# CS 360: Machine Learning

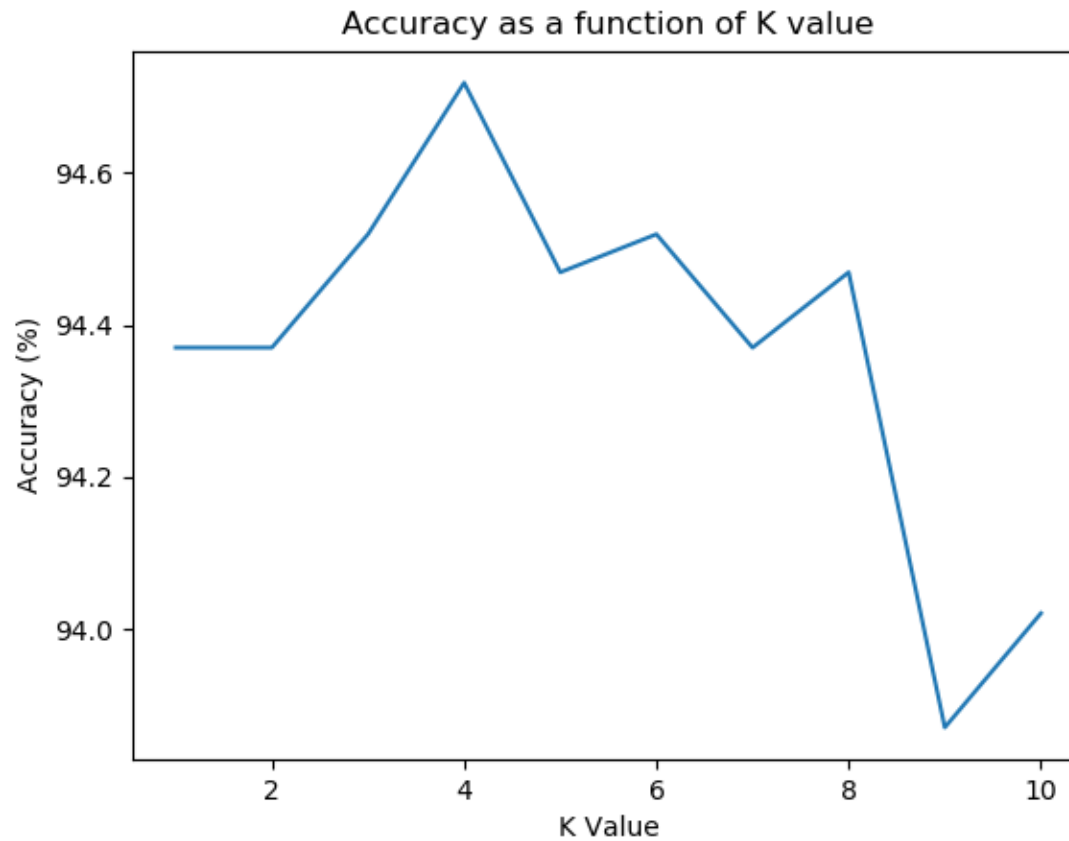Prof. Sara Mathieson

Fall 2020

HAVERFORD
COLLEGE

# Admin

- **Lab 1 graded**

- **Lab 2 collected 5pm today**
  - Last minute questions -> Piazza
  - I am unable to meet this afternoon

- **Lab 3** (shortened, step-by-step) due **Tuesday night** (grace period til Wed night)
  - Starter code: one place coeff -> coef
  - For part (a), easier to use train_data.plot(..)

