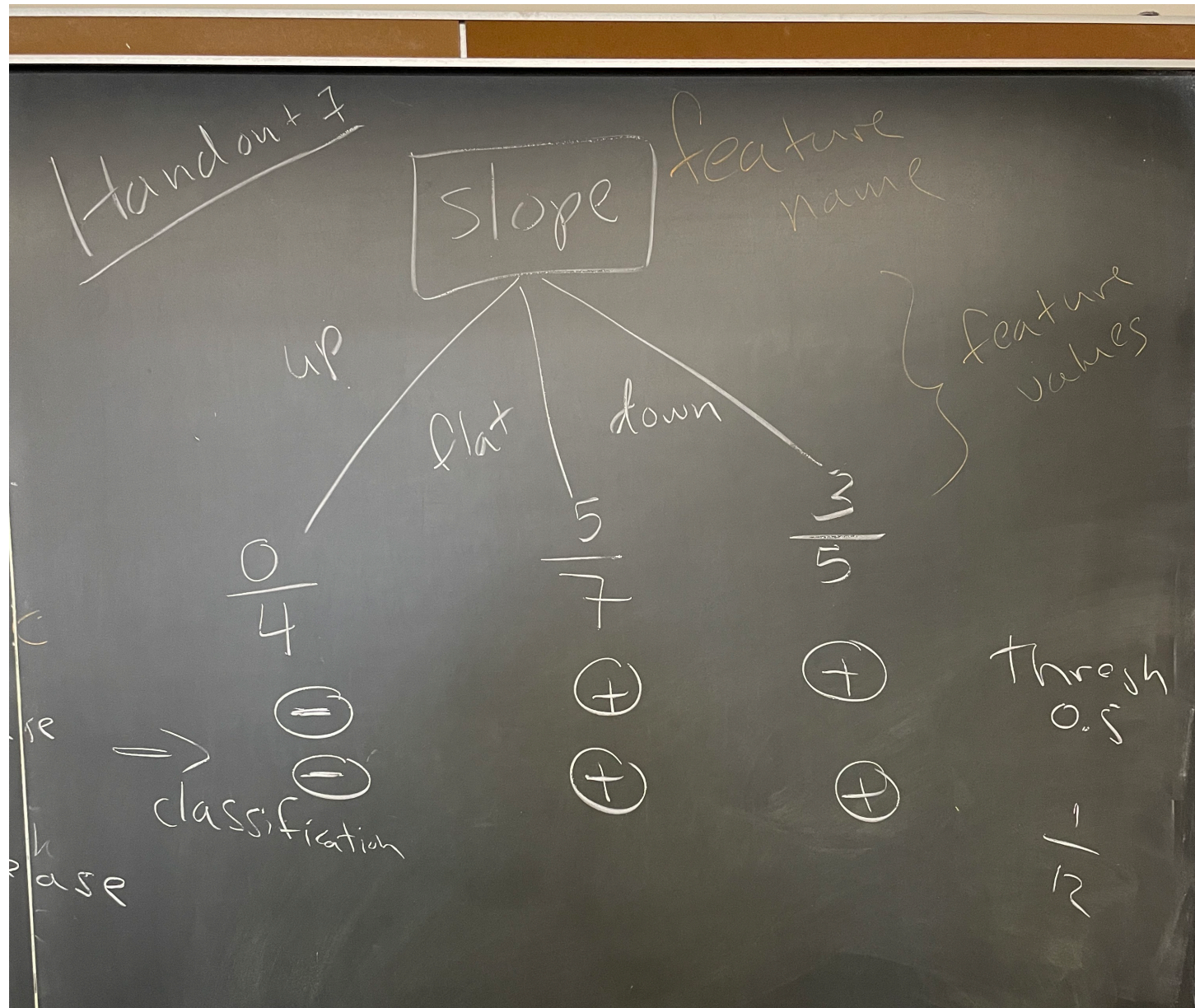
threshold 0.25  $\Rightarrow$

$$\hat{y}_{\text{test}} = \oplus$$

disease

$$P_{\text{pos}} \geq \text{threshold} \Rightarrow \text{classify } \oplus$$

# Handout 7



# Outline for September 26

- Recap SGD (stochastic gradient descent)
- Introduction to classification
  - Decision tree models
  - Probabilistic interpretation
- Evaluation Metrics
  - Confusion matrices
  - Precision and recall
  - ROC curves

Next time!