

CS 66: Machine Learning

Prof. Sara Mathieson

Spring 2019



Admin

- Office hours **TODAY!** 1-3pm (Sci Center 249)
- Lab 2 due **Tuesday night**
 - See Piazza for dictionary example
 - See Piazza for lab suggestions

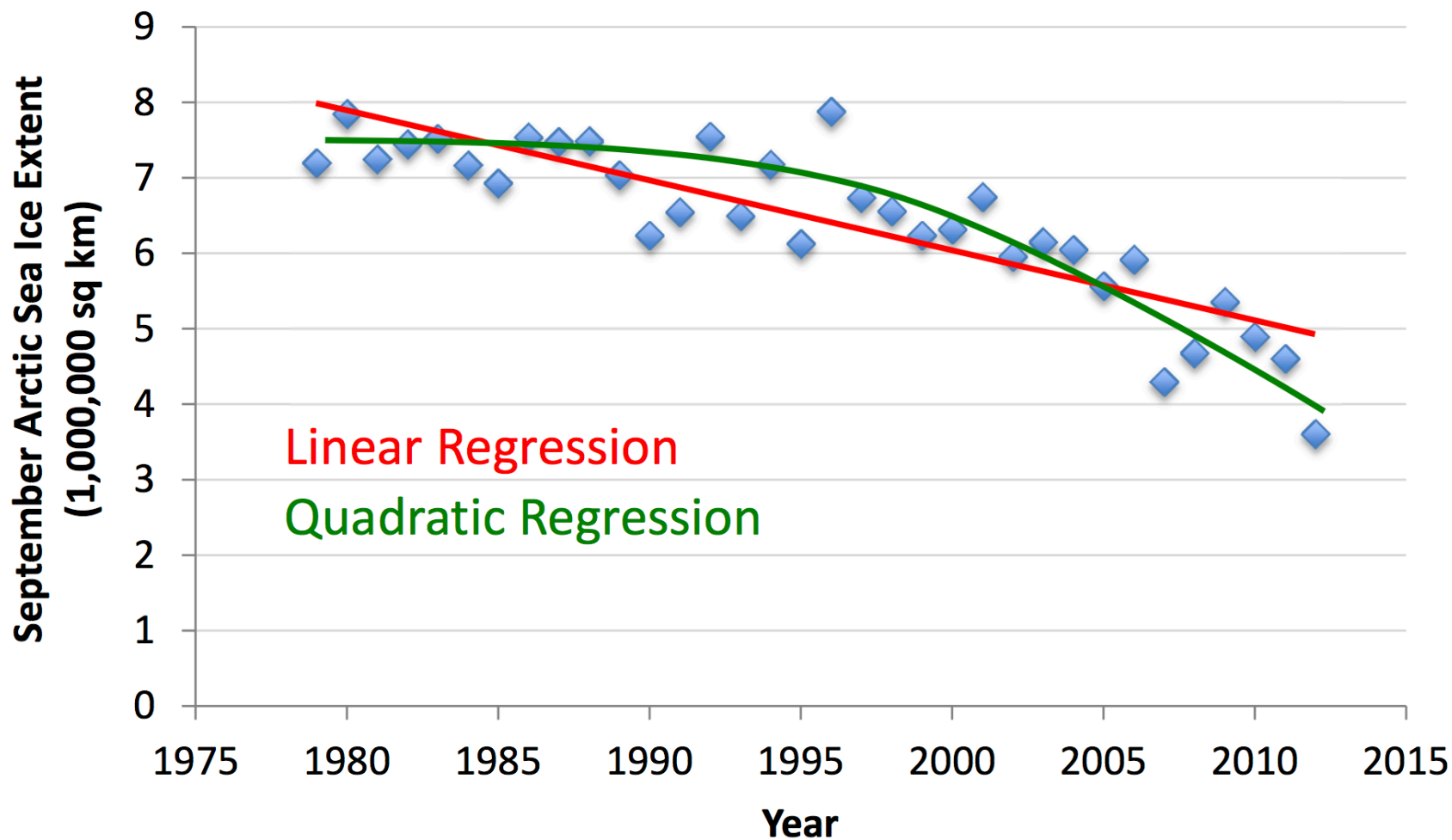
Outline for February 8

- Where we are with regression so far
- Cost function and multiple linear regression
- SGD (Stochastic Gradient Descent)
- Normal equations solution

Outline for February 8

- Where we are with regression so far
- Cost function and multiple linear regression
- SGD (Stochastic Gradient Descent)
- Normal equations solution

Regression Example



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 squared errors) or SSE (sum of squared errors)

Simple Linear Regression

- X only contains one feature
- Linear model can be described by a *slope* and a *y-intercept* (bias)

$$h_b(\mathbf{x}) = b_0 + b_1x$$

Simple Linear Regression

- X only contains one feature
- Linear model can be described by a *slope* and a *y-intercept* (bias)

$$h_b(\mathbf{x}) = b_0 + b_1 x$$

- Solution from Wed

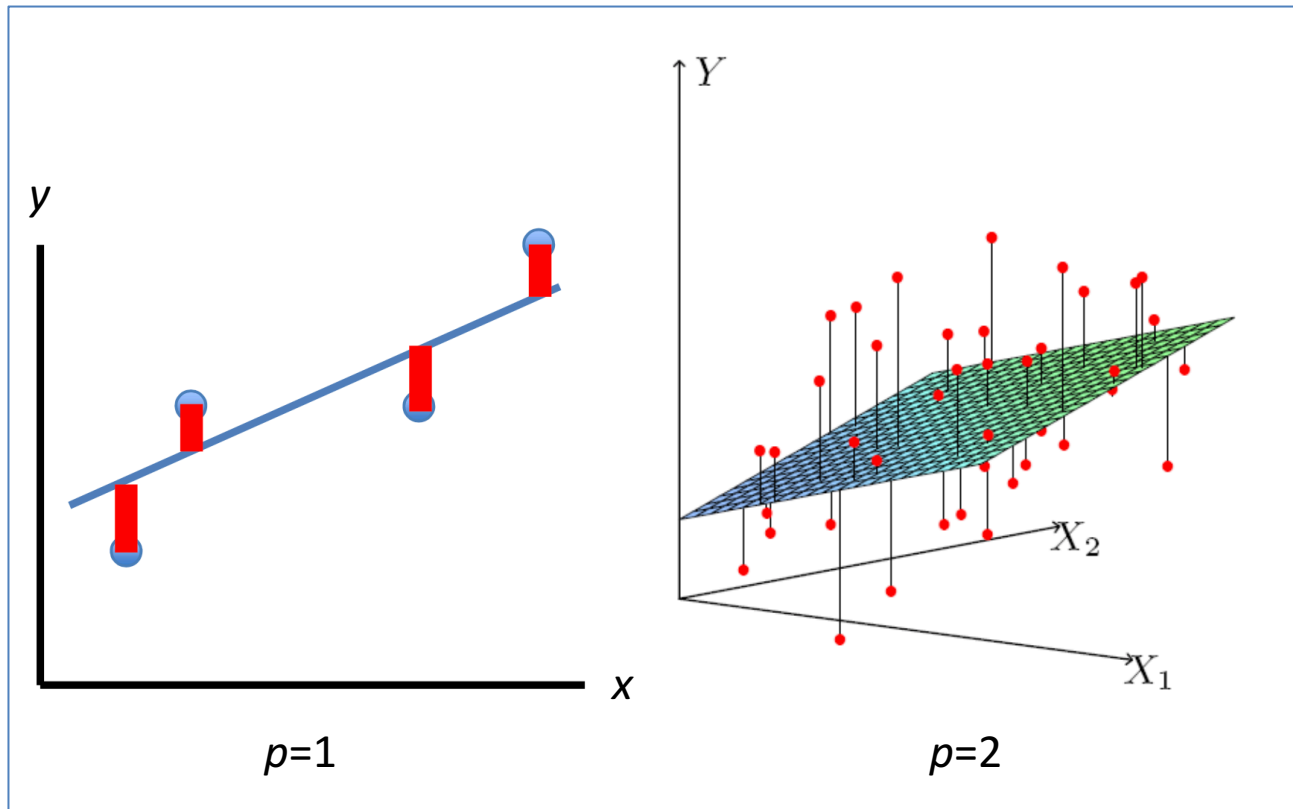
$$\hat{b}_1 = \frac{\text{Cov}(\mathbf{x}, \mathbf{y})}{\text{Var}(\mathbf{x})} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{b}_0 = \bar{y} - \hat{b}_1 \bar{x}$$