# 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!)

# Lab 1 Extensions: multi-class

# 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 $K$

- 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 25

- Recap bias/variance tradeoff

- Simple linear regression

- SGD (Stochastic Gradient Descent)

- Normal equations solution

# Outline for September 25

- Recap bias/variance tradeoff

- Simple linear regression

- SGD (Stochastic Gradient Descent)

- Normal equations solution
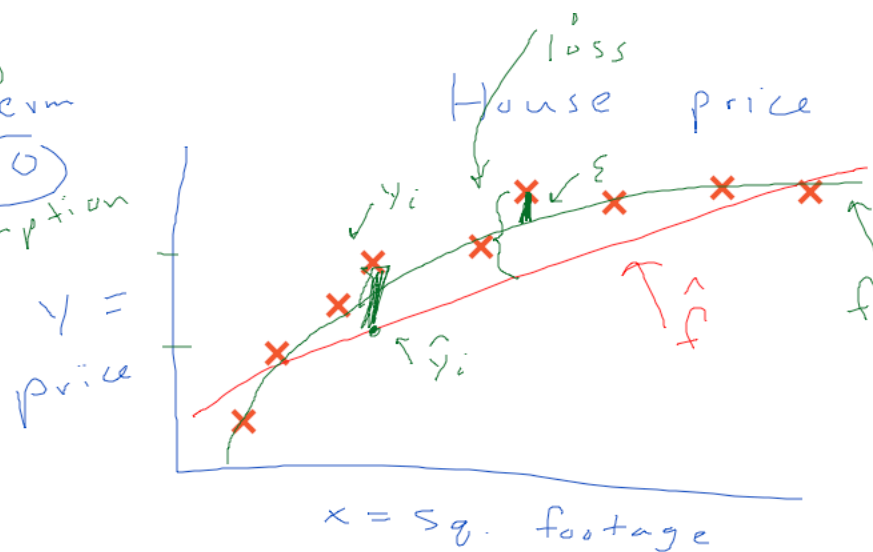
# Regression setup and Expected Value

Regression setup

model: true $y = f(x) + \varepsilon$

(noise) error term

mean 0 assumption

$\hat{f}$ = estimate of $f$

$\hat{y} = \hat{f}(x)$    prediction

loss

House price

$y = $ price

$y_i$

$\hat{y}_i$

$\hat{f}$

$f$

$x = sq.$ footage

GOAL    $\ell(y, \hat{y}) = (y - \hat{y})^2$

$E[(y - \hat{y})^2]$   expected loss

$E(D) = \frac{1}{10}(1 + 2 \cdots 5) + \frac{1}{2} \cdot 6$

$= 4.5$

expected value = weighted avg

$E[X] = \sum p(x = v) \cdot v$

$P(D = 6) = \frac{1}{2}$

$P(D \neq 6) = \frac{1}{10}$

$P(D)$

1 2 3 4 5 6

# Compute Expected Loss

$$E\left[(y-\hat{y})^2\right] = E\left[\left(\overbrace{y-f}^{a}+\overbrace{f-\hat{f}}^{b}\right)^2\right] \qquad \hat{f} = \hat{y}$$

$$\underbrace{y-f}_{\varepsilon} \qquad \underbrace{f-\hat{f}}_{0} \qquad \uparrow (x)$$

$$= \underbrace{\overbrace{Var(\varepsilon)}^{a^2}}_{\substack{noise \\ irreducible \\ error}} + \underbrace{E\left[\overbrace{(f-\hat{f})^2}^{b^2}\right]}_{\substack{reducible \\ error}}$$

$$(a+b)^2$$
$$= a^2 + b^2 + 2ab \quad 0$$
$$Var(x) = E\left[(x-\mu)^2\right]$$

$$E\left[(f-\hat{f})^2\right] = E\left[\left(\overbrace{f - E[\hat{f}]} + \underbrace{E[\hat{f}] - \hat{f}}_{0}\right)^2\right]$$



bias

too simple

variance

too complex

# Bias/variance big picture

$$\downarrow \underbrace{E\left[(y - \hat{y})^2\right]} = Var(\varepsilon) + \underbrace{E\left(f - \overbrace{E(\hat{f})}\right)^2}_{bias(\hat{f})^2} + \underbrace{E\left(E[\hat{f}] - \hat{f}\right)^2}_{Var(\hat{f})}$$

$$y = \overset{\downarrow}{f} + \textcircled{$\varepsilon$}$$
$$(true)$$

noise — irreducible

mean

polynomial

reducible



Y

bias: high

Variance: low

bias: low

Variance: high

overfitting

Y

X

X

# Assessing Model Accuracy



**FIGURE 2.9.** Left: *Data simulated from f, shown in black. Three estimates of f are shown: the linear regression line (orange curve), and two smoothing spline fits (blue and green curves).* Right: *Training MSE (grey curve), test MSE (red curve), and minimum possible test MSE over all methods (dashed line). Squares represent the training and test MSEs for the three fits shown in the left-hand panel.*

# Outline for September 25

- Recap bias/variance tradeoff

- 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?

