

CS 260: Foundations of Data Science

Prof. Sara Mathieson

Fall 2021



Admin

- **Lab 3** due Tuesday night
 - Pair programming option
- **Note-taker:** Matthew

Outline for September 16

- Handout 5 (analytic solution example)
- SGD (Stochastic Gradient Descent)
- Handout 6 (SGD solution example)
- Analytic vs. SGD (pros and cons)

Outline for September 16

- Handout 5 (analytic solution example)
- SGD (Stochastic Gradient Descent)
- Handout 6 (SGD solution example)
- Analytic vs. SGD (pros and cons)

PLEASE LEAVE COMPUTERS ON

Handout 5

B must be square

$$BB^{-1} = B^{-1}B$$

- transpose: switch rows & cols

$$A = \begin{bmatrix} a & b & c \\ d & e & f \end{bmatrix}_{2 \times 3}, \quad A^T = \begin{bmatrix} a & d \\ b & e \\ c & f \end{bmatrix}_{3 \times 2}$$

$$\begin{aligned} A^T A &\Rightarrow 3 \times 3 \\ A A^T &\Rightarrow 2 \times 2 \end{aligned} \left. \vphantom{\begin{aligned} A^T A &\Rightarrow 3 \times 3 \\ A A^T &\Rightarrow 2 \times 2 \end{aligned}} \right\} \text{square!}$$

- inverse:

$$B^{-1}B = I \leftarrow \text{identity matrix ("1")}$$

$$B^{-1}B = \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix} \begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} ad-bc & db-db \\ -ac+ac & -cb+cd \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \star \text{identity}$$

$$I_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

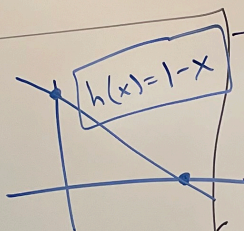
$$B = \begin{bmatrix} a & b \\ c & d \end{bmatrix}_{2 \times 2}$$
$$B^{-1} = \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

$$a \begin{bmatrix} 3 & 5 \\ -1 & 2 \end{bmatrix} = \begin{bmatrix} 3a & 5a \\ -a & 2a \end{bmatrix}$$

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

$$(X^T X)^{-1} = \left(\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \right)^{-1} \\ = \begin{bmatrix} 2 & 1 \\ 1 & 1 \end{bmatrix}^{-1} = \frac{1}{2-1} \begin{bmatrix} 1 & -1 \\ -1 & 2 \end{bmatrix}$$

$$\vec{w} = \begin{bmatrix} 1 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \vec{w}$$



~~$$(X^T)^{-1} X^T X \vec{w} = \vec{y}$$~~

solving for \vec{w}

X^T is not square
 \Rightarrow not invertible!

$$(X^T X)^{-1} X^T X \vec{w} = (X^T X)^{-1} X^T \vec{y}$$

$$\hat{\vec{w}} = (X^T X)^{-1} X^T \vec{y}$$

Annotations for the equation above:

- $(p+1) \times 1$ (pointing to $\hat{\vec{w}}$)
- $(p+1) \times n$ (pointing to $(X^T X)^{-1}$)
- $n \times 1$ (pointing to $X^T \vec{y}$)
- variance of X (pointing to $(X^T X)^{-1}$)
- cov of X & \vec{y} (pointing to $X^T \vec{y}$)
- $n \times (p+1)$ (pointing to $X^T \vec{y}$)

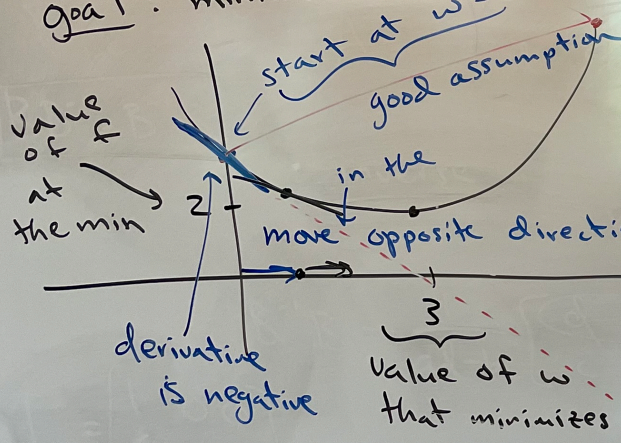
Outline for September 16

- Handout 5 (analytic solution example)
- **SGD (Stochastic Gradient Descent)**
- Handout 6 (SGD solution example)
- Analytic vs. SGD (pros and cons)

PLEASE LEAVE COMPUTERS ON

Gradient Descent

Goal: minimize a function



$$f(w) = (w-3)^2 + 2$$

$$f(w) = w^2 - 6w + 11$$

$$f'(w) = 2w - 6$$

$$w \leftarrow w - \alpha f'(w)$$

update
curr w
gradient
descent
step size

$$\alpha = 0.1$$

① (iteration)

$$w \leftarrow 0 - 0.1(2 \cdot 0 - 6)$$

$$w \leftarrow 0.6$$

②

$$w \leftarrow 0.6 - 0.1(2(0.6) - 6)$$

$$w \leftarrow 1.08$$

if $|f(\omega^t) - f(\omega^{t-1})| < \epsilon \rightarrow 1 \times 10^{-8}$
 \Rightarrow stop # converged

$$\alpha = \frac{1}{t}$$

Way of changing alpha adaptively

PLEASE LEAVE COMPUTERS ON

Gradient Descent for Linear Regression

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

want to minimize

wrt one weight

$$\frac{\partial J}{\partial w_j}$$

=

$$\sum_{i=1}^n (\vec{w} \cdot \vec{x}_i - y_i) x_{ij}$$

very slow when n is large

wrt one data point

$$\frac{\partial J_{x_i}}{\partial w_j} = (\vec{w} \cdot \vec{x}_i - y_i) x_{ij}$$

same for all w_j

=>

$$\nabla J_{x_i} = (\vec{w} \cdot \vec{x}_i - y_i) \vec{x}_i$$

$$\vec{w} \cdot \vec{x}_i = w_0 + w_1 x_{i1} + \dots + w_p x_{ip}$$

$\nabla J =$
gradient

$$\begin{bmatrix} \frac{\partial J}{\partial w_0} \\ \frac{\partial J}{\partial w_1} \\ \vdots \\ \frac{\partial J}{\partial w_p} \end{bmatrix}$$

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

derivative

derivative of "inside"