Outline for February 8

- Where we are with regression so far
- Cost function and multiple linear regression
- SGD (Stochastic Gradient Descent)
- Normal equations solution

Cost Function: sum of squared errors



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

Multiple Linear Regression

$$\hat{y} = h_{\vec{b}}(\vec{x}) = b_0 x_0 + b_1 x_1 + b_2 x_2 + \dots + b_p x_p$$

$x_0 = 1$

$$h_{\vec{b}}(\vec{x}) = \vec{b}^T \vec{x}$$

minimize

$$J(\vec{b}) = \frac{1}{2} \sum_{i=1}^n (\vec{b}^T \vec{x}_i - y_i)^2$$

cost function

$$[b_0 \ b_1 \ \dots \ b_p]$$

$$\begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix}$$

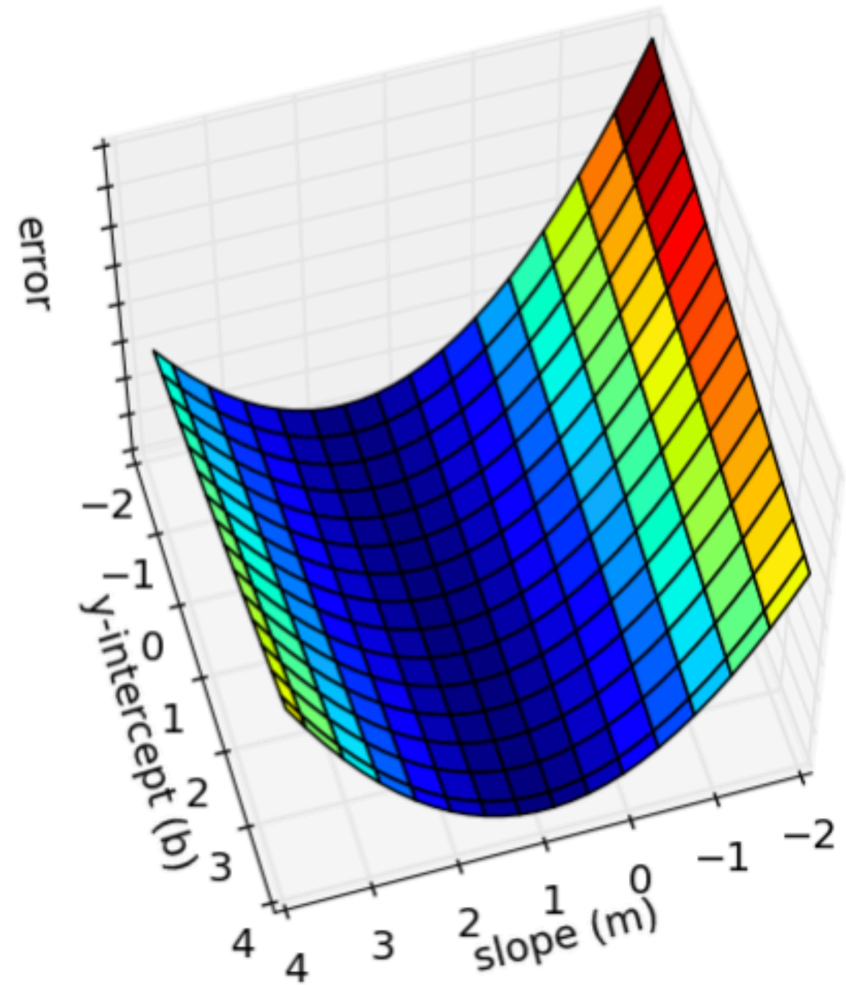
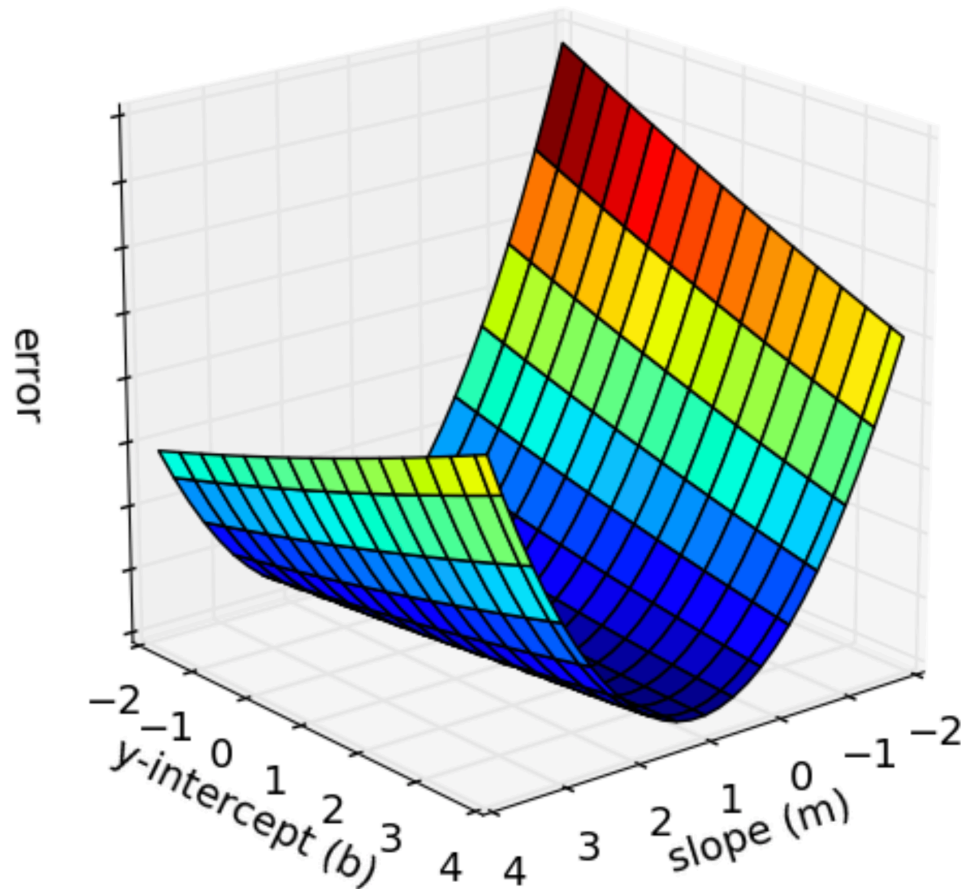
$|x(p+1)|$
 $(p+1) \cdot 1$

