

CS 360: Machine Learning

Prof. Sara Mathieson

Fall 2019



Admin

- Office hours **Friday**! 3-5pm (Hilles 110)
- Start Lab 3 in lab today (I will do some instruction during lab)
- Lab 3 due **Tuesday night**

Outline for September 19

- Reading quiz #3
- Simple linear regression
- SGD (Stochastic Gradient Descent)
- Normal equations solution
- Week 3 feedback

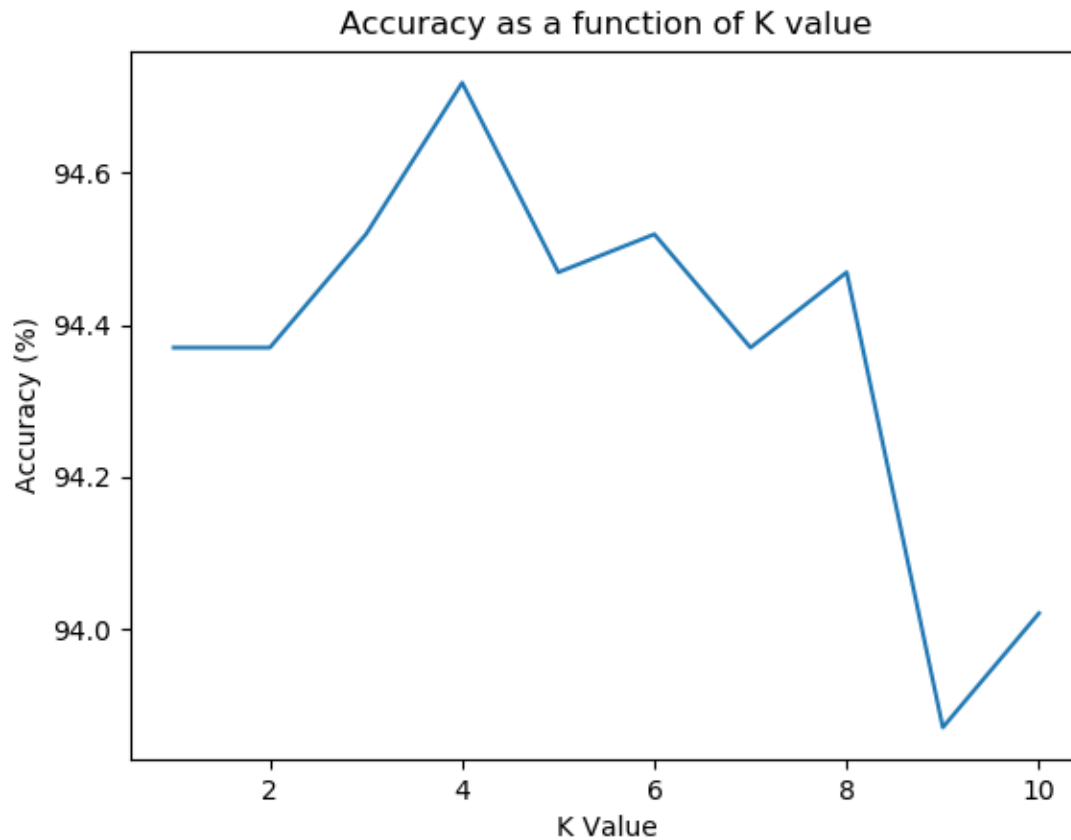
Lab 1 feedback

- Lab 1 grades posted (10pts for knn, 3 for math)
- Keep large data out of repos (unless it is given as part of the starter)
- Make sure code is in functions and function decomposition makes sense
- Keep lines less than 80 characters
- Official python style for methods/functions
 - `my_function`
 - ~~`MyFunction`~~ (uppercase is reserved for classes!)
- Some people are new to Python: I am happy to take some time during lab to go over syntax/common issues

Lab 1 Extensions

- **Ahmed**: implemented a weighting voting scheme for multi-class classification
- Amazing style: **Fiona, Jason**
- One of the fastest solutions: **Emile**
- **Jiaping**: command line options and automatic plot creation

Lab 1 Extensions: multi-class



Brian and Gareth created similar plots

Lab 1: how to make KNN faster?

- 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 for September 19

- Reading quiz #3
- Simple linear regression
- SGD (Stochastic Gradient Descent)
- Normal equations solution
- Week 3 feedback

Reading Quiz #3

1. Name one goal of linear regression. Why would we want to fit a linear model?

Reading Quiz #3

1. Name one goal of linear regression. Why would we want to fit a linear model?

Learn which features contribute to the response

Predict future responses based on new features

2. In *simple* linear regression, how many features/predictors does each example have? How many parameters does our linear model have in this case?

- # features (p) =
- # model params =

Reading Quiz #3

1. Name one goal of linear regression. Why would we want to fit a linear model?

Learn which features contribute to the response

Predict future responses based on new features

2. In *simple* linear regression, how many features/predictors does each example have? How many parameters does our linear model have in this case?

- # features (p) = 1
- # model params = 2

3. What loss function are we trying to minimize in classical linear regression?

Reading Quiz #3

1. Name one goal of linear regression. Why would we want to fit a linear model?

Learn which features contribute to the response

Predict future responses based on new features

2. In *simple* linear regression, how many features/predictors does each example have? How many parameters does our linear model have in this case?

- # features (p) = 1
- # model params = 2

3. What loss function are we trying to minimize in classical linear regression?

squared error

4. Does a closed-form (analytic) solution exist for the parameters of a linear model?

Reading Quiz #3

1. Name one goal of linear regression. Why would we want to fit a linear model?

Learn which features contribute to the response

Predict future responses based on new features

2. In *simple* linear regression, how many features/predictors does each example have? How many parameters does our linear model have in this case?

- # features (p) = 1
- # model params = 2

3. What loss function are we trying to minimize in classical linear regression?

squared error

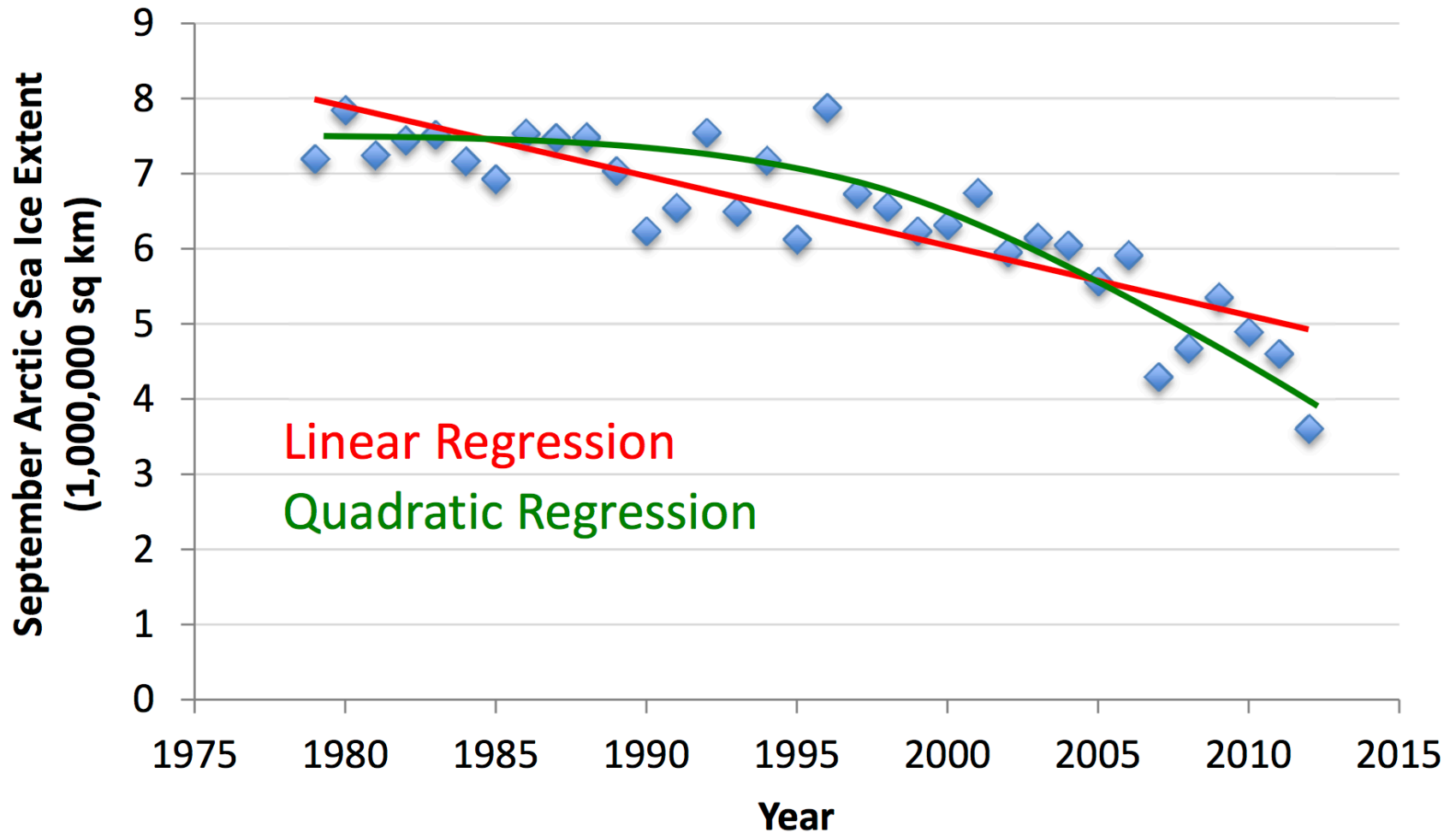
4. Does a closed-form (analytic) solution exist for the parameters of a linear model?