Stochastic Gradient Descent Algorithm

for epoch $t \dots$ until convergence

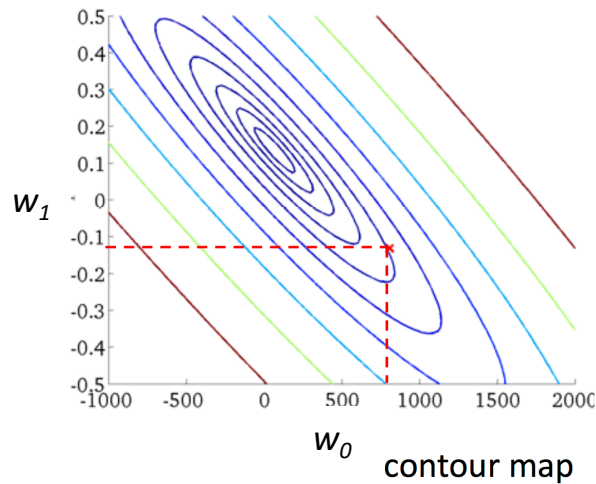
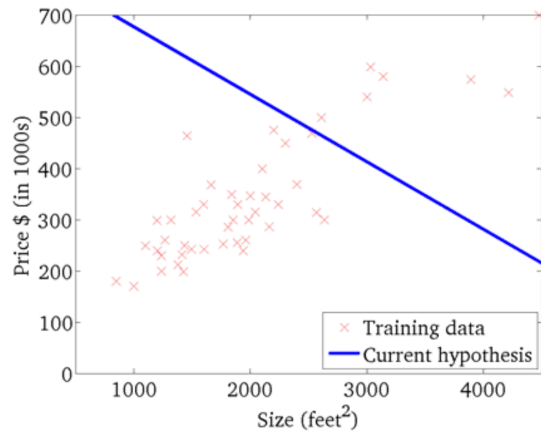
for $i = 1, 2, 3 \dots n$:

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

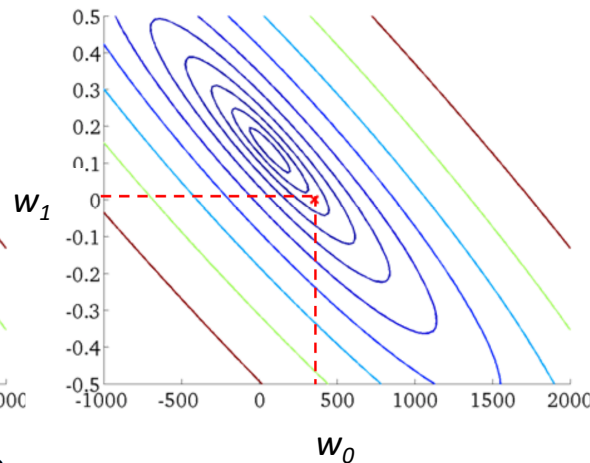
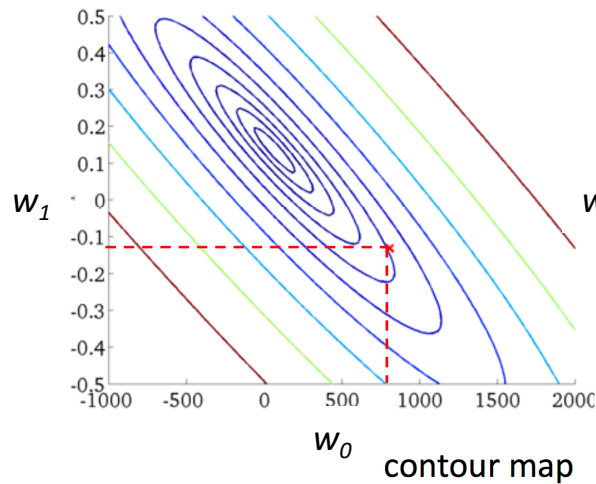
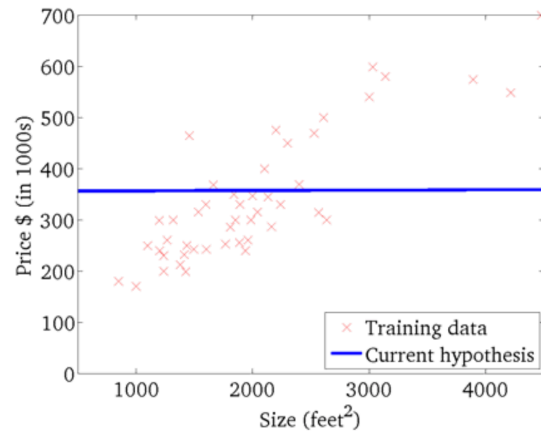
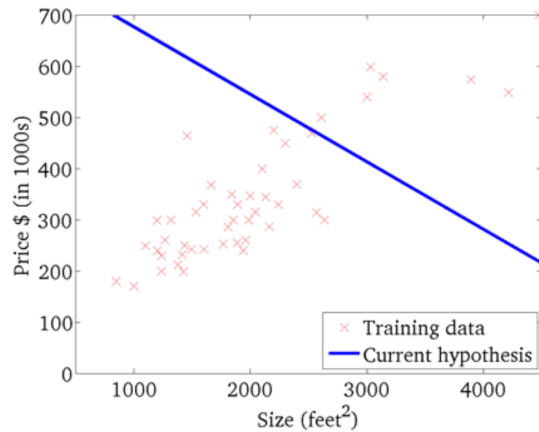
check for convergence: $|\mathcal{J}(\vec{w}^t) - \mathcal{J}(\vec{w}^{t-1})| < \epsilon$