# Regression Example

# Example: predict sales from TV advertising budget

# Cost Function: sum of squared errors



$p=1$

$p=2$

# Linear Regression

- 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

intercept

$$\hat{f}(x) = \hat{w}_0 + \hat{w}_1 x = \hat{y} = h_{\hat{w}}(x) \longleftarrow \text{hypothesis}$$

minimize: $\dfrac{1}{n} \displaystyle\sum_{i=1}^{n} (\hat{y}_i - y_i)^2 = MSE$ $\Bigg\}$ $\approx$ $\begin{array}{l} RSS \\ SSE \end{array}$

cost: $J(\vec{w}) = \dfrac{1}{2} \displaystyle\sum_{i=1}^{n} \left( w_0 + w_1 x_i - y_i \right)^2$  $[w_0, w_1]$  $\Bigg\}$  $\dfrac{1}{2}$ is for nice derivative

(a)

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

$$\Rightarrow) \quad \dfrac{n w_0}{n} = \dfrac{1}{n}\sum y_i - w_1 \dfrac{1}{n}\sum x_i$$

$\bar{y}$    $\bar{x}$

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

gradient    partial derivatives

$$\nabla J_{\vec{w}} = \begin{bmatrix} \dfrac{\partial J}{\partial w_0} \\[2mm] \dfrac{\partial J}{\partial w_1} \end{bmatrix} = \vec{0}$$

(a)

(b)

# "Simple" linear regression

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

$$\frac{\partial J}{\partial w_1} = \sum (w_0 + w_1 x_i - y_i) x_i = 0$$

Solve for $\boxed{w_1}$

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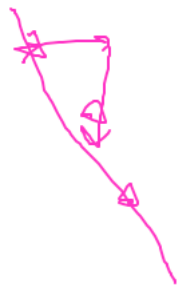
feature true output

mean

input orig data

pred
$$\hat{y}_i = \hat{w}_0 + \hat{w}_1 x_i$$

$Cov(x,y)$ moderate

high slope $\hat{w}_1$

low $Var(x)$

$\tilde{w}_1 = low$

$Cov(x,y)$ low

$Var(x)$ high

Handout 4
(on piazza, try on your own!)

# Outline for September 25

- Recap bias/variance tradeoff

- Simple linear regression

- SGD (Stochastic Gradient Descent)

- Normal equations solution

# Multiple linear regression

$$\hat{f}(\vec{x}) = w_0 + w_1 x_1 + w_2 x_2 + \cdots w_p x_p$$

$$= \vec{w} \cdot \vec{x}$$

$$[w_0 \quad w_1 \quad \cdots \quad w_p] \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix}$$

*fake 1*

$\vec{x}_i$ example

$x_j$ feature

$x_{ij}$

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

[SGD] Stocartic gradient descent

$$\frac{\partial J}{\partial w_j} = \sum (\vec{w} \cdot \vec{x}_i - y_i) \cdot x_{ij} \leftarrow$$

opp. direction of gradient

*for all* (i)