Yes!

Outline for September 19

- Reading quiz #3
- Simple linear regression
- SGD (Stochastic Gradient Descent)
- Normal equations solution
- Week 3 feedback

Regression Example



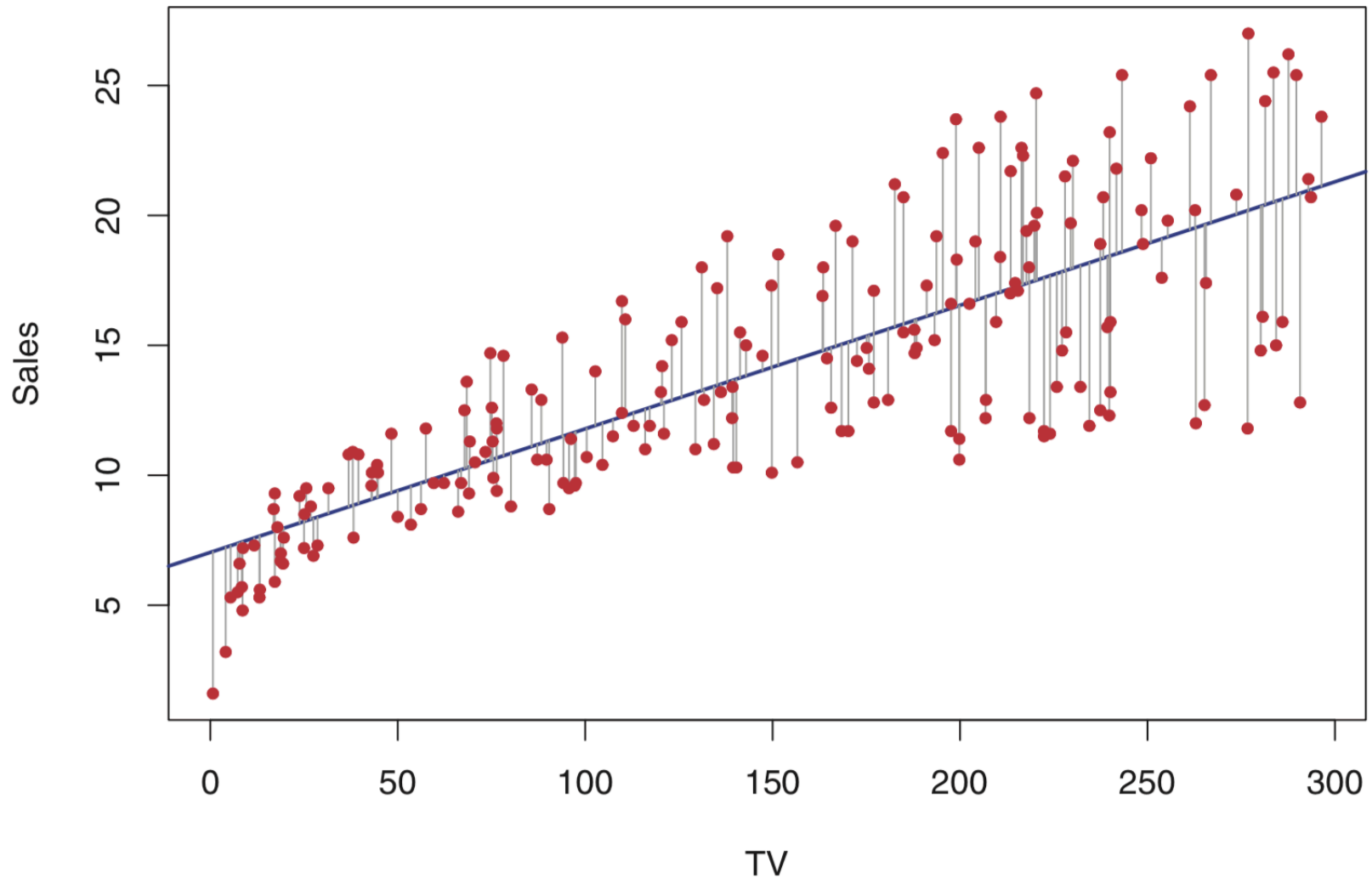
Goals of Inference

- 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 ?

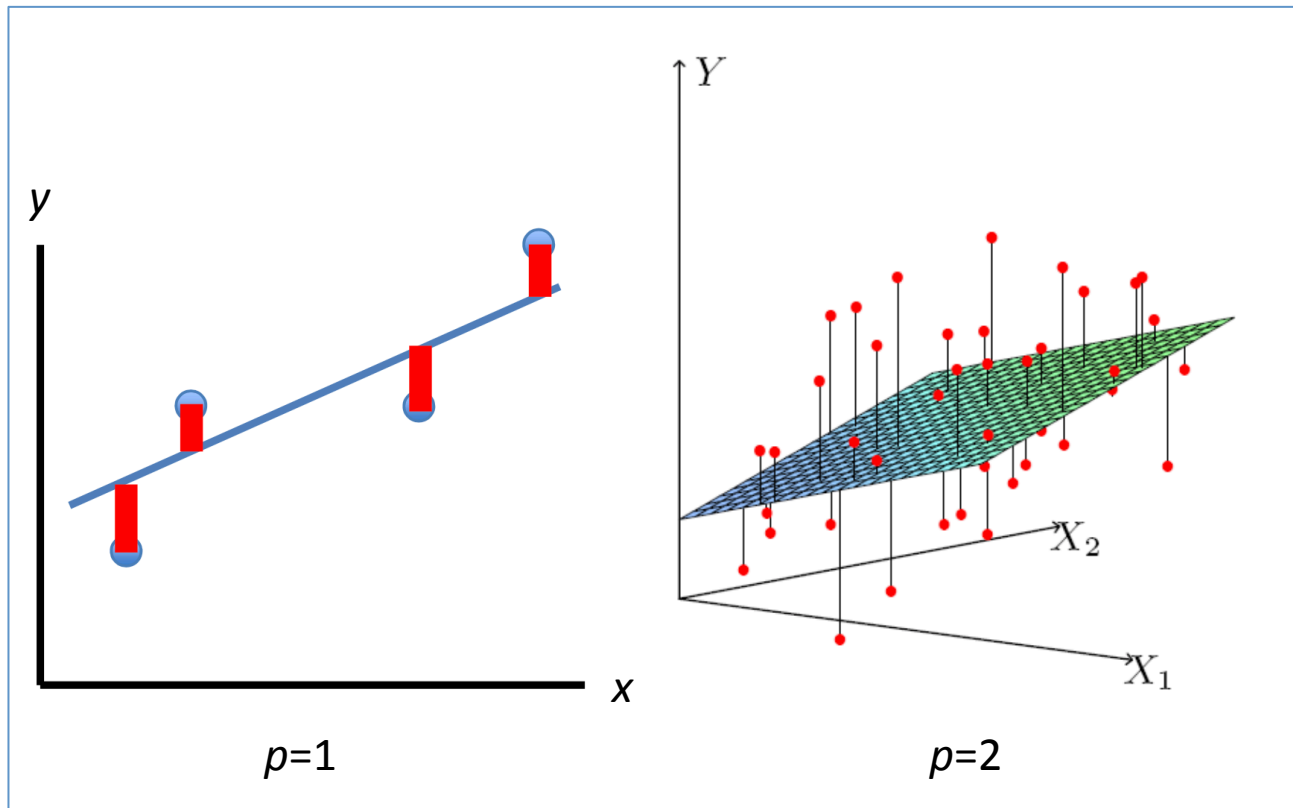
Linear Regression so far...