gradient (of all weights)
wrt one data point \vec{x}_i

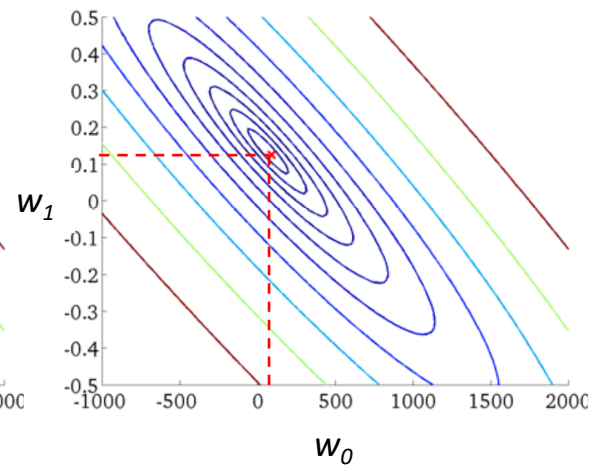
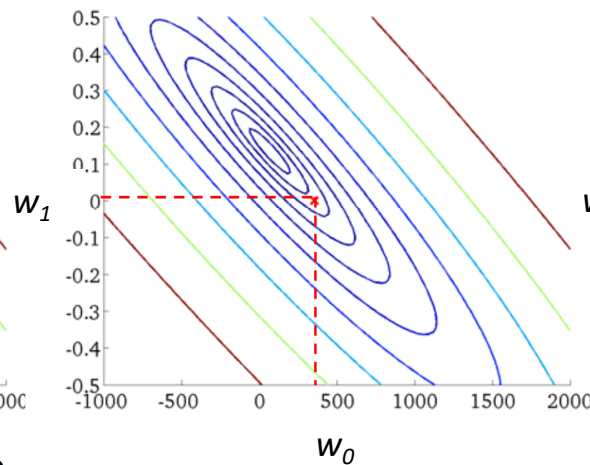
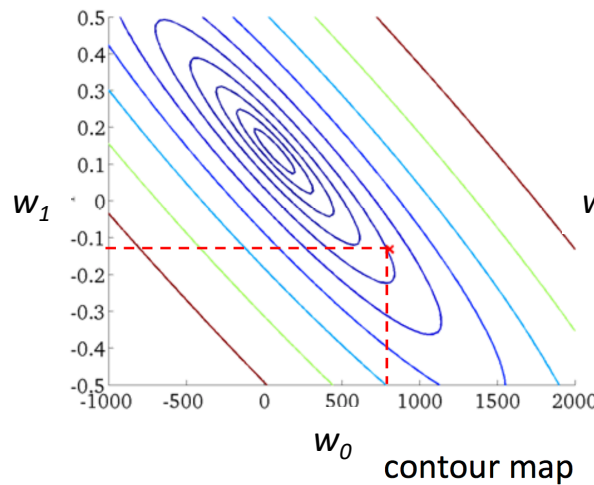
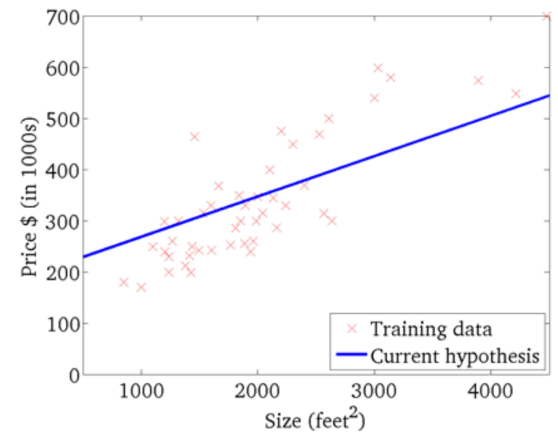
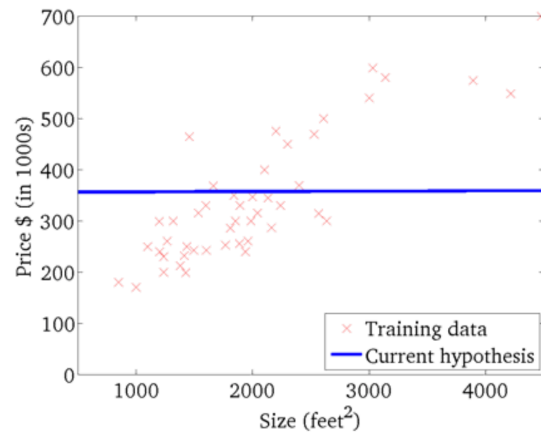
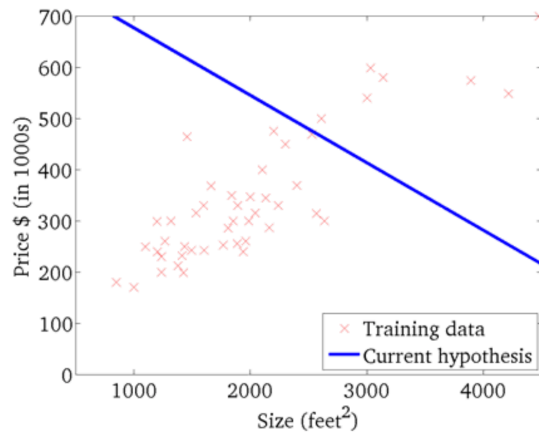
Linear Model and Cost Function J



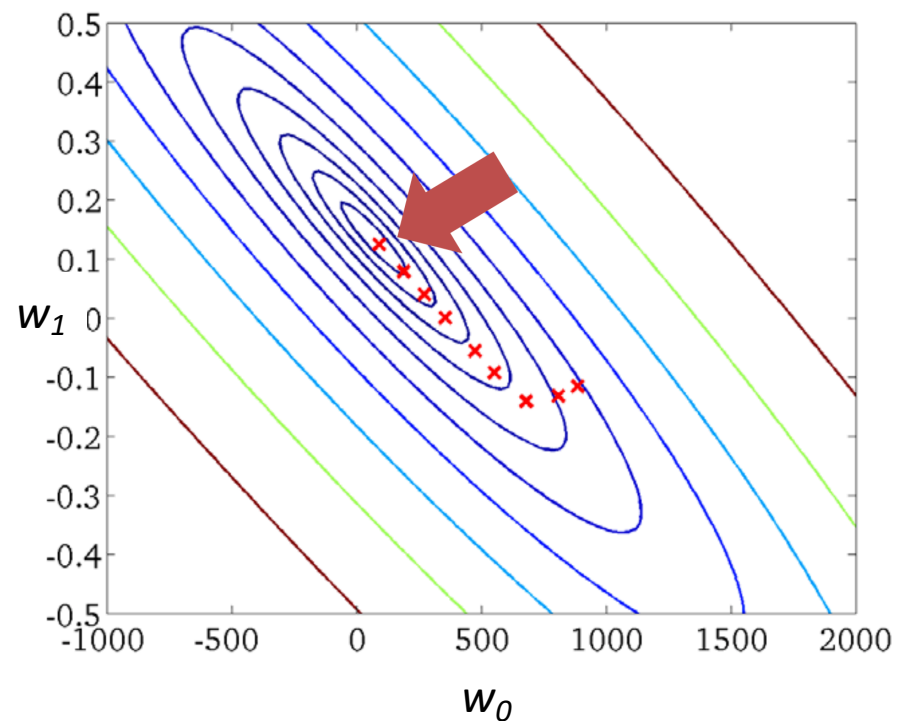
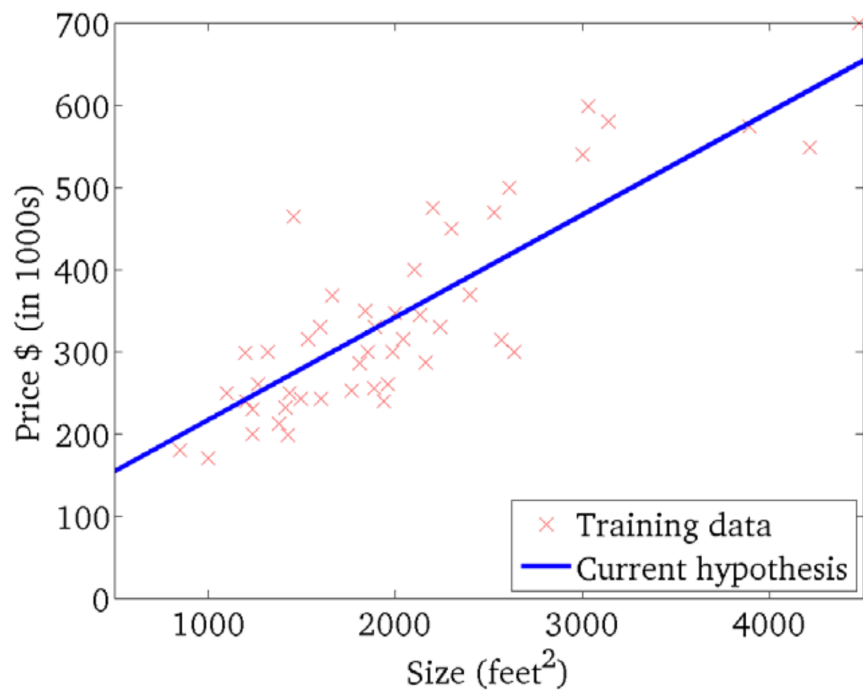
Linear Model and Cost Function J