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

step size

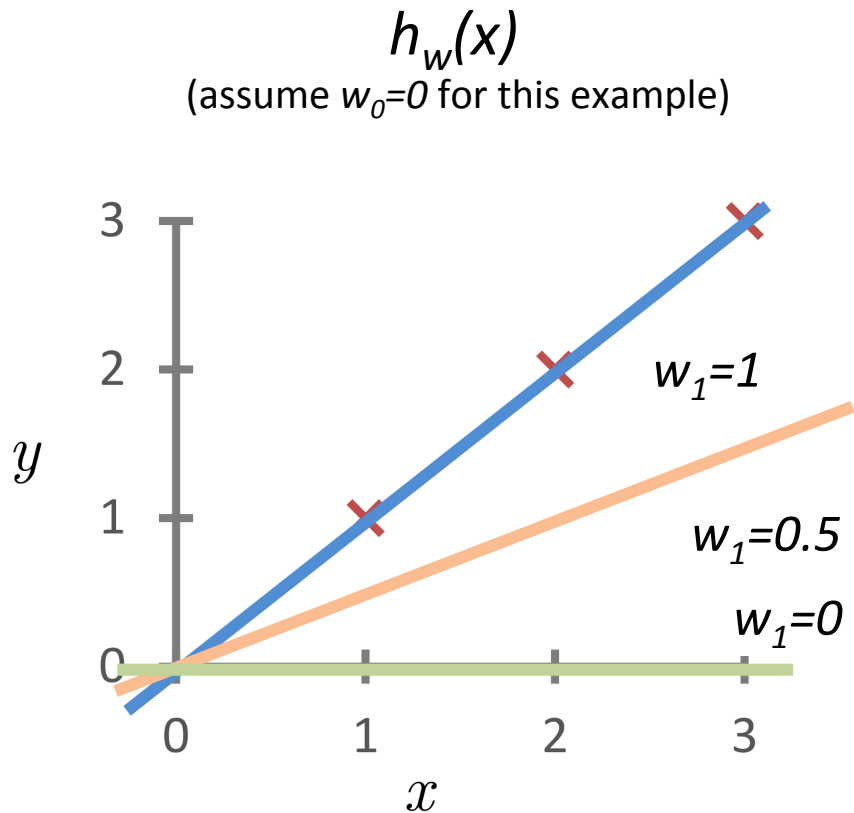# Cost Function

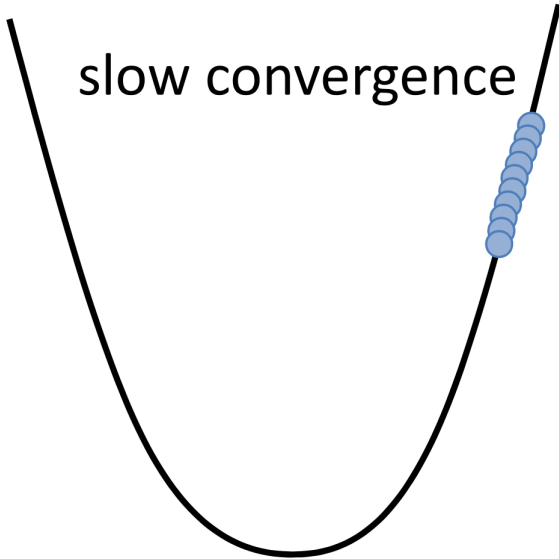# Gradient Descent: walking toward the minimum

# Cost Function (extra practice)



$h_w(x)$
(assume $w_0=0$ for this example)

$J(w_1)$

$w_1=1$

$w_1=0.5$

$w_1=0$

$J(0) = 7$

$J(w_1)$

$J(0.5) = \frac{1}{2}\left[(0.5-1)^2 + (1-2)^2 + (1.5-3)^2\right] = 1.75$
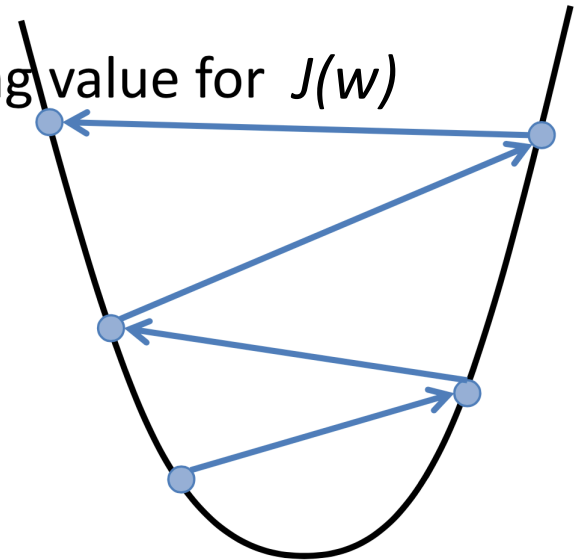
# Choosing step size