- 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 squared errors) or *SSE* (sum of squared errors)

Example: predict sales from TV advertising budget



Cost Function: sum of squared errors



"Simple" linear regression

$p=1$

"best fit" $\left\{ \begin{aligned} \hat{f}(x) &= \hat{w}_0 + \hat{w}_1 x \\ &= \hat{y} \\ &= h_{\vec{w}}(x) \end{aligned} \right.$

Goal: minimize

$$MSE = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

RSS
SSE

Cost function:

$$J(\vec{w}) = \frac{1}{2} \sum_{i=1}^n (\underbrace{w_0}_{\text{bias}} + \underbrace{w_1 x_i}_{\text{weight}} - y_i)^2$$

make derivative
w.r.t
gradient

we go in
this
direction

gradient

$$\nabla J = \begin{bmatrix} \frac{\partial J}{\partial w_0} \\ \frac{\partial J}{\partial w_1} \end{bmatrix}$$

partial
derivatives

$$\frac{\partial J}{\partial w_0} = \sum_{i=1}^n (w_0 + w_1 x_i - y_i) = 0$$

solve for
 w_0

$$\frac{n w_0}{n} = -w_1 \underbrace{\left(\frac{1}{n} \sum_{i=1}^n x_i \right)}_{\text{avg/mean } \bar{x}} + \underbrace{\left(\frac{1}{n} \sum_{i=1}^n y_i \right)}_{\bar{y} \text{ (mean of } y\text{'s)}}$$

$$\hat{\omega}_0 = \bar{y} - \omega_1 \bar{x}$$

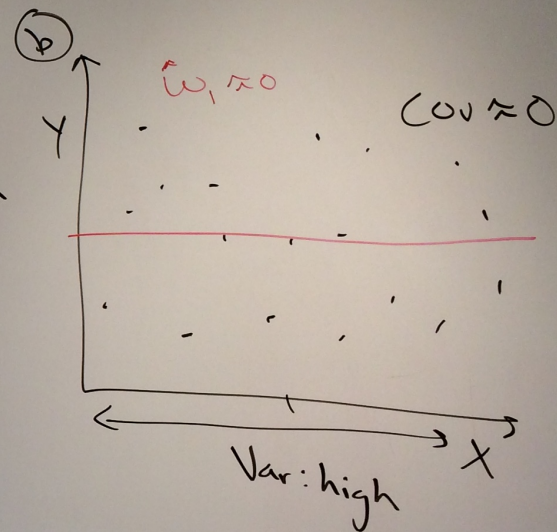
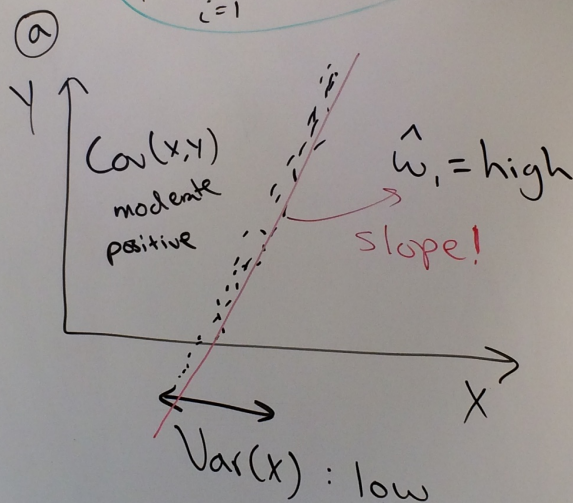
$$\frac{\partial J}{\partial \omega_1} = \sum_{i=1}^n (\hat{\omega}_0 + \omega_1 x_i - y_i) x_i$$

$$\Rightarrow \omega_1 = \frac{\sum_{i=1}^n (x_i y_i - x_i \bar{y})}{\sum_{i=1}^n (x_i^2 - x_i \bar{x})}$$

note: $\sum_{i=1}^n (\bar{x} \bar{y} - y_i \bar{x}) = 0$

$$\sum_{i=1}^n (\bar{x}^2 - x_i \bar{x}) = 0$$

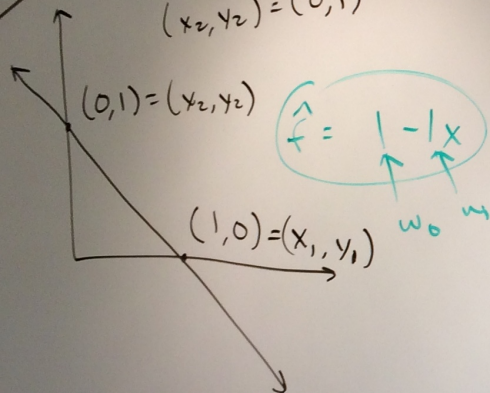
$$\hat{\omega}_1 = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$



Handout 5

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

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



$$\hat{w}_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

$$\bar{x} = \frac{1}{2}$$

$$\bar{y} = \frac{1}{2}$$

$$= \frac{\frac{1}{2} \left[\left(1 - \frac{1}{2}\right) \left(0 - \frac{1}{2}\right) + \left(0 - \frac{1}{2}\right) \left(1 - \frac{1}{2}\right) \right]}{\frac{1}{2} \left[\left(1 - \frac{1}{2}\right)^2 + \left(0 - \frac{1}{2}\right)^2 \right]}$$

$$\hat{w}_1 = -1$$

$$w_0 = \frac{1}{2} - (-1) \cdot \frac{1}{2}$$

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

Outline for September 19

- Reading quiz #3
- Simple linear regression
- SGD (Stochastic Gradient Descent)
- Normal equations solution
- Week 3 feedback

$p > 1$

multiple linear regression

$$\hat{y} = \hat{f}(\vec{x}) = h_{\vec{w}}(\vec{x}) = w_0 \overset{x_0=1}{\bullet} + w_1 x_1 + w_2 x_2 \dots + w_p x_p$$
$$= \vec{w}^T \vec{x} = \vec{w} \cdot \vec{x}$$

dot product.

$$\vec{a} \cdot \vec{b}$$
$$= \sum_{i=1}^{\text{len}} a_i b_i$$

Cost function (minimize!)

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

row x col.

$$\underbrace{[w_0 \ w_1 \ \dots \ w_p]}_{\vec{w}^T \quad 1 \times (p+1)} \underbrace{\begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_p \end{bmatrix}}_{(p+1) \times 1} = \text{Scalar}$$

Method 1

Stochastic Gradient Descent (SGD)

$$\frac{\partial J}{\partial w_j} = \sum_{i=1}^n (\vec{w}^T \vec{x}_i - y_i) x_{ij}$$

slow!! when n is large.

algorithm

while not converged:

shuffle all data points

for $i = 1 \dots n$ # go through one at a time!

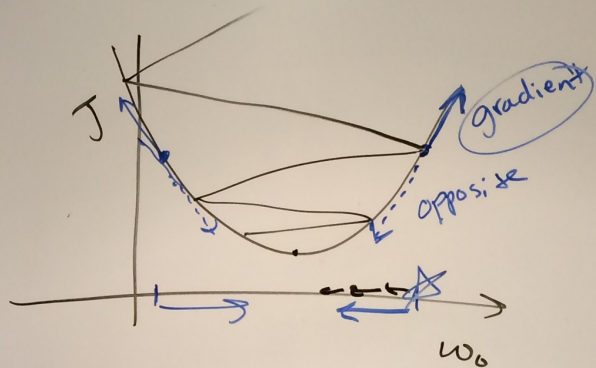
$$\vec{w} \leftarrow \vec{w} - \eta (\vec{w}^T \vec{x}_i - y_i) \vec{x}_i$$

all weights!