Linear Model and Cost Function J



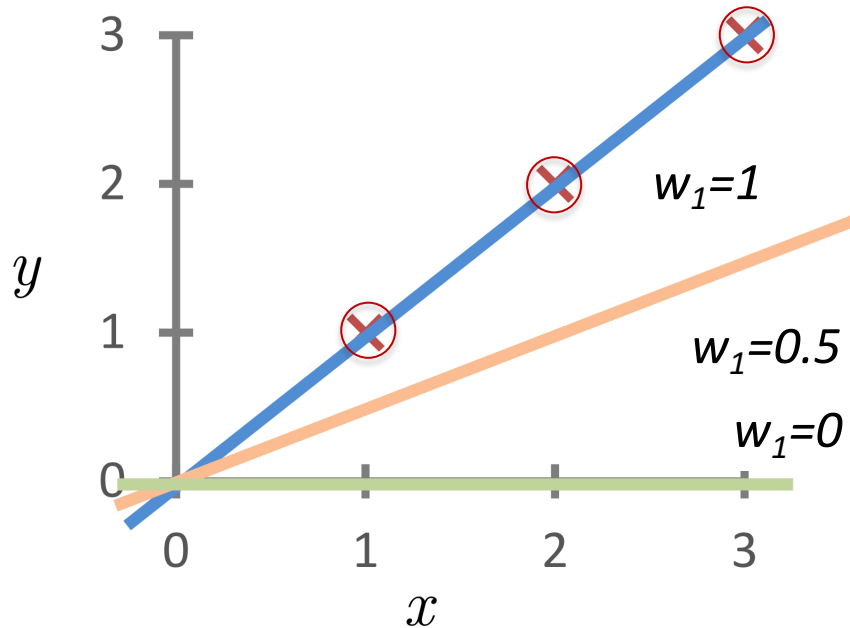
Gradient Descent: walking toward the minimum



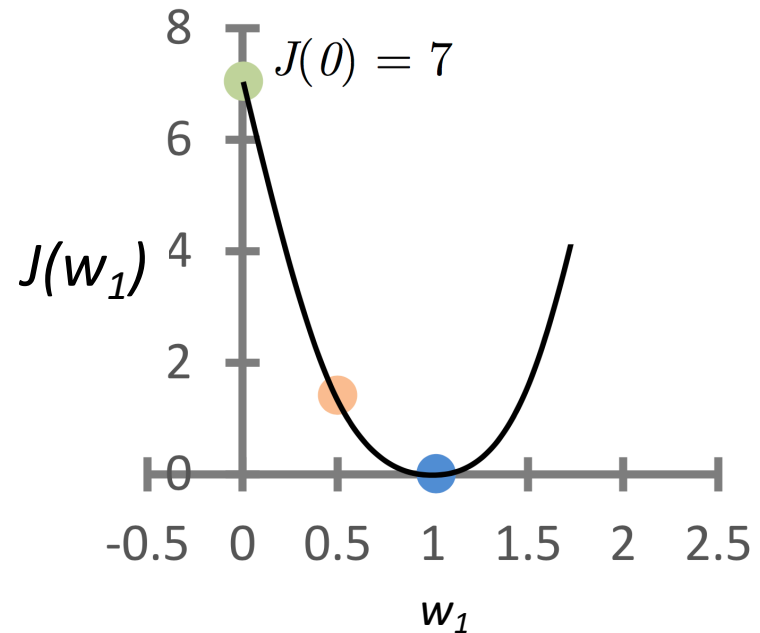
Cost Function (extra practice)

$$h_w(x) = w_1 x$$

(assume $w_0=0$ for this example)



$$J(w_1)$$

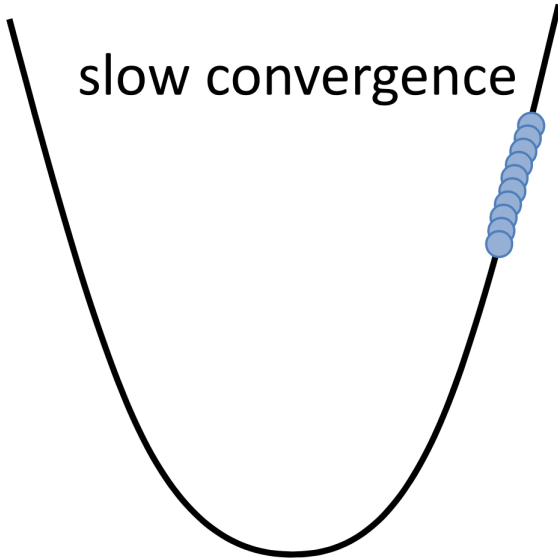


$$J(0.5) = \frac{1}{2} [(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2] = 1.75$$