$$X^T = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}$$

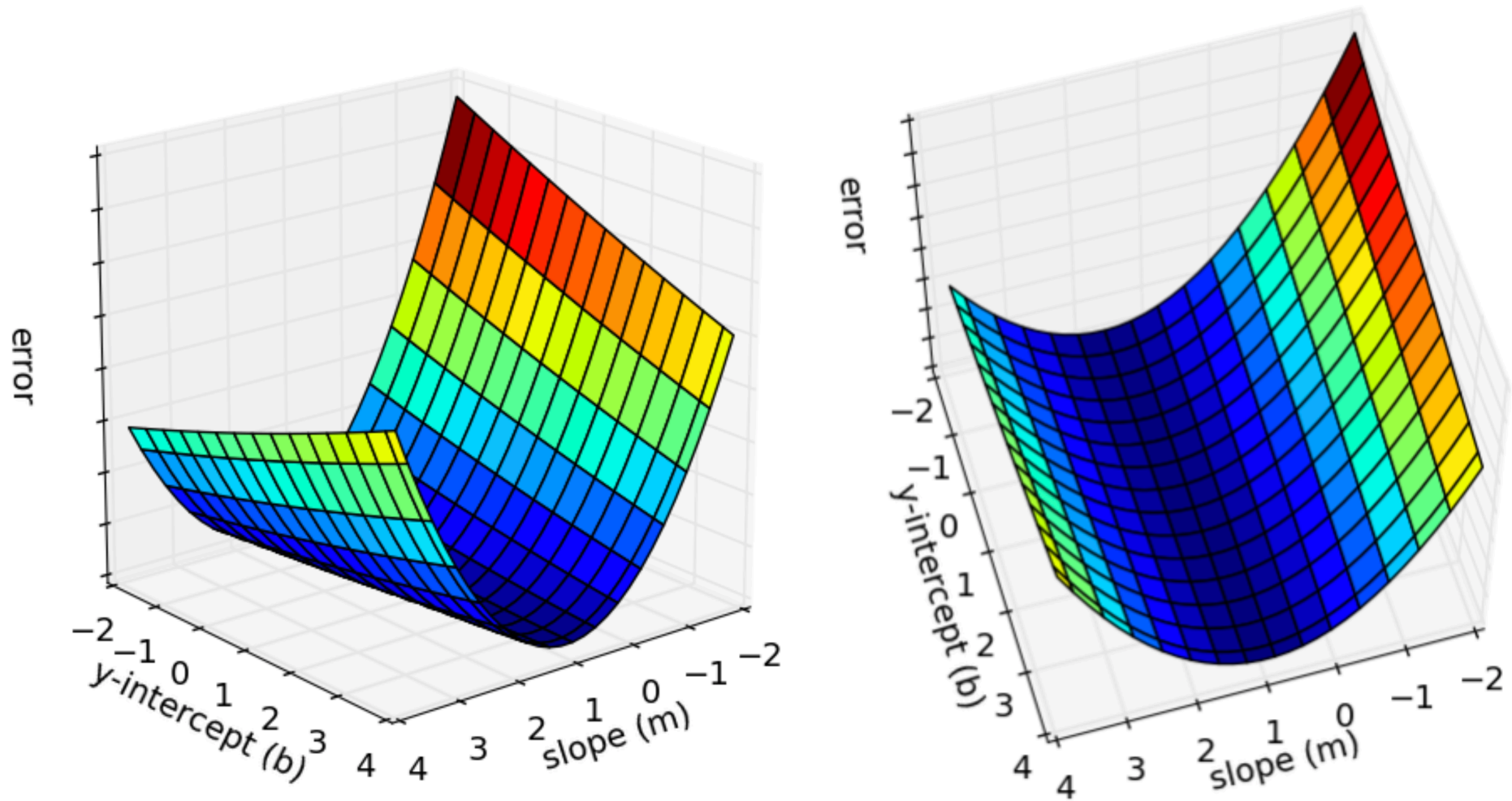
n rows, $p+1$ columns

$$= [\quad]$$

Error as a function of slope & y-intercept

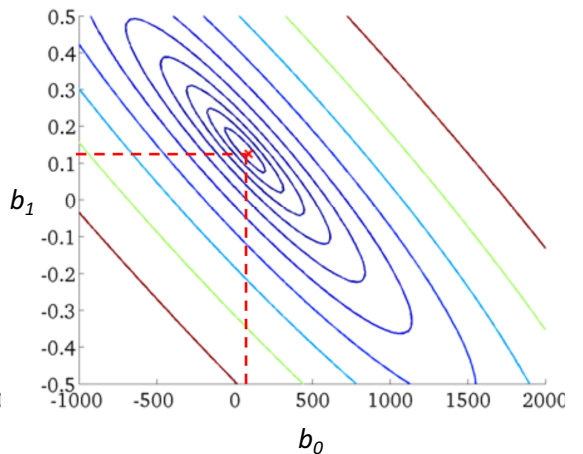
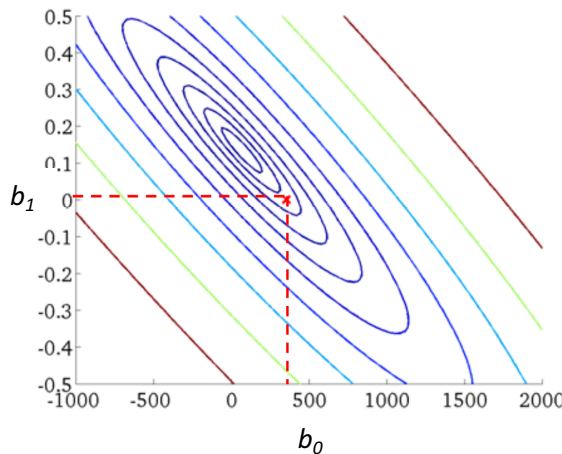
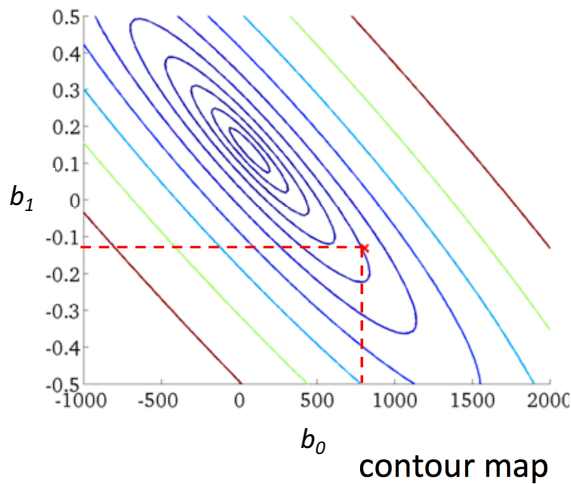
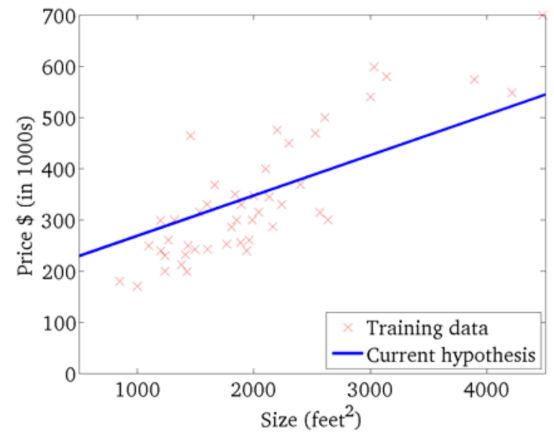
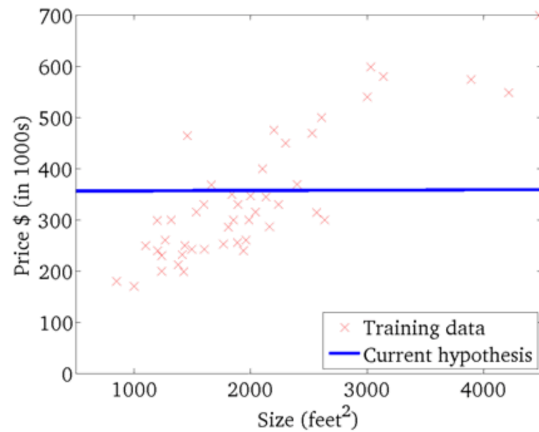
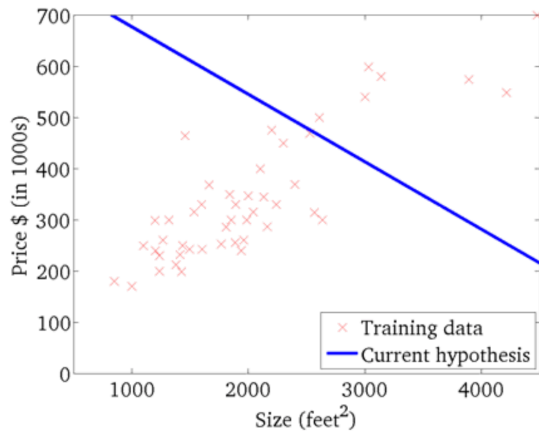