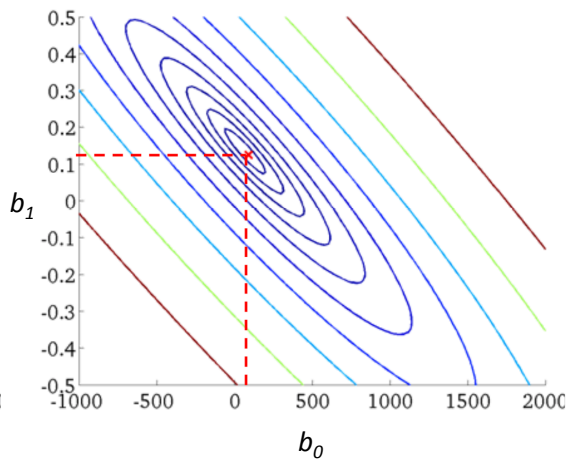
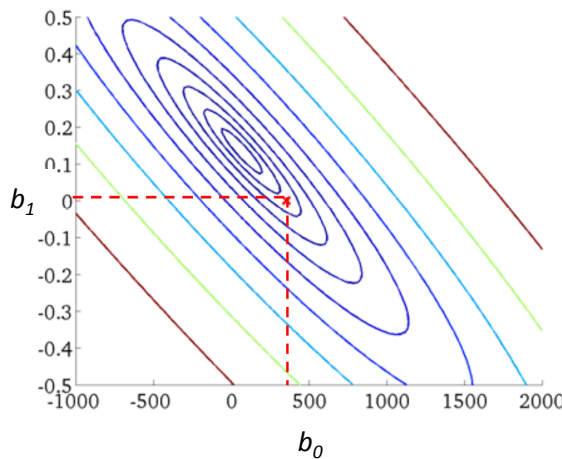
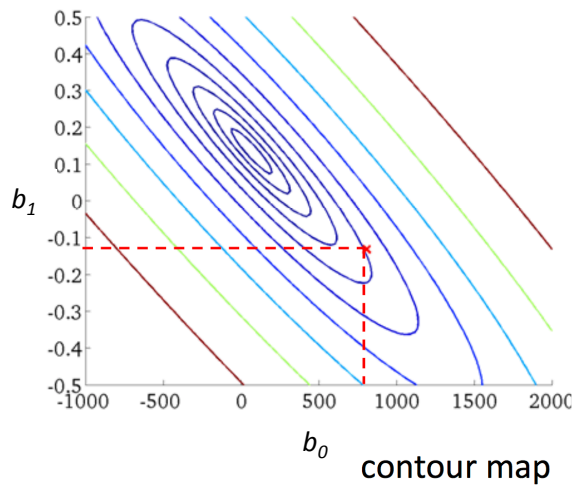
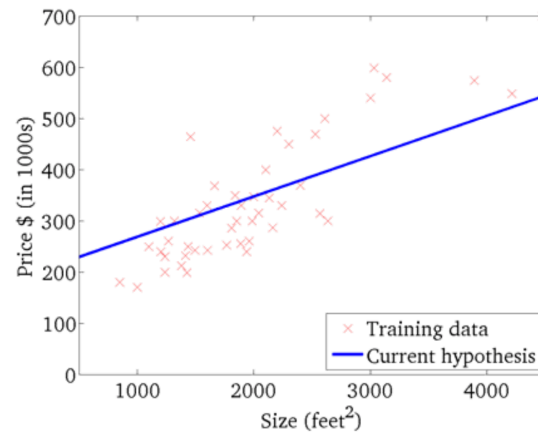
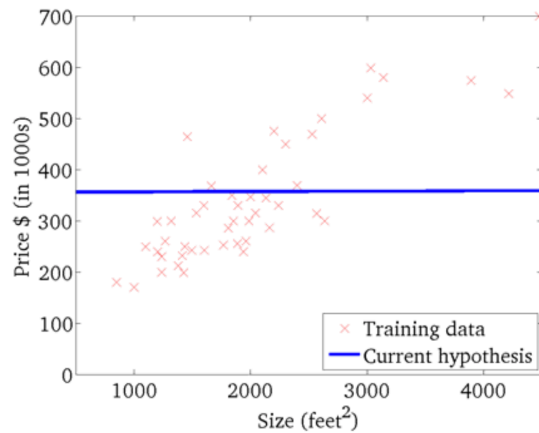
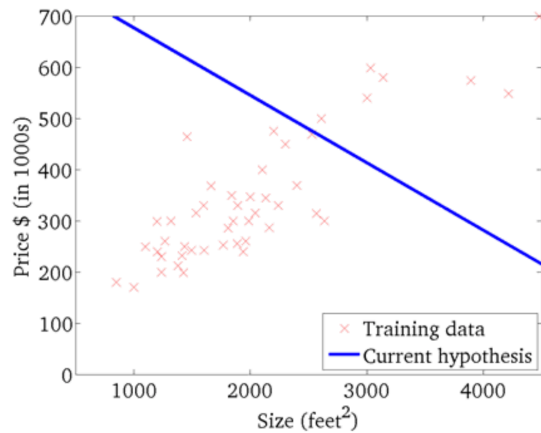
opposite

step
size.

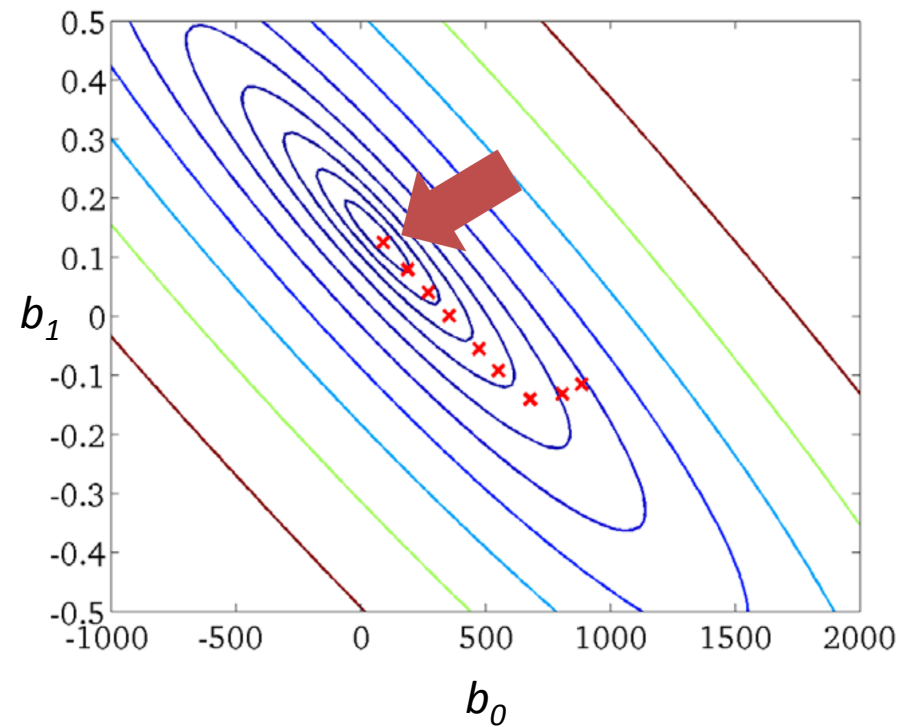
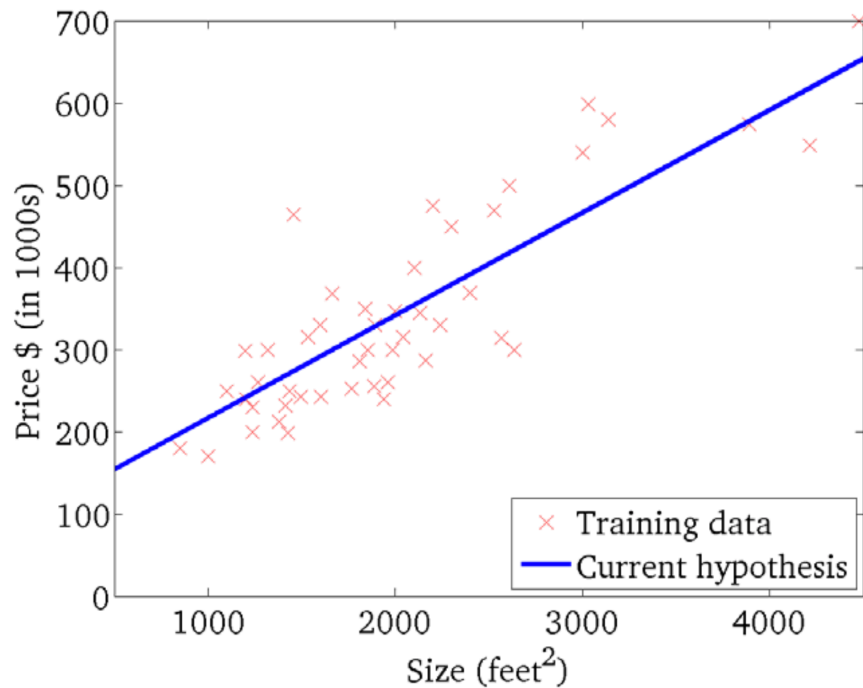
gradient of J
wrt one data point.