$\alpha$ **too small**

slow convergence

$\alpha$ **too large**

increasing value for  *J(w)*

- may overshoot minimum
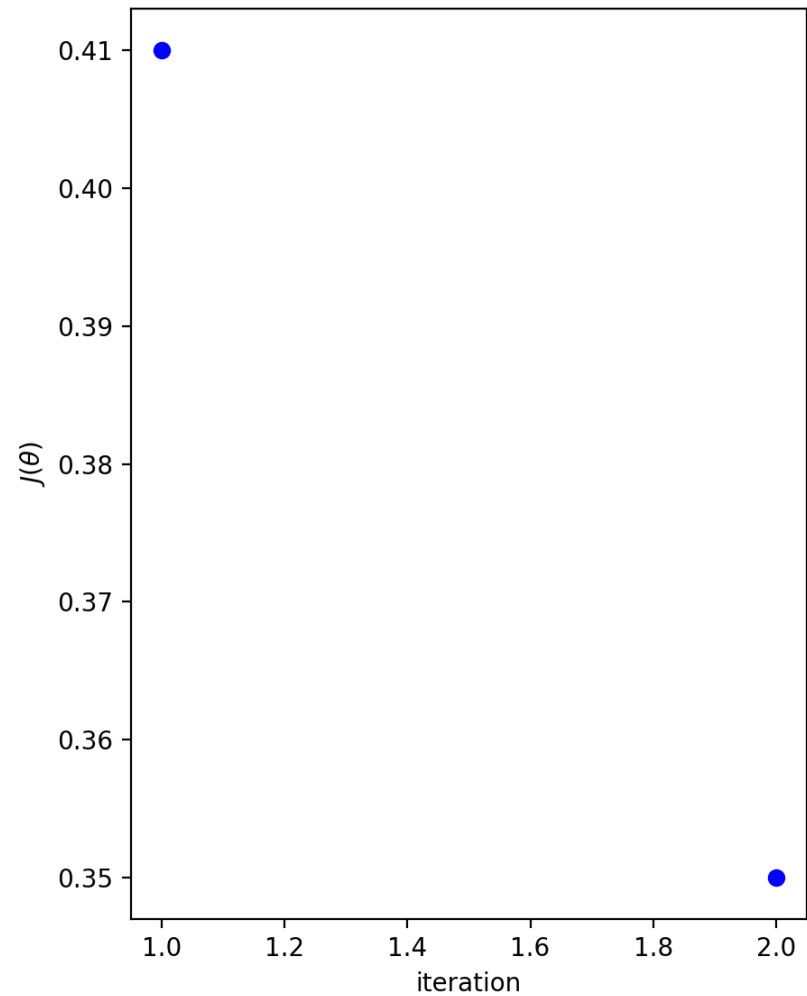- may fail to converge (may even diverge)
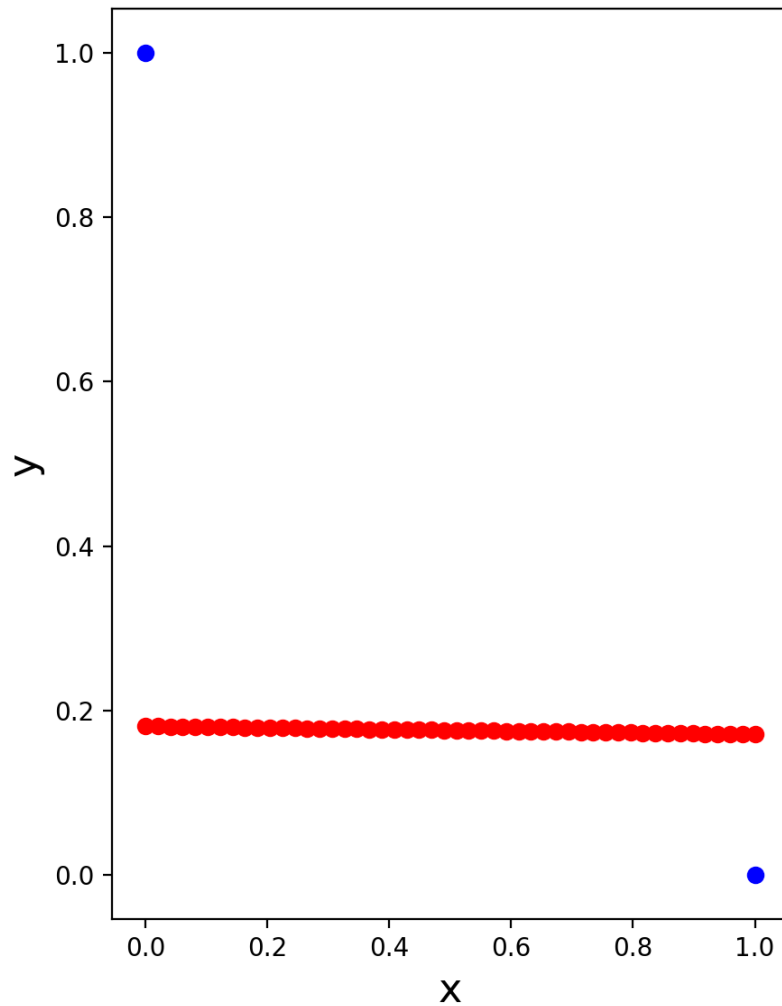
# Lab 3 applied to Handout 4