Choosing the step size alpha

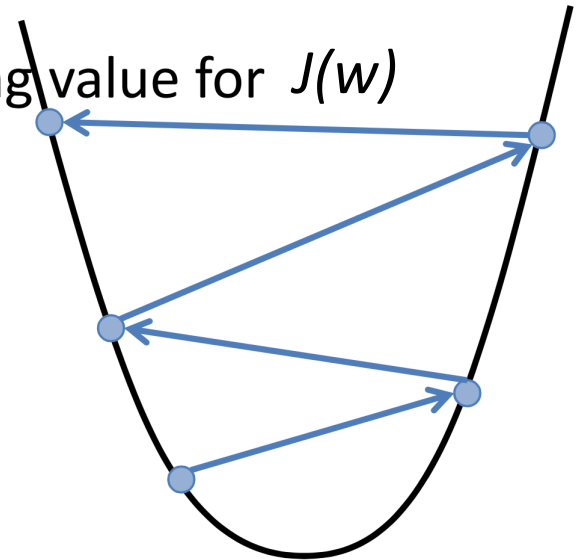
α too small

slow convergence



α too large

increasing value for $J(w)$



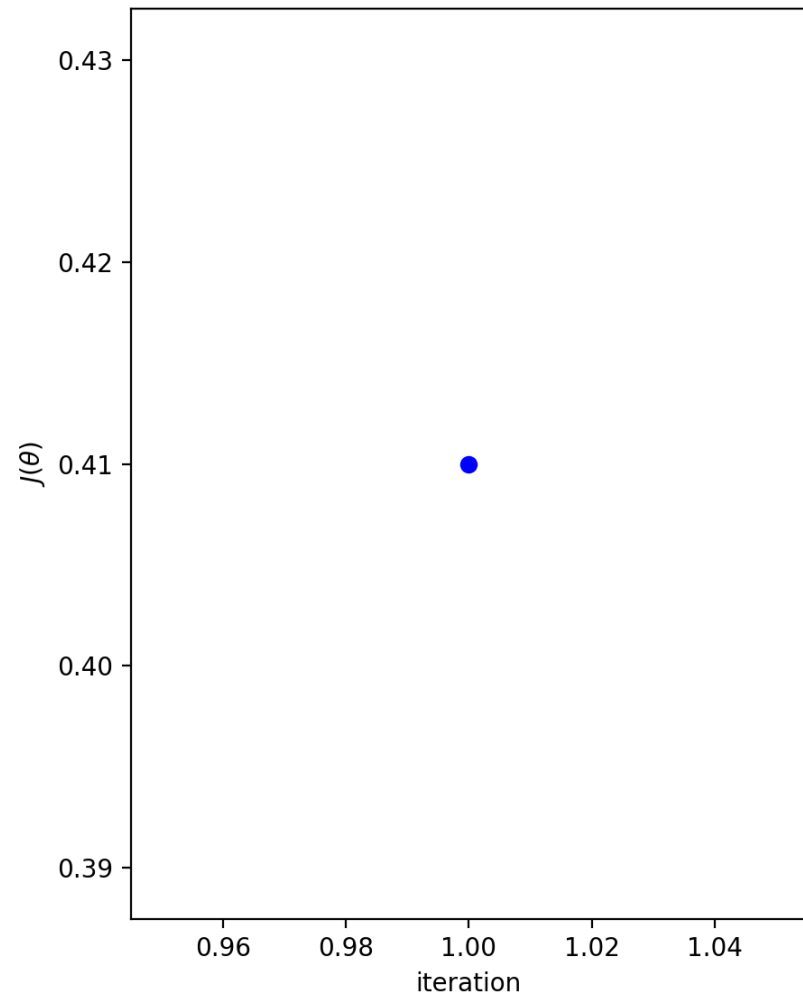
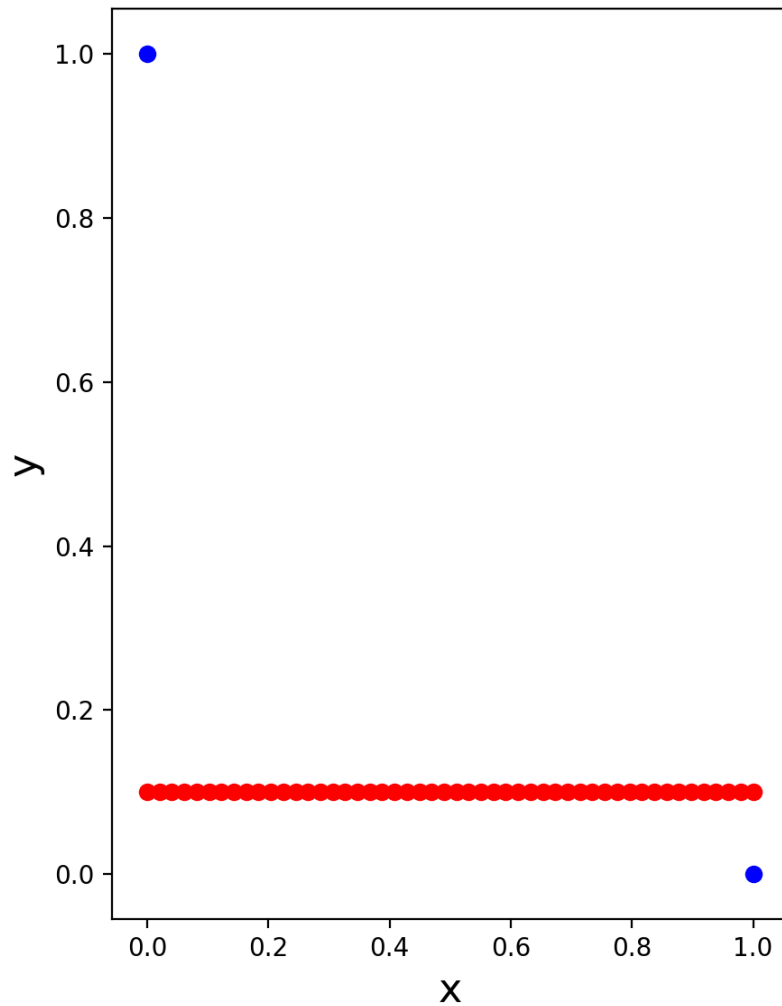
- may overshoot minimum
- may fail to converge (may even diverge)