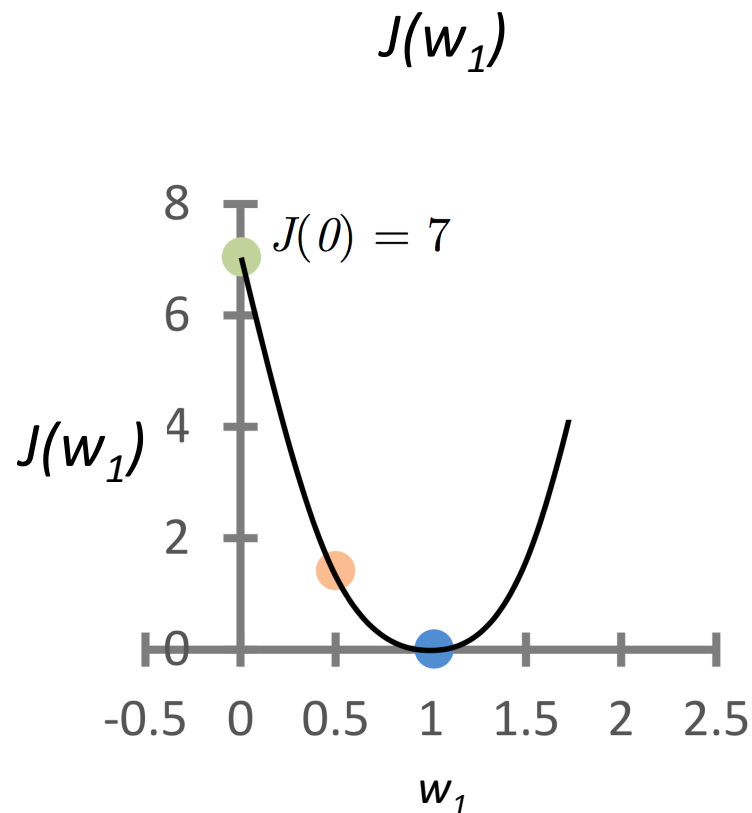
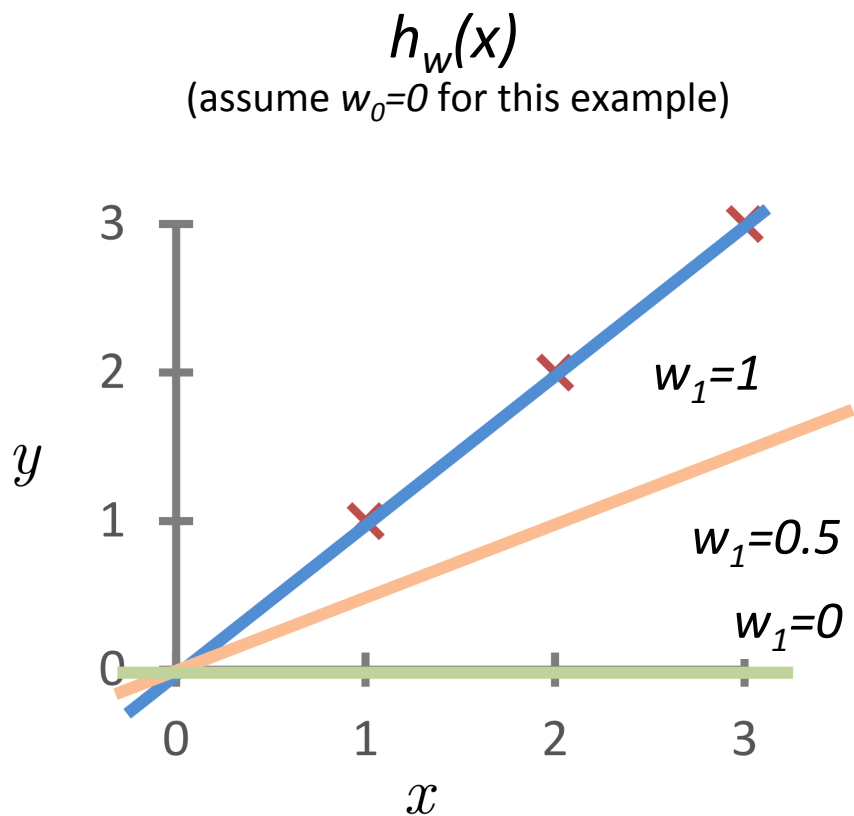
Cost Function



Gradient Descent: walking toward the minimum



Cost Function (extra practice)

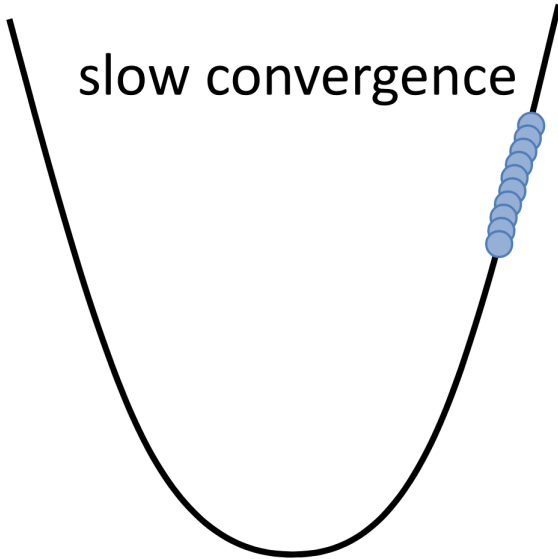


$$J(0.5) = \frac{1}{2} [(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2] = 1.75$$

Choosing step size α

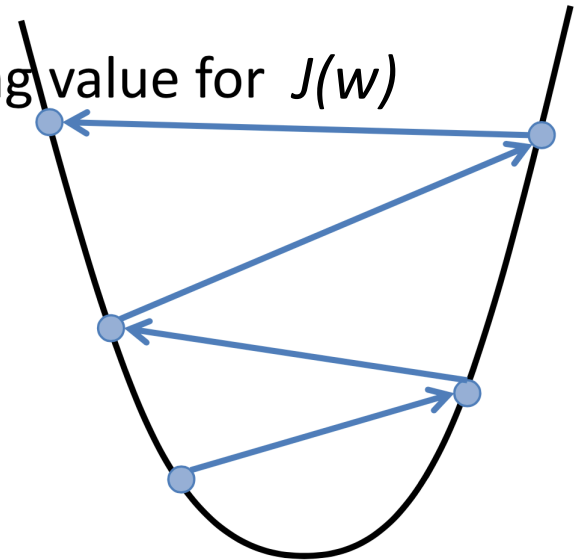
α too small

slow convergence



α too large

increasing value for $J(w)$



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

Outline for September 19

- Reading quiz #3
- Simple linear regression
- SGD (Stochastic Gradient Descent)
- Normal equations solution
- Week 3 feedback

Normal Equations

cost function

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

Scalar

pred

truth

technically,

$$J(\vec{w}) = \frac{1}{2} (X\vec{w} - \vec{y})^T (X\vec{w} - \vec{y})$$

$(1 \times n)$ $(n \times 1)$

X $n \times (p+1)$

\vec{w} $(p+1) \times 1$

\vec{y} $n \times 1$

$1 \cdot w_0 + x_1 \cdot w_1 + x_2 \cdot w_2 \dots x_p \cdot w_p$

take derivative
 \Rightarrow set equal to 0

normal equation

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

$O(p^3) \Rightarrow$ slow.

closed form solution.

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{\det} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

det.

$$\begin{bmatrix} 1 & 3 \\ 1 & 3 \end{bmatrix}^{-1} = \frac{1}{3-3} ()$$

AHH!!!

polynomial regression

d degree polynomial.