Error as a function of slope & y-intercept

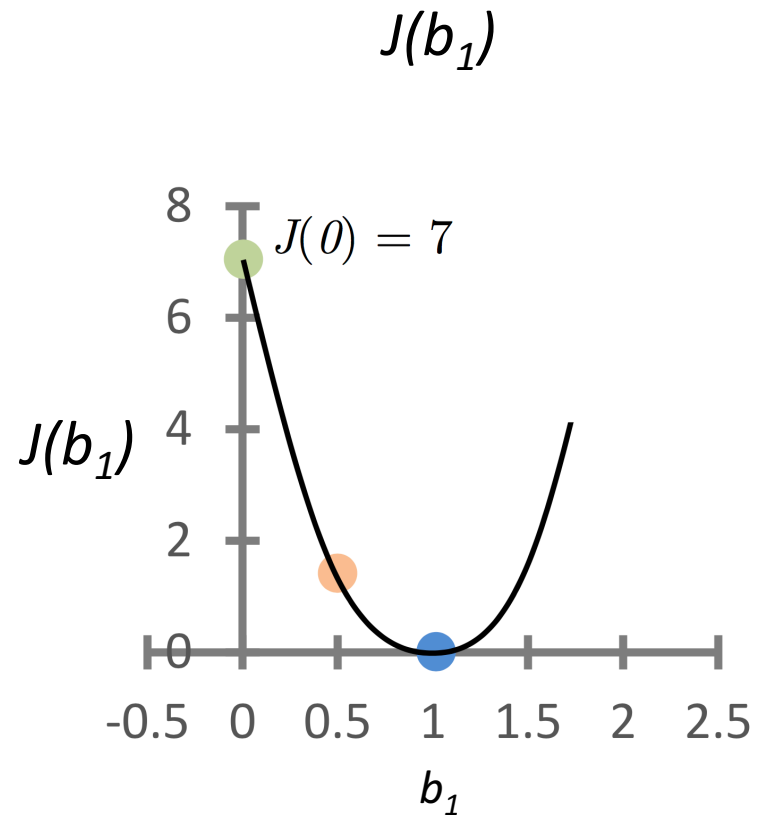
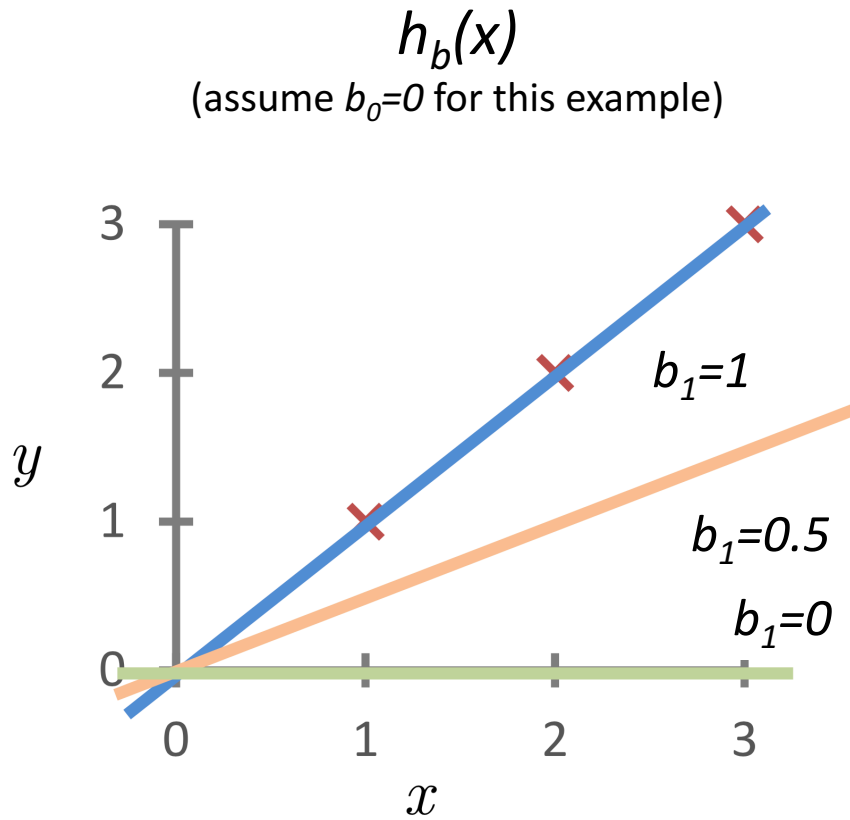


Convex => single global optimum (can prove it is a minimum with second derivative!)

Cost Function



Cost Function (extra example)



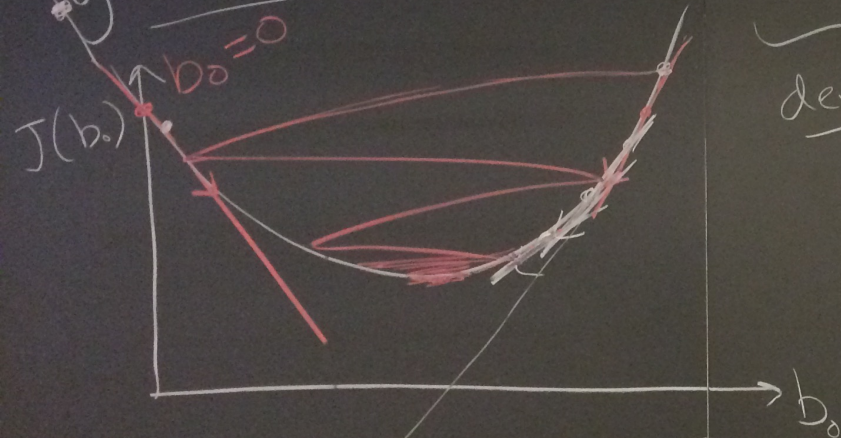
$$J(0.5) = \frac{1}{2} [(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2] = 1.75$$

Outline for February 8

- Where we are with regression so far
- Cost function and multiple linear regression
- **SGD (Stochastic Gradient Descent)**
- Normal equations solution

Method 1

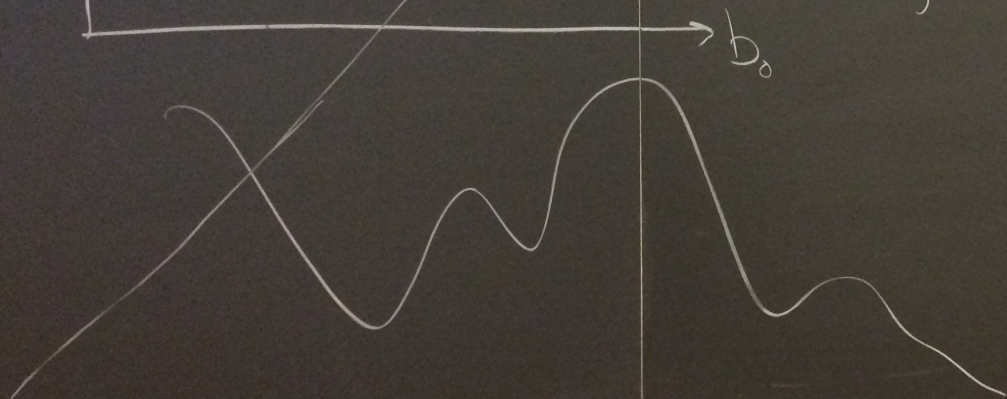
gradient descent



$$\underbrace{J(\vec{b})}_{\text{derivative}} \Big|_{b_j} = \sum_{i=1}^n (\vec{b}^T \vec{x}_i - y_i) x_{ij}$$

batch gradient descent for all $j=0 \dots p$

$$b_j \leftarrow b_j - \alpha \underbrace{\sum_{i=1}^n (\vec{b}^T \vec{x}_i - y_i) x_{ij}}_{\text{Slow!}}$$



Stochastic gradient descent

shuffle the n data points

for $i = 1 \dots n$:

for $j = 0 \dots p$:

$$b_j \leftarrow b_j - \alpha (\vec{b}^T \vec{x}_i - y_i) x_{ij}$$

$$\alpha = \frac{1}{\text{iteration}}$$

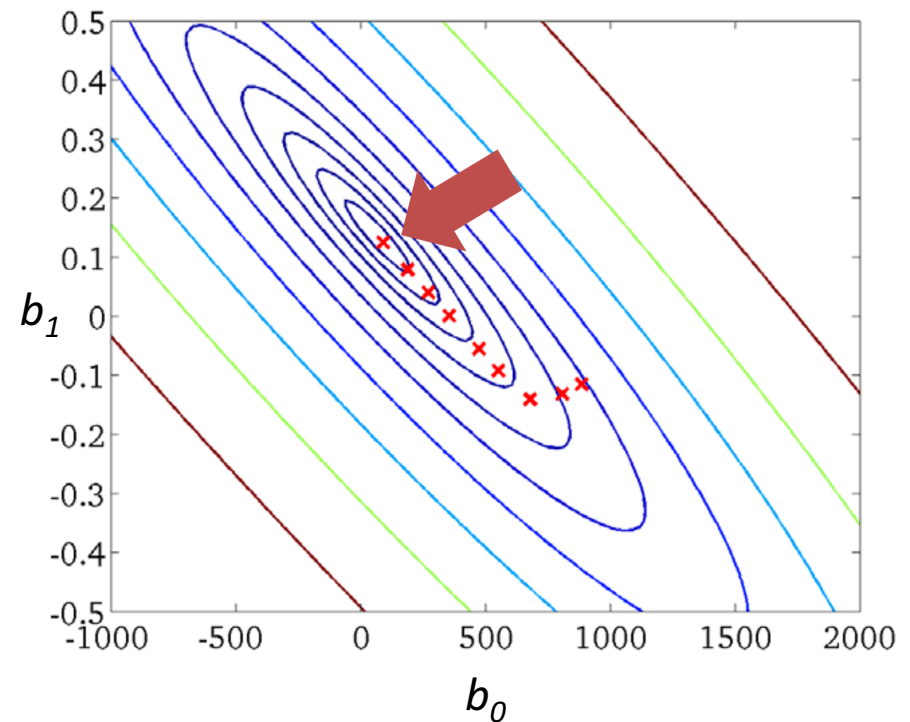
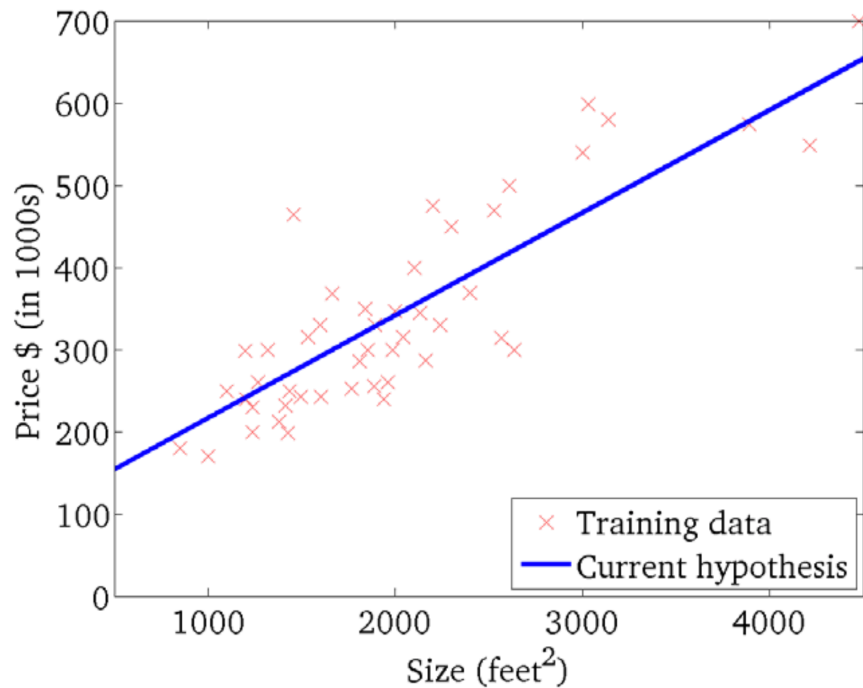
iterate til convergence
"J doesn't change"

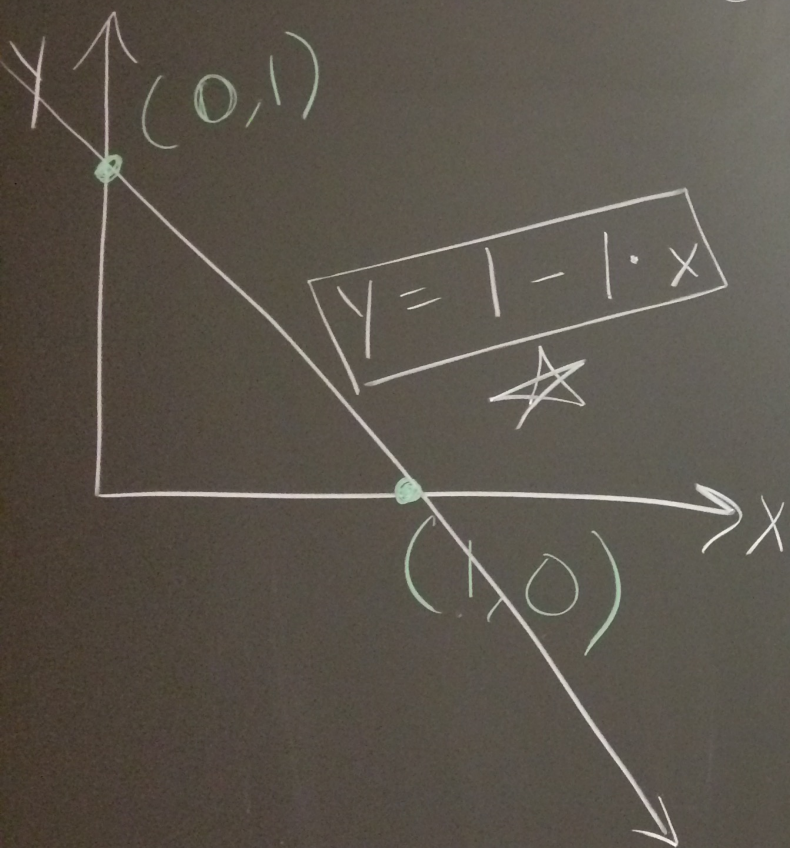
$$\alpha = 0.01$$

for all
 $j = 0 \dots p$

$$- y_i) x_{ij}$$

Gradient Descent: walking toward the minimum





①

$$\frac{\text{COV}(X, Y)}{\text{Var}(X)} = \frac{\frac{1}{2}[(1 - \frac{1}{2})(0 - \frac{1}{2}) + (0 - \frac{1}{2})(1 - \frac{1}{2})]}{\frac{1}{2}[(1 - \frac{1}{2})^2 + (0 - \frac{1}{2})^2]}$$

$$\hat{b}_1 = \frac{-\frac{1}{4} - \frac{1}{4}}{\frac{1}{4} + \frac{1}{4}} = -1$$

$$\hat{b}_0 = \frac{1}{2} - (-1) \cdot \frac{1}{2}$$

$$\hat{b}_0 = 1$$

$$)(1-\frac{1}{2})\}$$

$$)^2\}$$

$$i=1: (x_1, y_1) = (1, 0)$$

$$\left. \begin{aligned} b_0 &\leftarrow 0 - 0.1(0 - 0) \cdot 1 \\ b_1 &\leftarrow 0 - 0.1(0 - 0) \cdot 1 \end{aligned} \right\} \text{all together}$$

$$\boxed{b_0 = 0, b_1 = 0}$$

$$i=2: (x_2, y_2) = (0, 1)$$

$$b_0 \leftarrow 0 - 0.1(0 - 1) \cdot 1$$

$$b_1 \leftarrow 0 - 0.1(0 - 1) \cdot 0$$

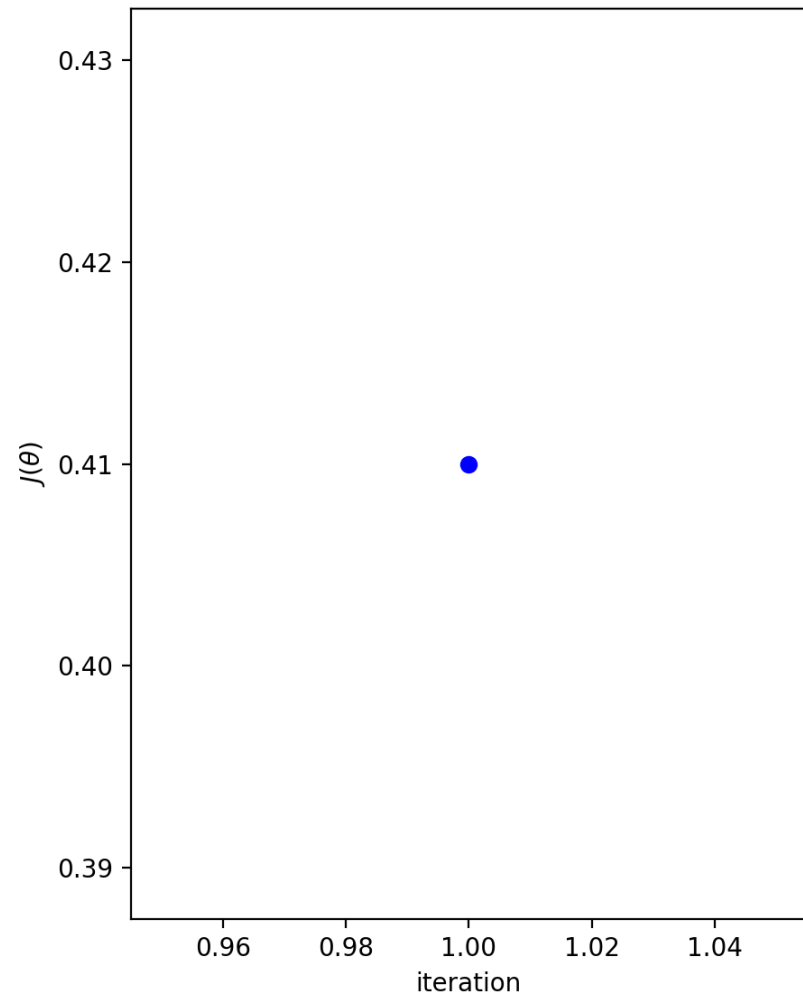
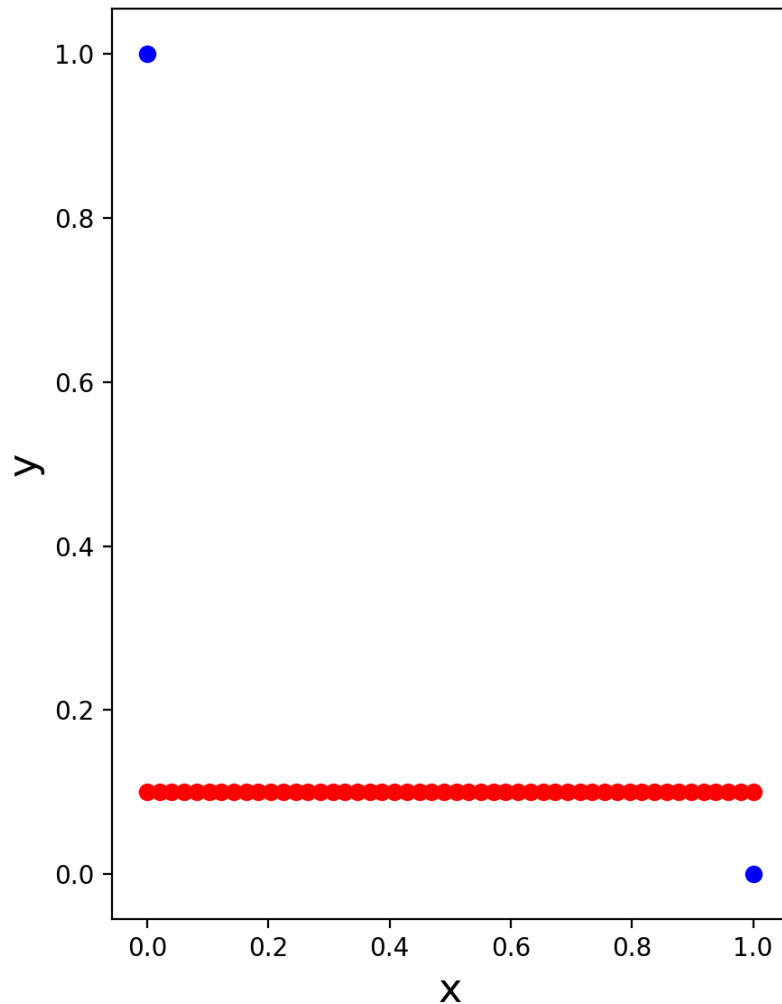
$$\boxed{\begin{aligned} b_0 &= 0.1 \\ b_1 &= 0 \end{aligned}}$$