SGD with our small dataset from the handouts

Note: this is with the original order of the points

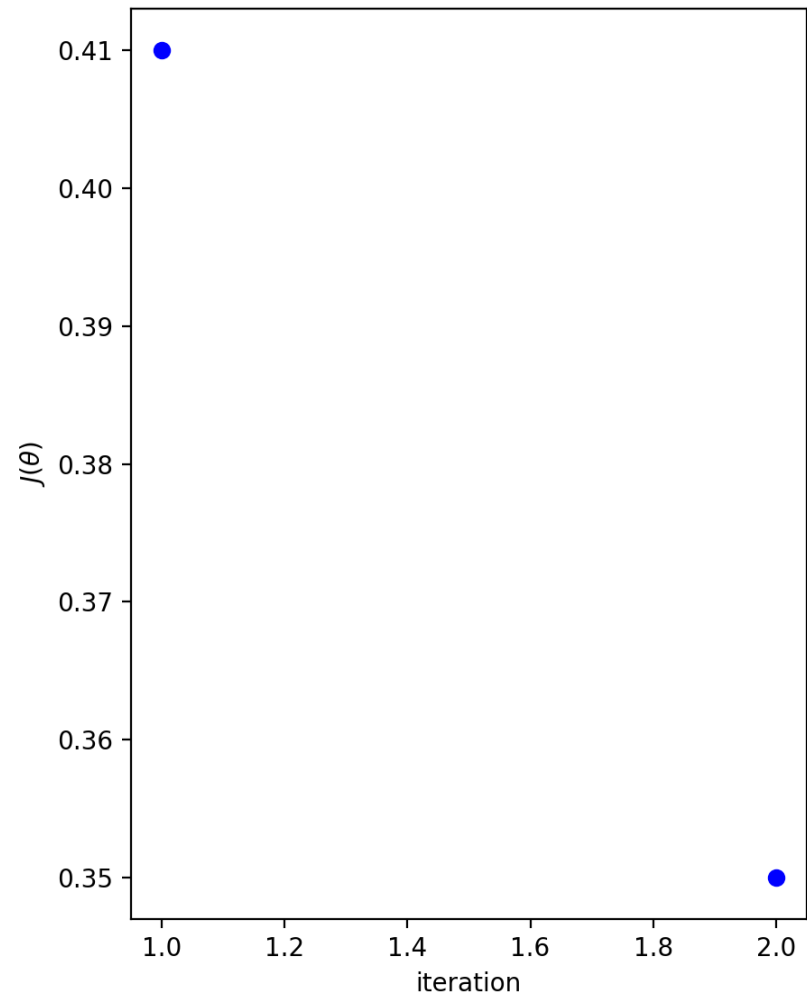
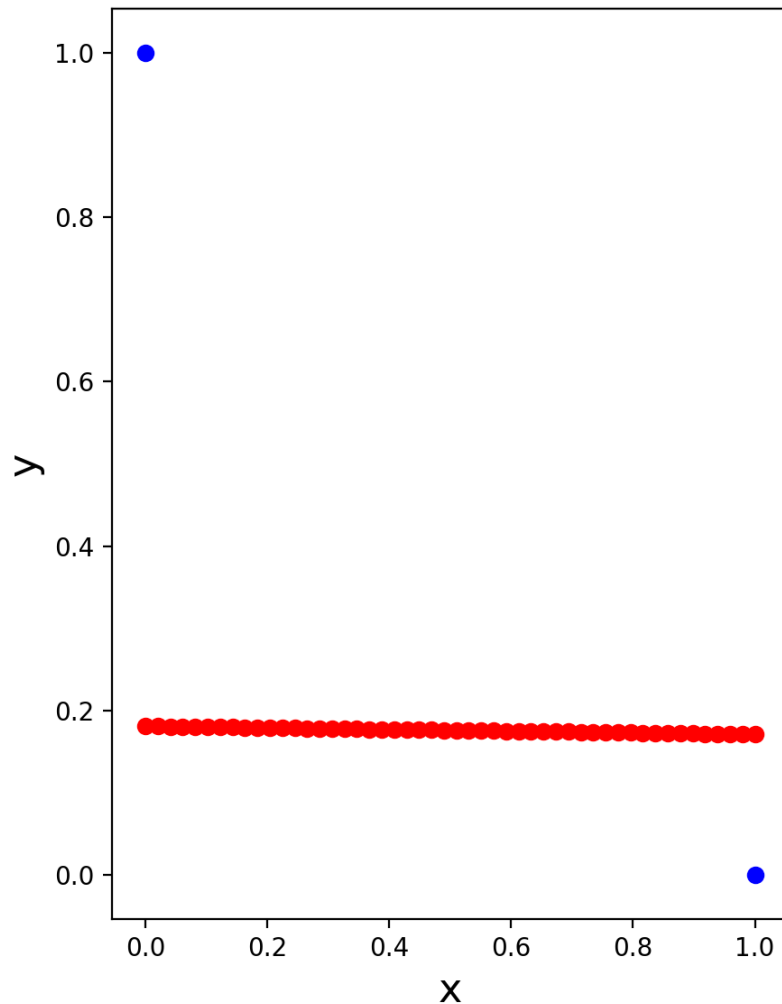
Small example, iteration 1

iteration: 1, cost: 0.410000



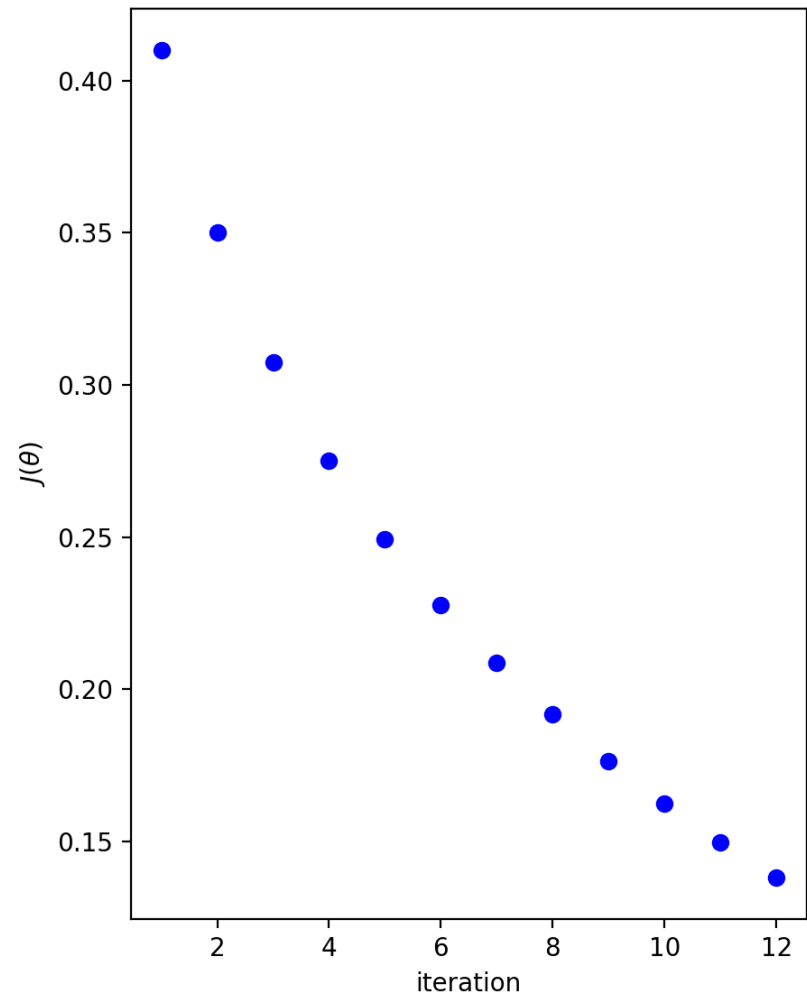
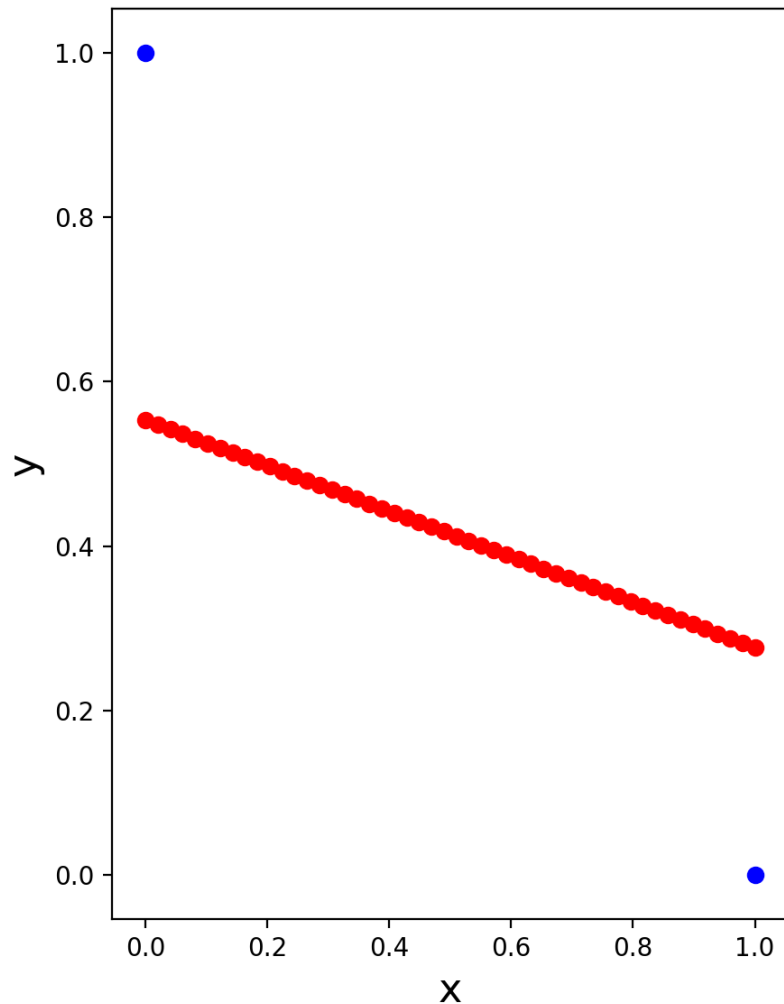
Small example, iteration 2