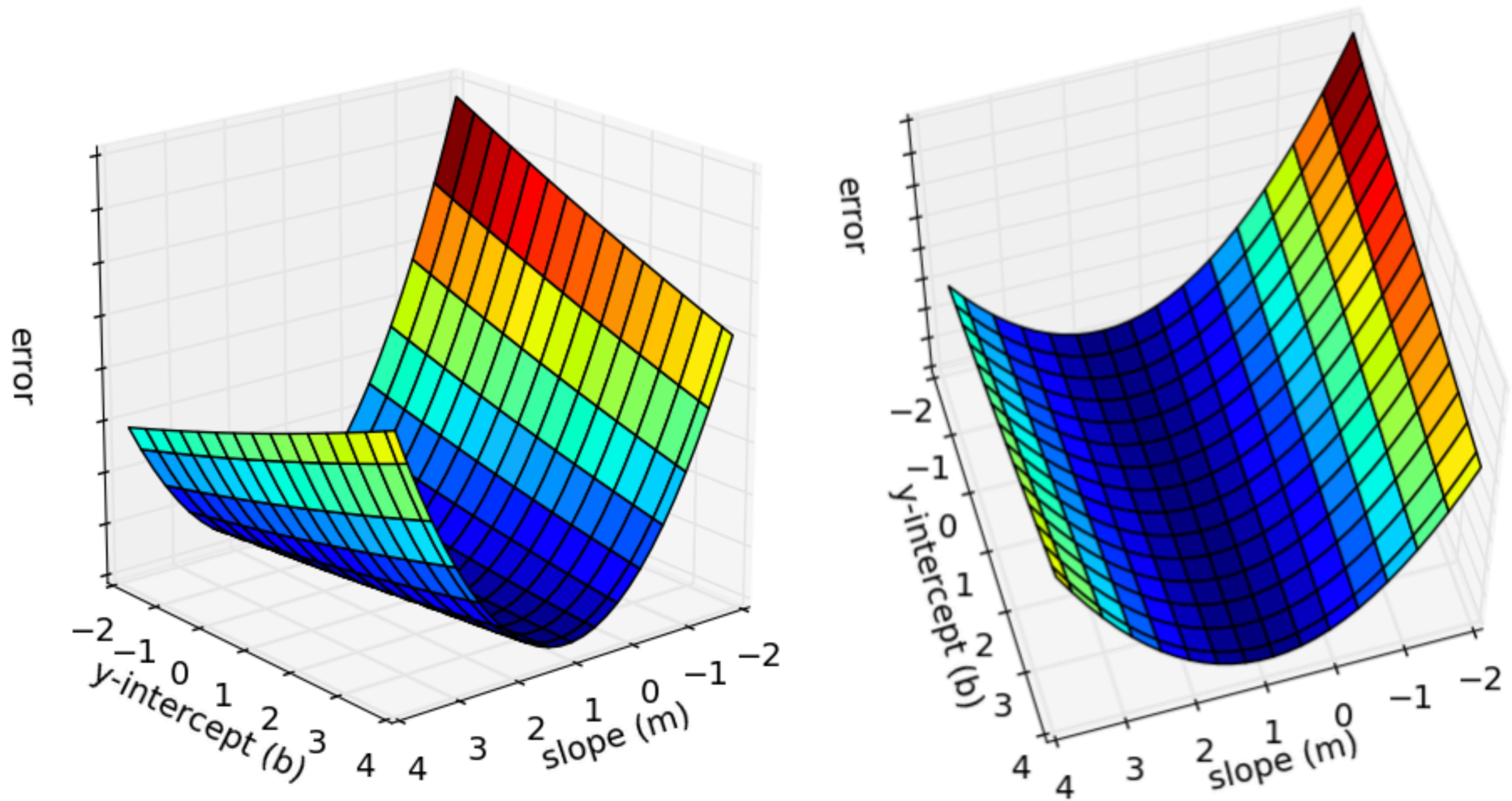
$p=1$

$$h_{\vec{w}}(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3 + \dots + w_d x^d$$

$$X = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 & \dots & x_1^d \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \end{bmatrix}$$

$$J = 0.41$$

Error as a function of slope & y-intercept



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

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 $O(p^3)$

Outline for September 19

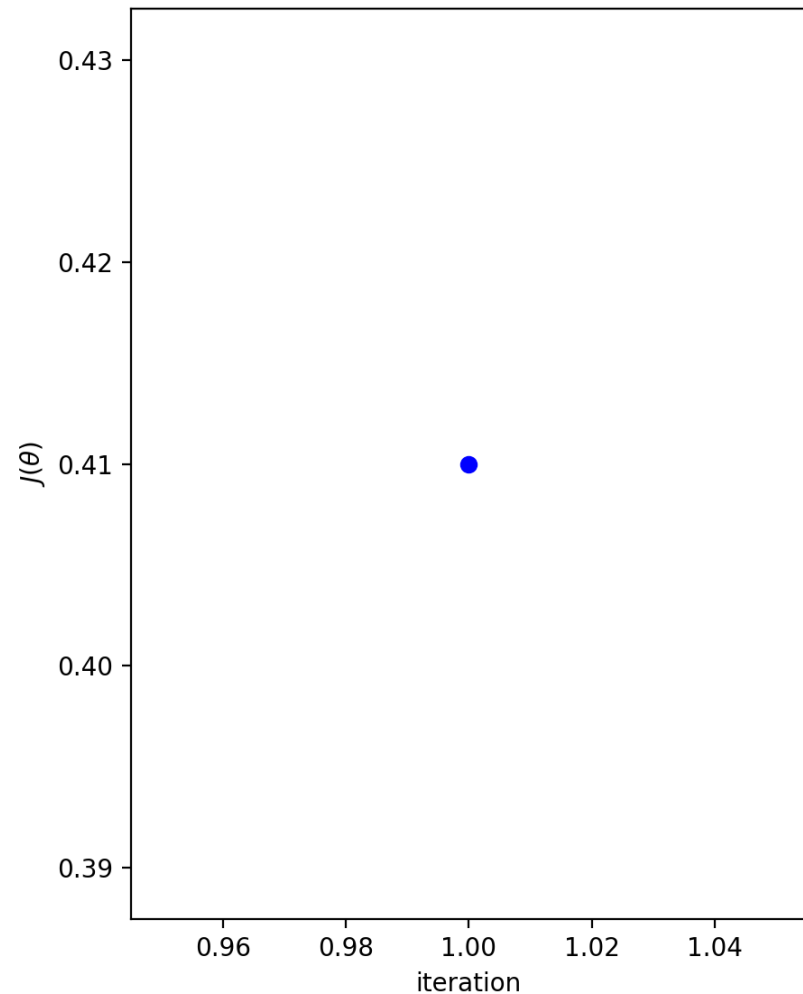
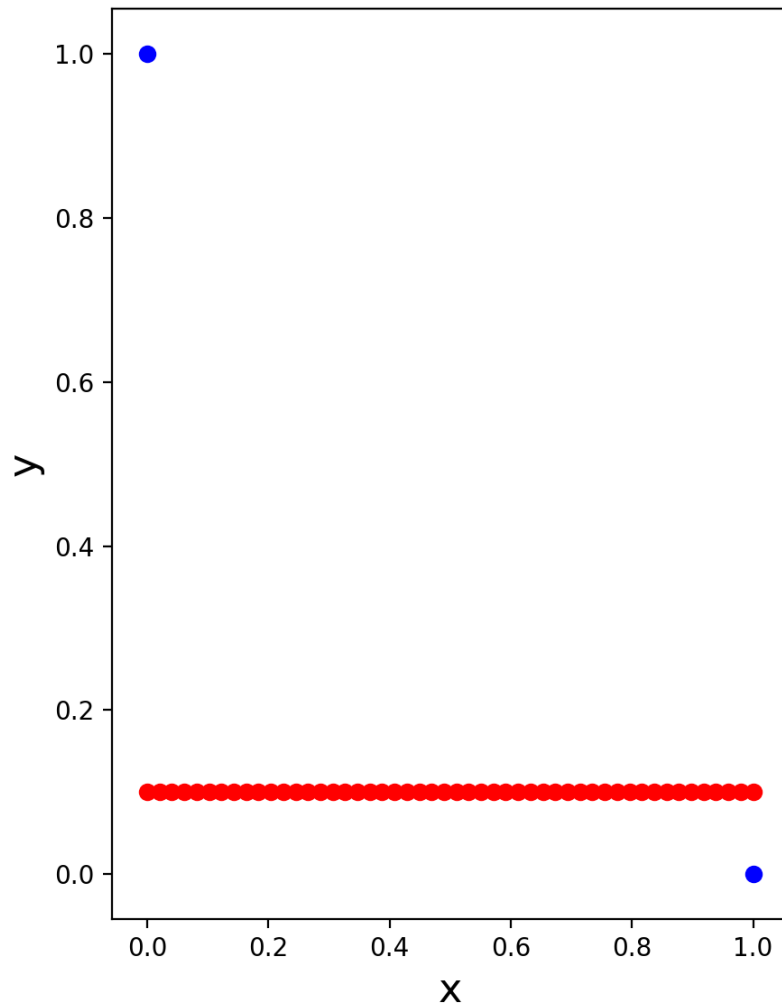
- Reading quiz #3
- Simple linear regression
- SGD (Stochastic Gradient Descent)
- Normal equations solution
- **Week 3 feedback**