Handout 2, Question 2(a)

Toy example, iteration 1

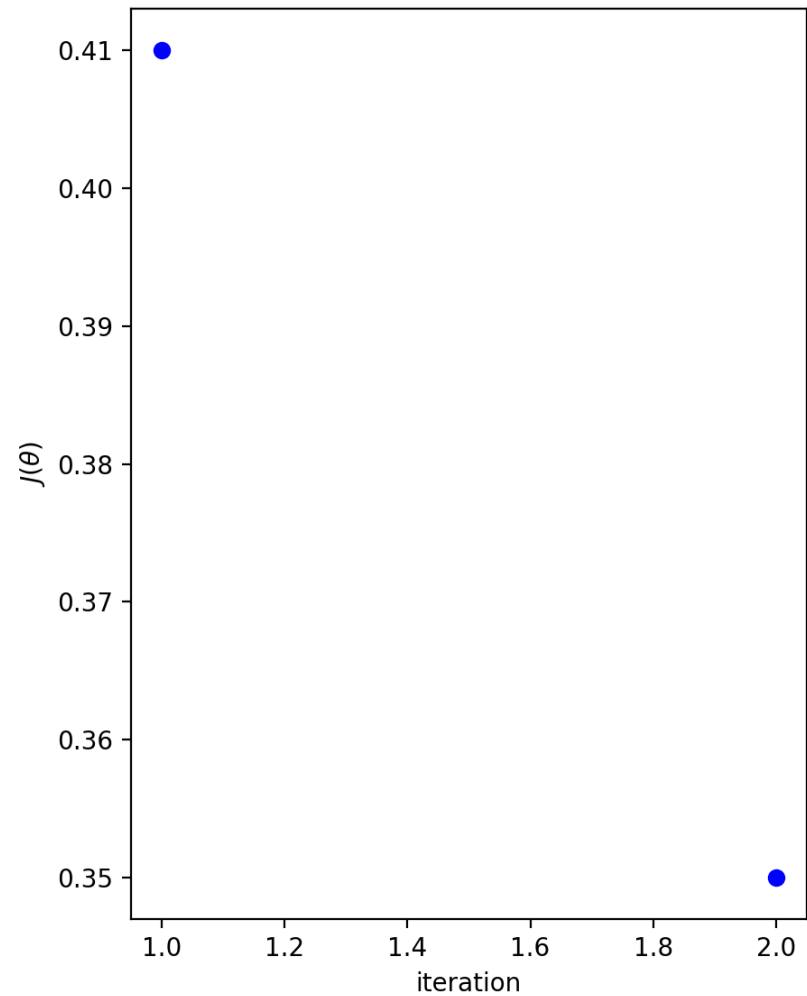
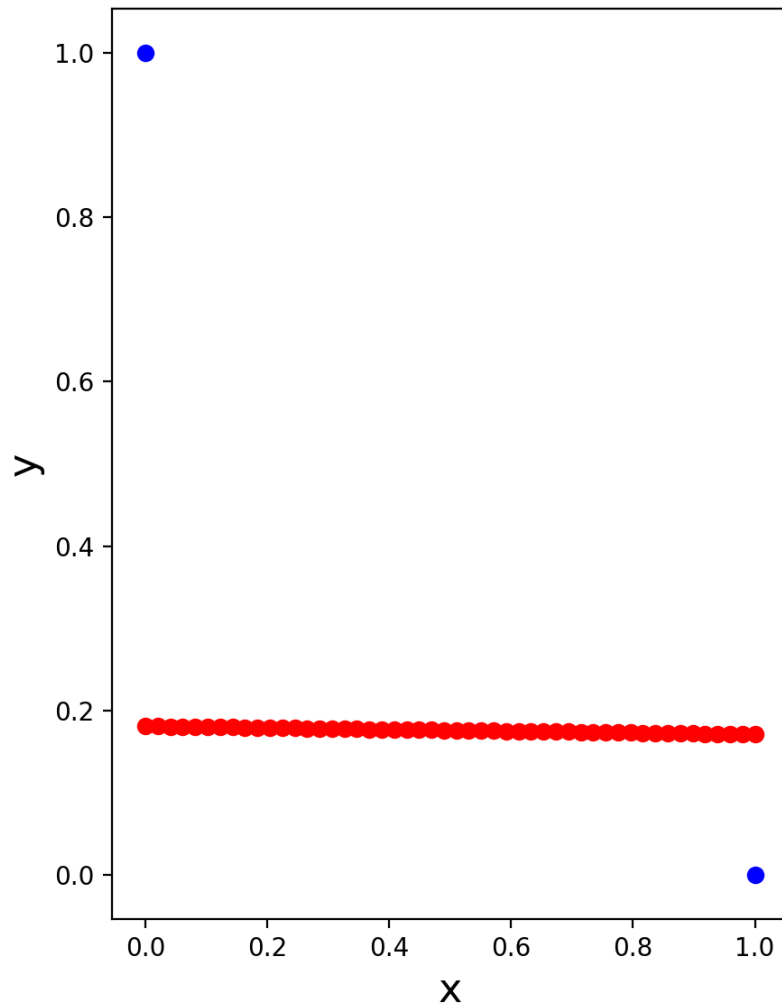
This is what you
should have obtained
in Handout 2!