iteration: 2, cost: 0.350001



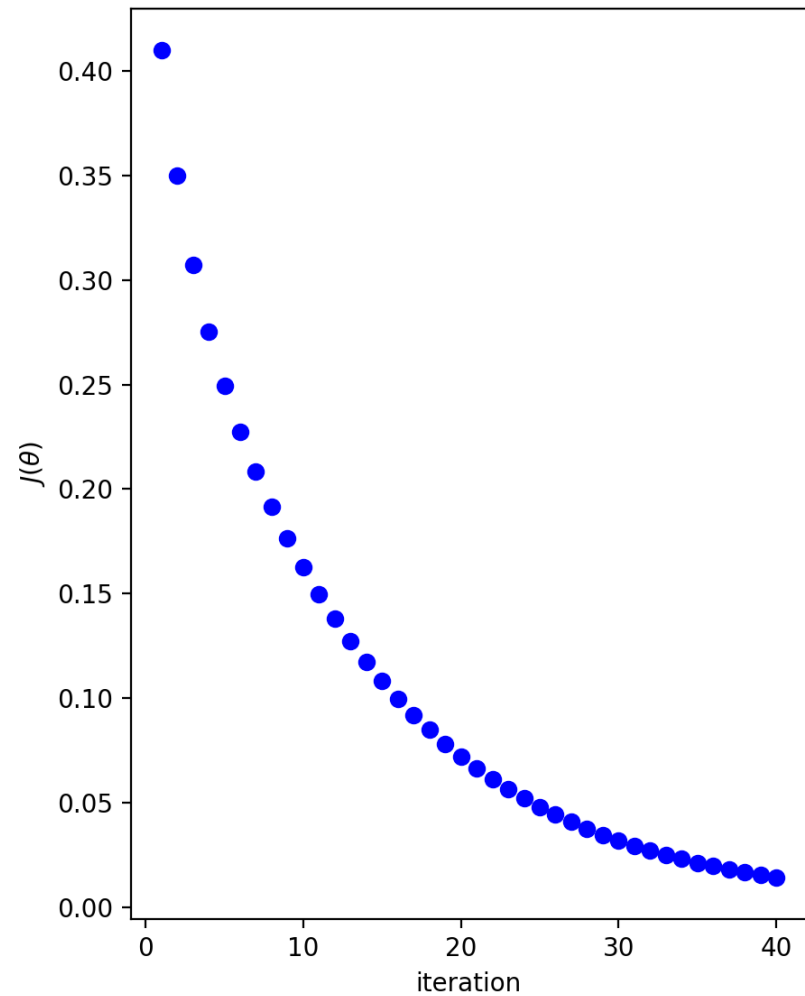
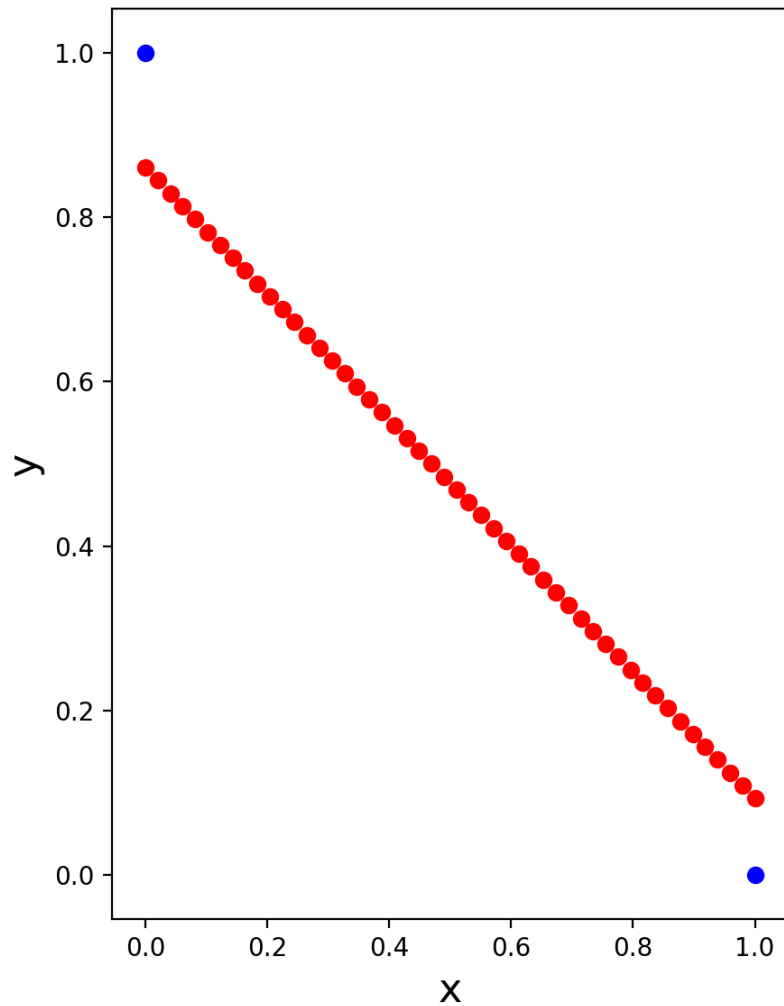
Small example, iteration 12