Lab 3 applied to Handout 5

Toy example, iteration 1

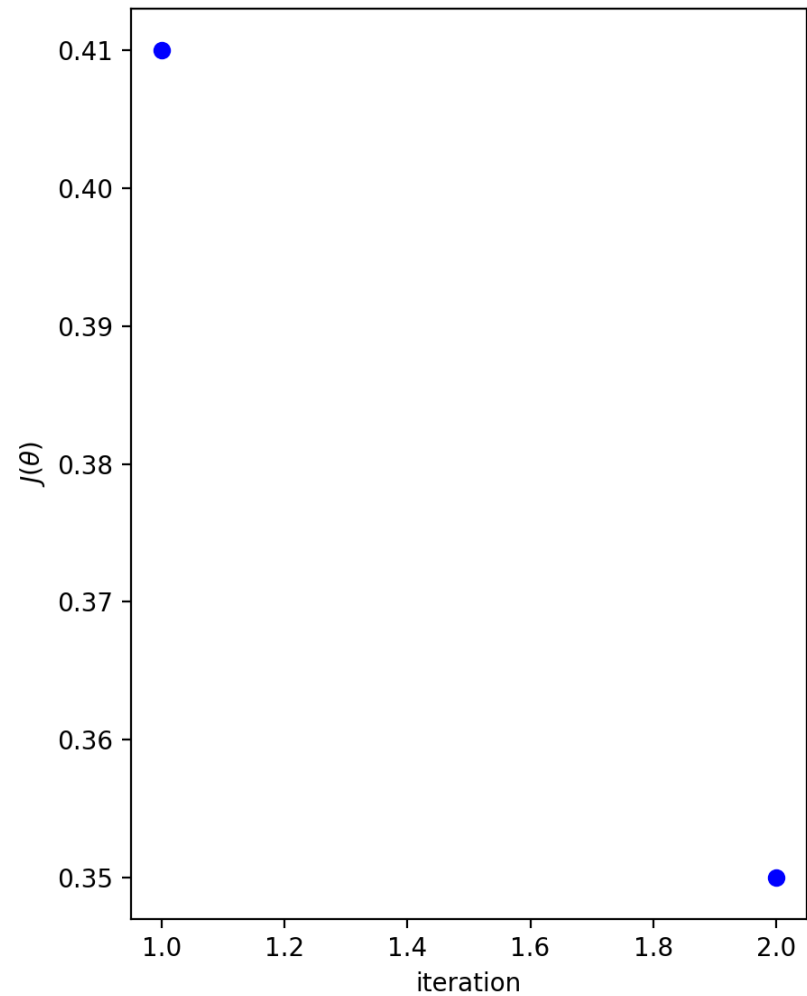
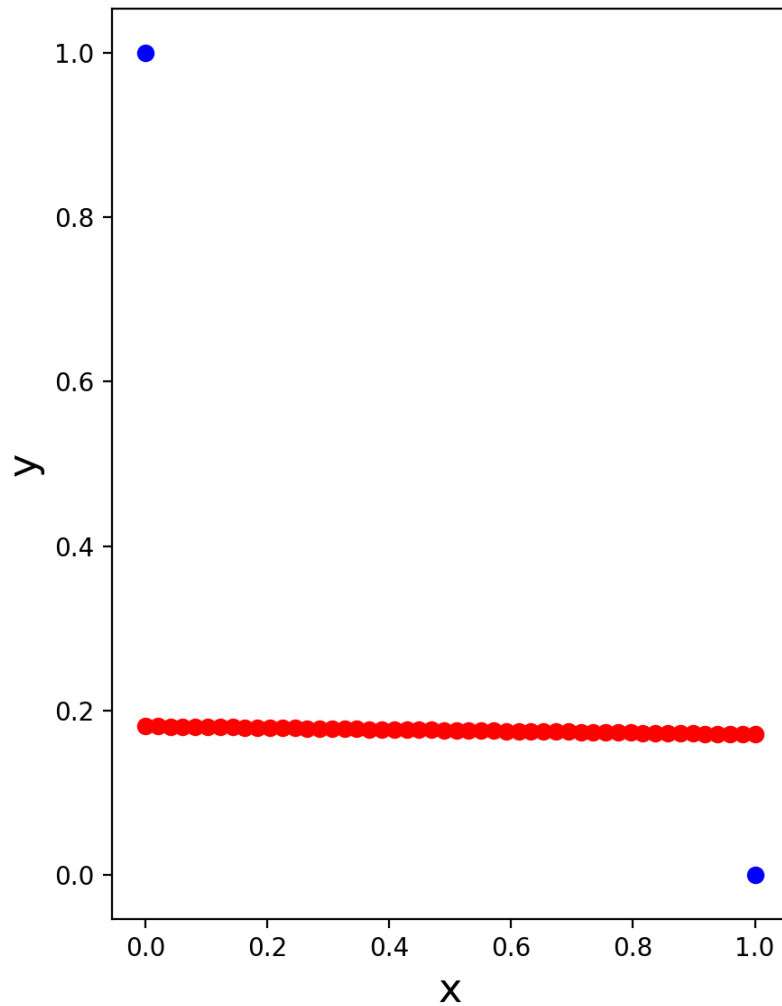
This is what you
should have obtained
in Handout 5!