iteration: 1, cost: 0.410000



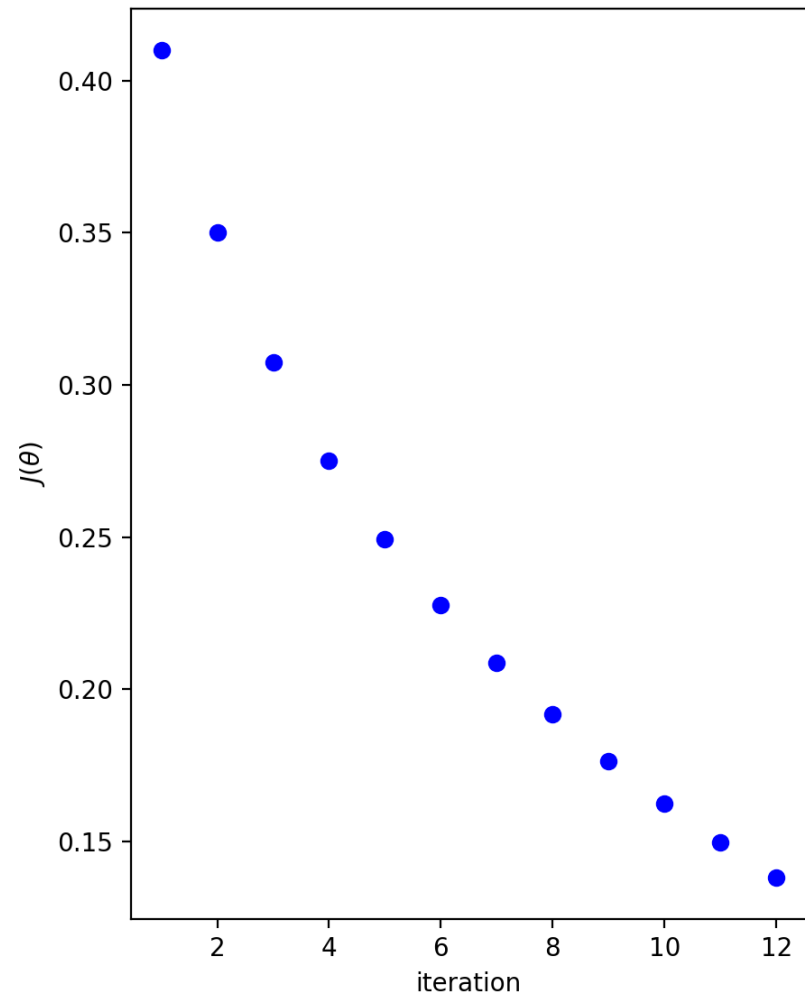
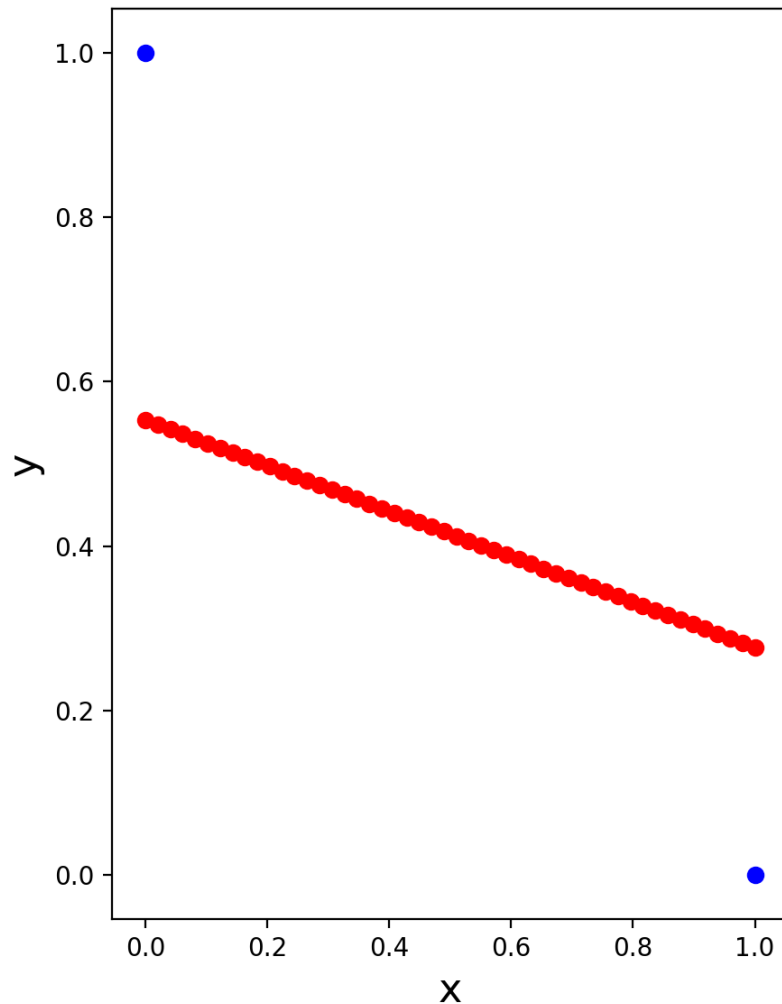
Toy example, iteration 2

iteration: 2, cost: 0.350001



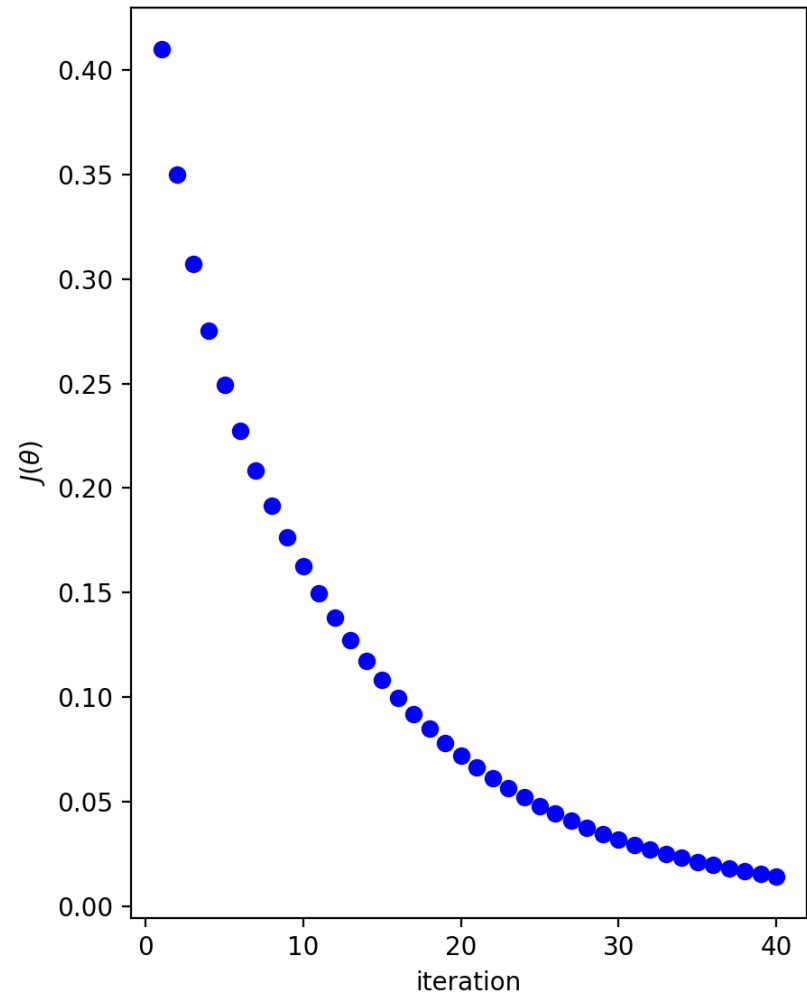
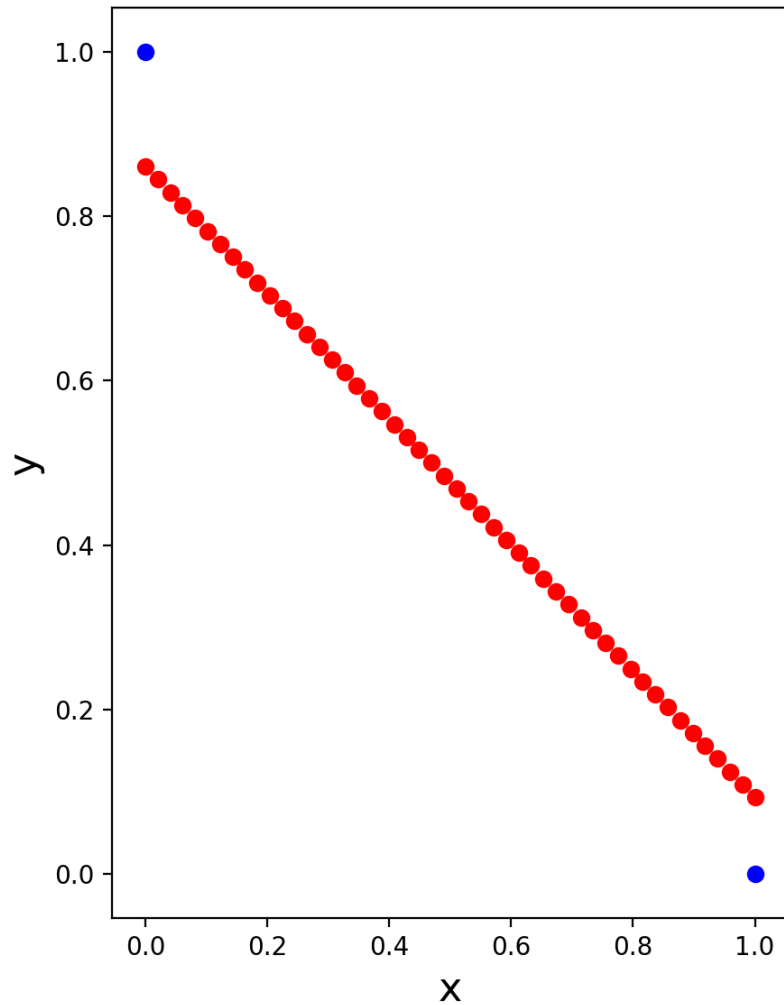
Toy example, iteration 12