iteration: 12, cost: 0.138047



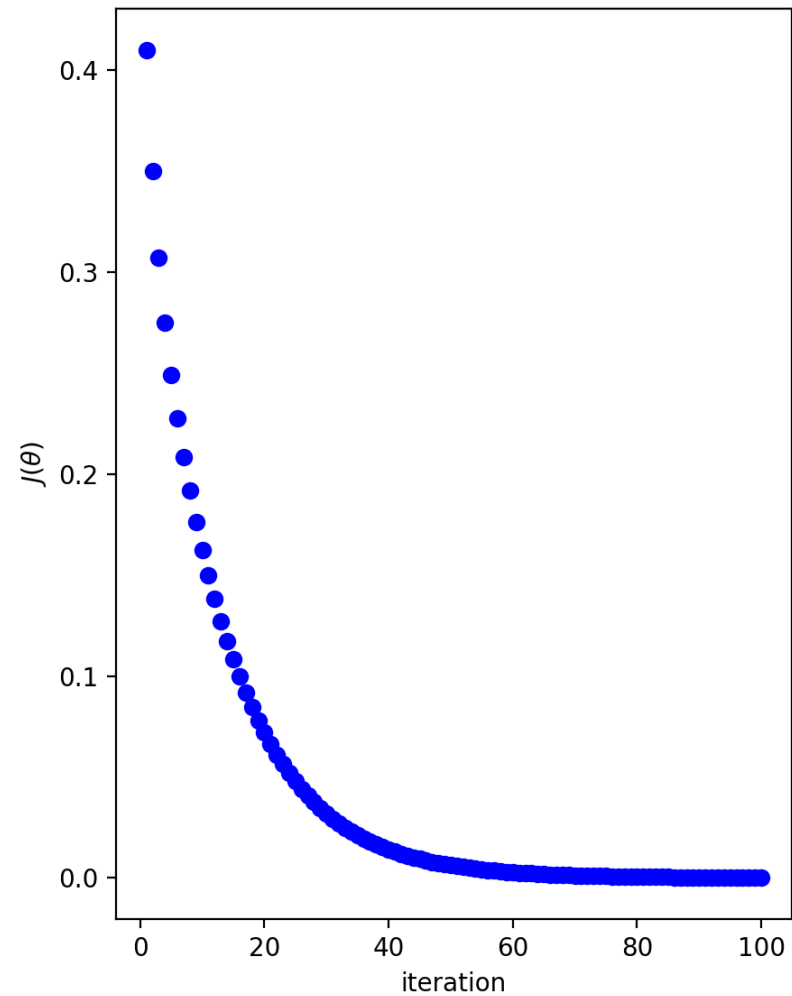
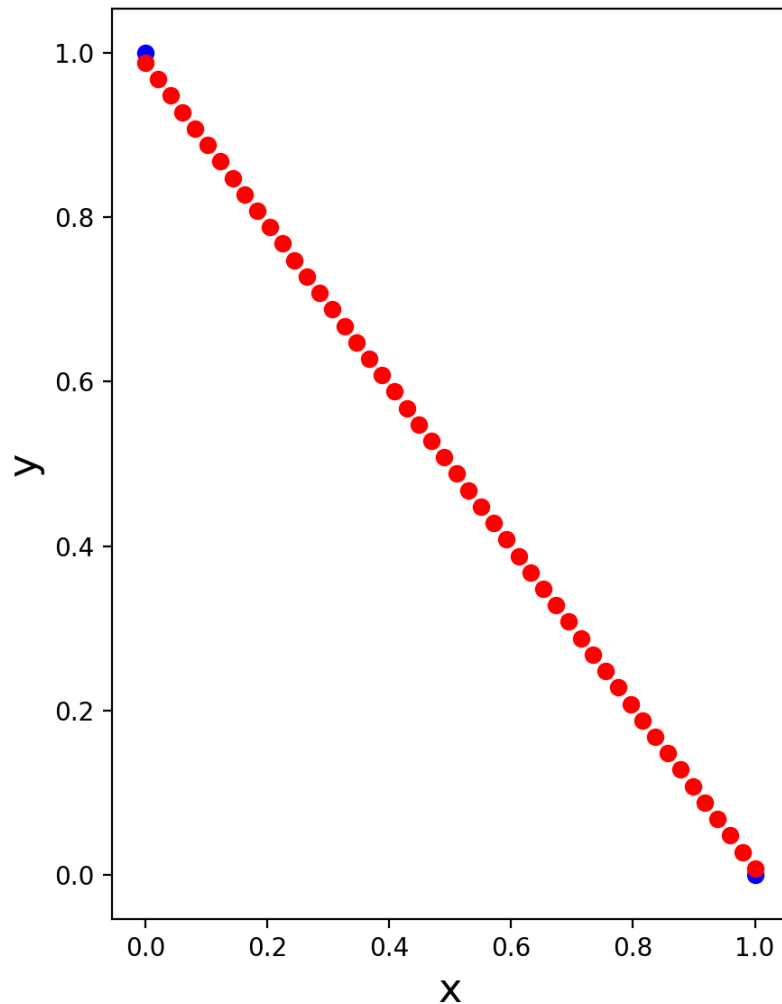
Small example, iteration 40

iteration: 40, cost: 0.014064



Small example, iteration 100

iteration: 100, cost: 0.000105



Outline for September 16

- Handout 5 (analytic solution example)
- SGD (Stochastic Gradient Descent)
- Handout 6 (SGD solution example)
- Analytic vs. SGD (pros and cons)

opt + parse

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

1st point

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

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

2nd point

$$\begin{aligned} \vec{w} &\leftarrow \begin{bmatrix} 0.1 \\ 0 \end{bmatrix} - 0.1 \left(\begin{bmatrix} 0.1 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} - 0 \right) \begin{bmatrix} 1 \\ 1 \end{bmatrix} \\ &\leftarrow \begin{bmatrix} 0.09 \\ -0.01 \end{bmatrix} \end{aligned}$$

fake ones

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

$$y = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\vec{x}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\vec{x}_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

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

$$y = 0.09 - 0.01x$$

$$(x_2, y_2) = (1, 0)$$

$$y = 0.1$$

$$y = 0 + 0x = 0$$

Outline for September 16

- Handout 5 (analytic solution example)
- SGD (Stochastic Gradient Descent)
- Handout 6 (SGD solution example)
- **Analytic vs. SGD (pros and cons)**

Pros and Cons

(Analytic Solution)

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 $O(p^3)$

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?

Extra topic: Polynomial Regression

- Can be thought of as regular linear regression with a change of basis