iteration: 1, cost: 0.410000



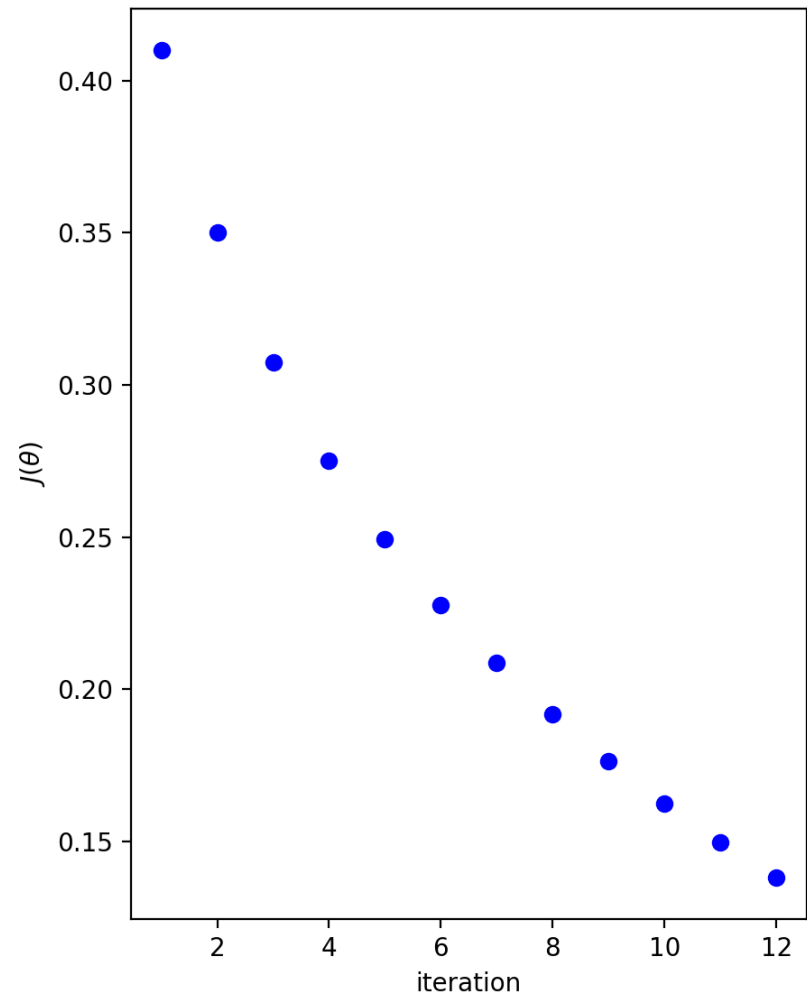
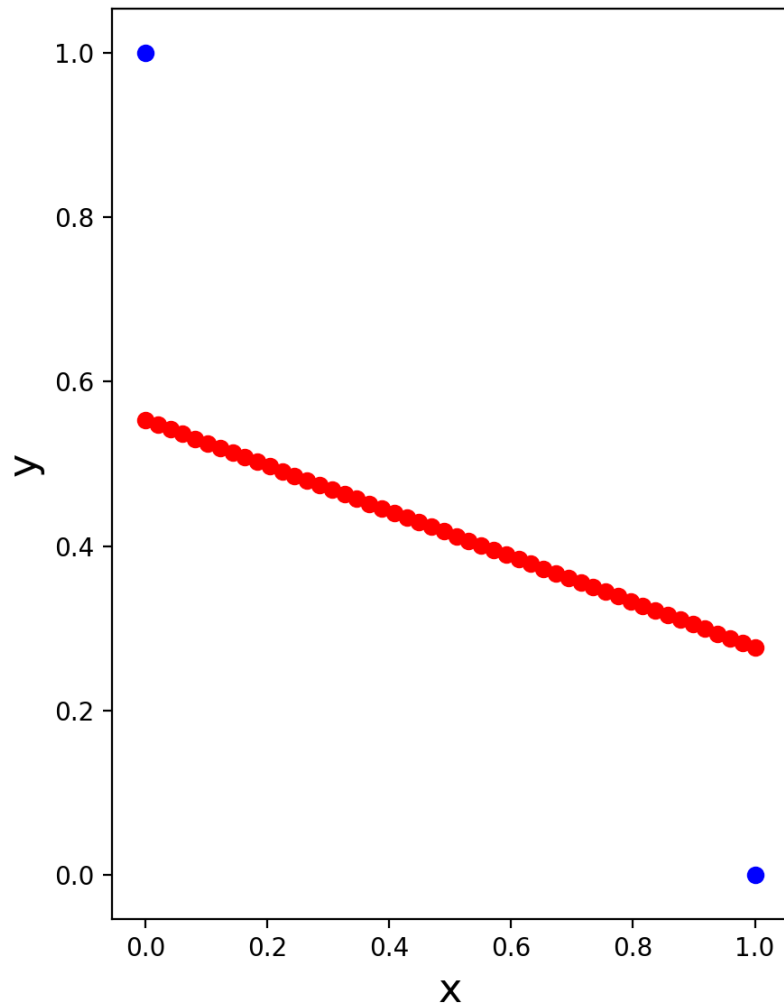
Toy example, iteration 2

iteration: 2, cost: 0.350001



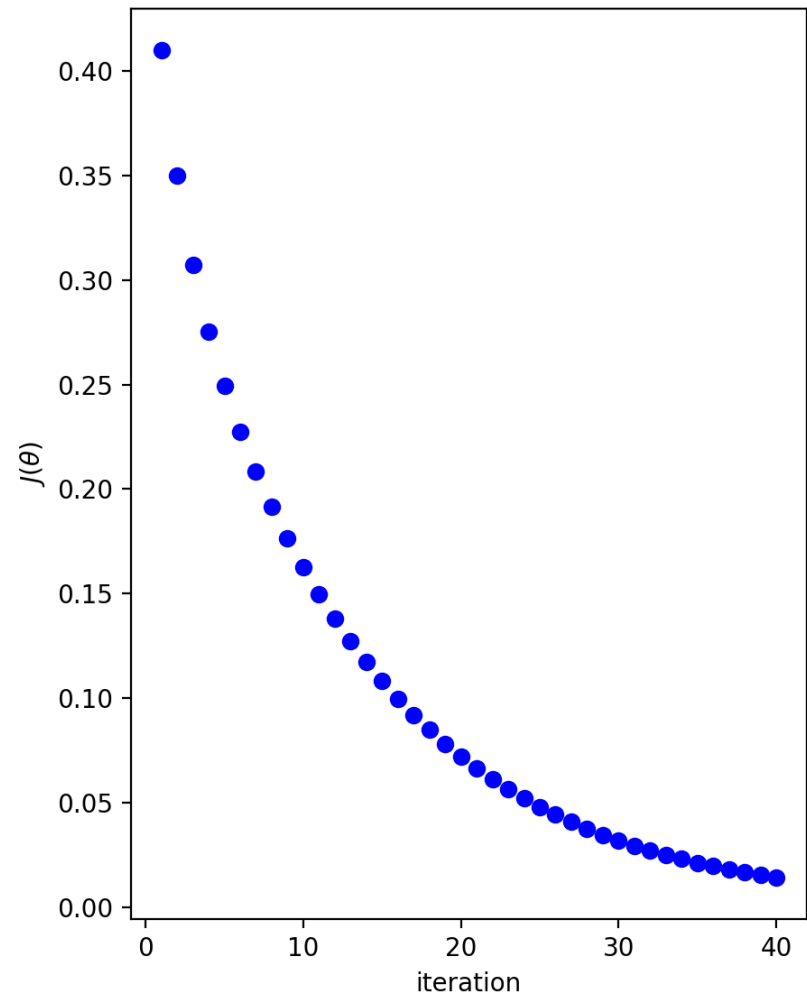
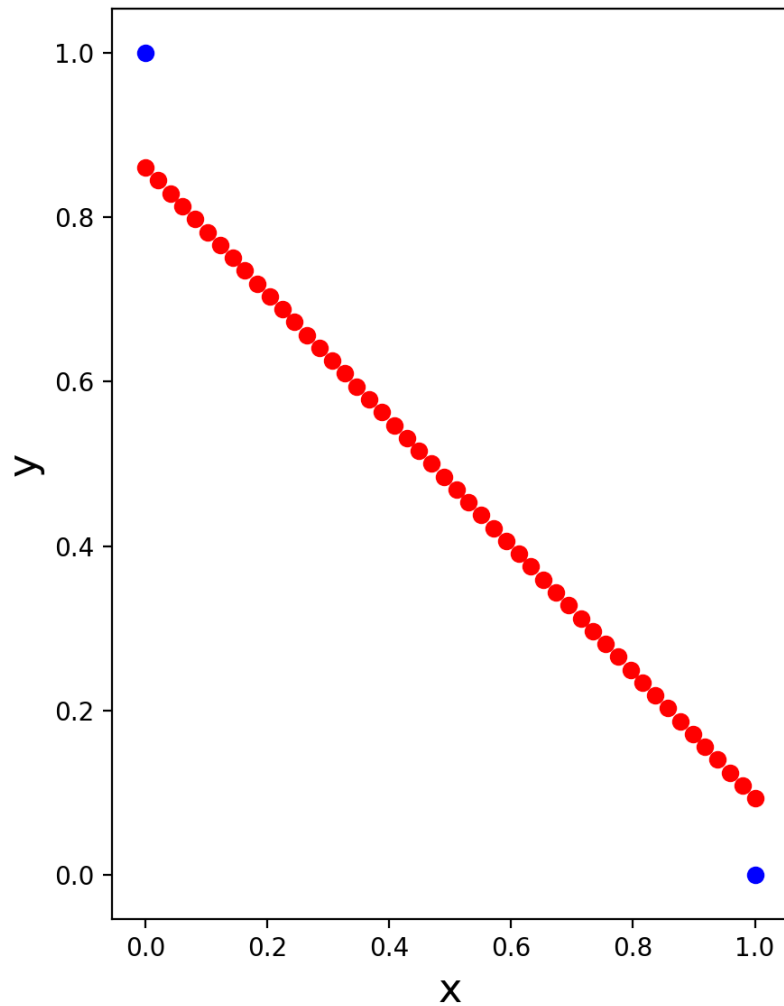
Toy example, iteration 12

iteration: 12, cost: 0.138047



Toy example, iteration 40

iteration: 40, cost: 0.014064



Toy example, iteration 100

iteration: 100, cost: 0.000105