iteration: 12, cost: 0.138047



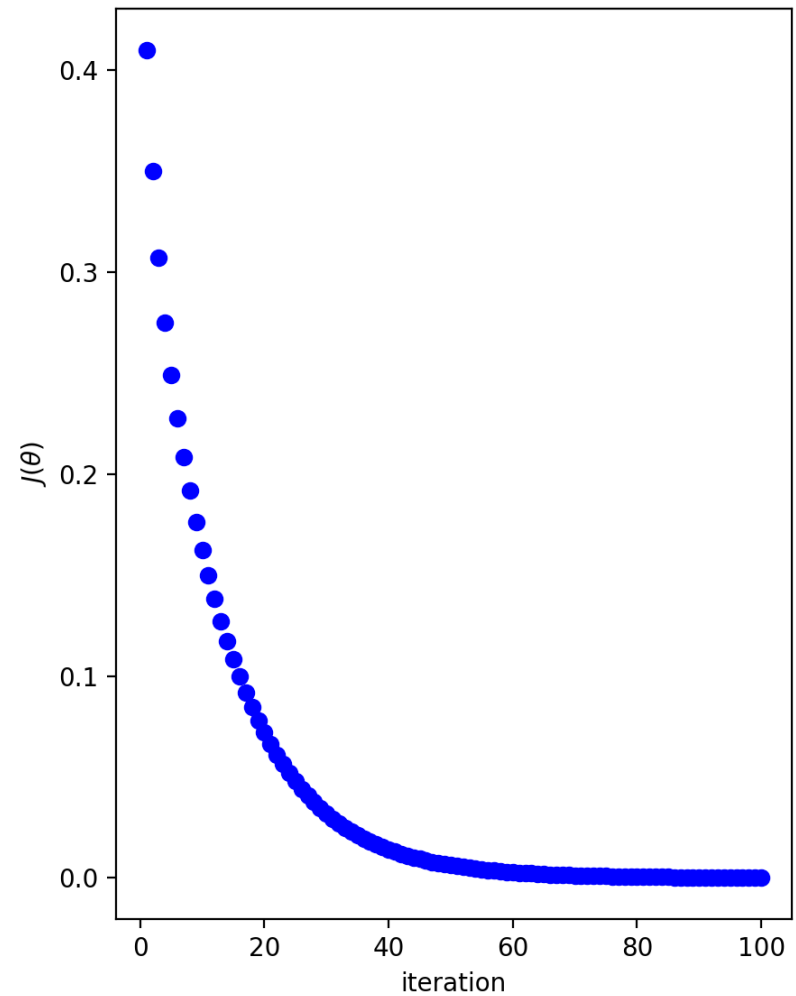
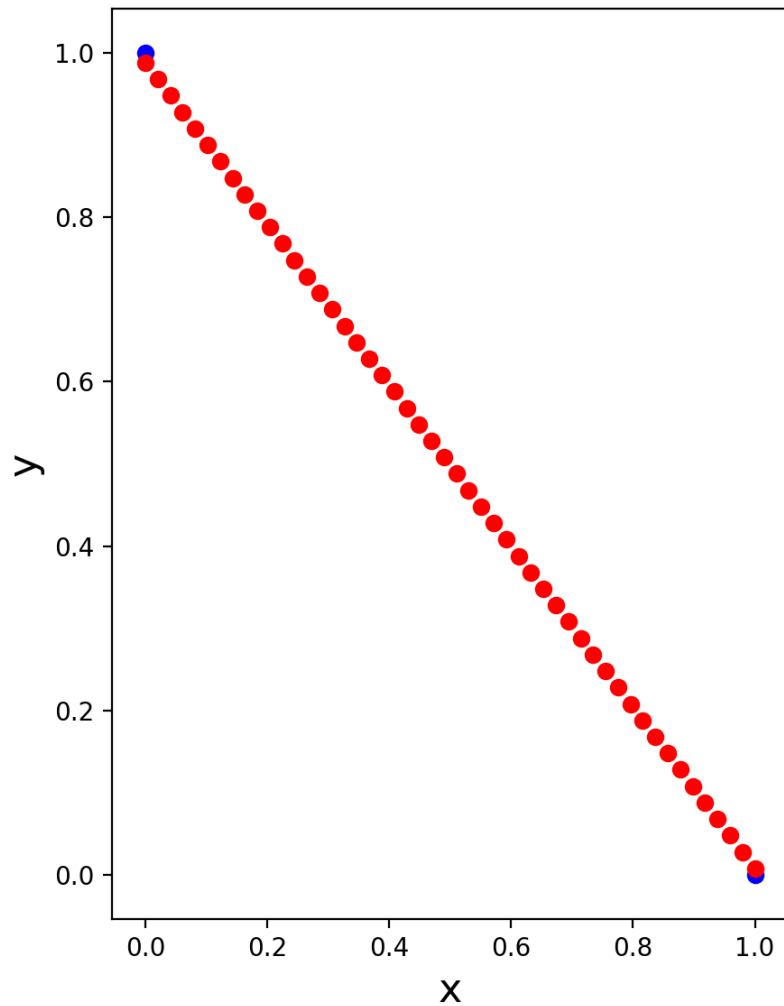
Toy example, iteration 40

iteration: 40, cost: 0.014064



Toy example, iteration 100

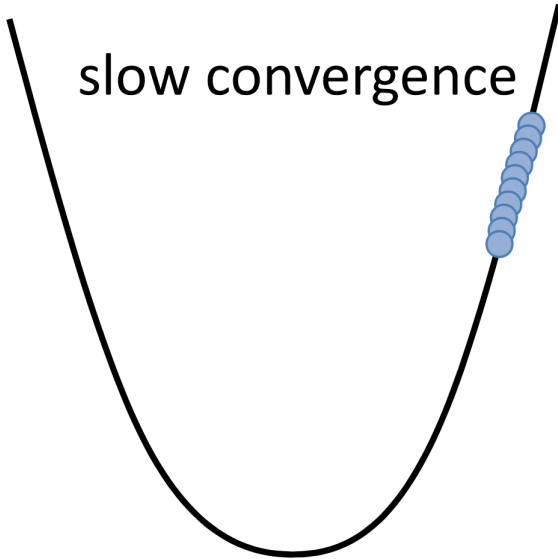
iteration: 100, cost: 0.000105



Choosing step size α

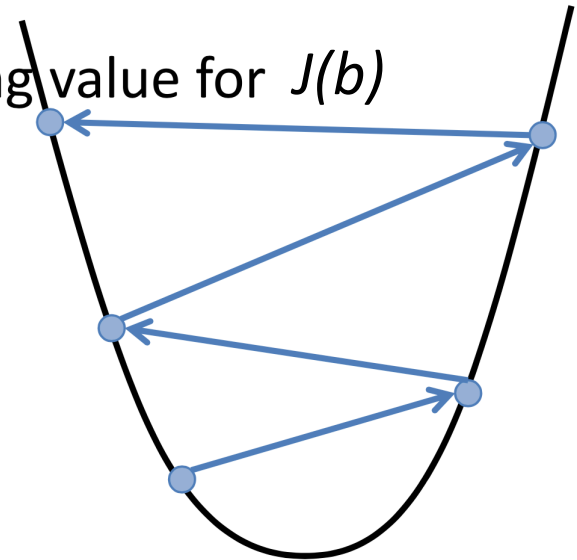
α too small

slow convergence



α too large

increasing value for $J(b)$



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

Outline for February 8

- Where we are with regression so far
- Cost function and multiple linear regression
- SGD (Stochastic Gradient Descent)
- Normal equations solution

Pros and Cons

Gradient Descent

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

Normal Equations

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