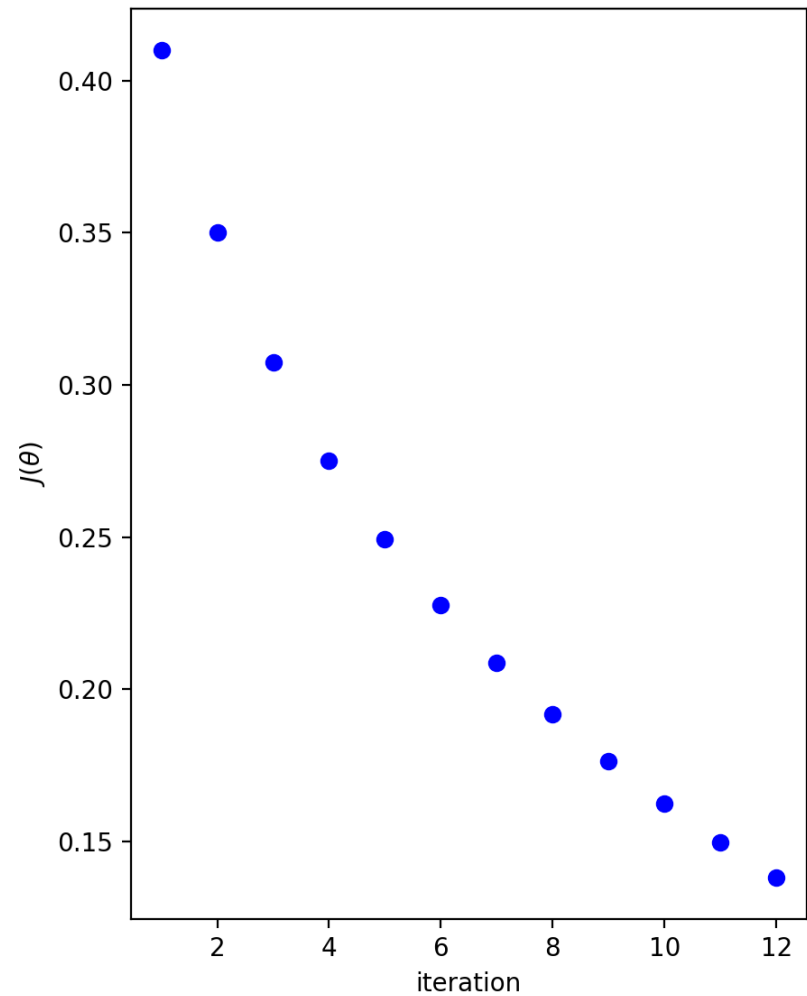# Toy example, iteration 1
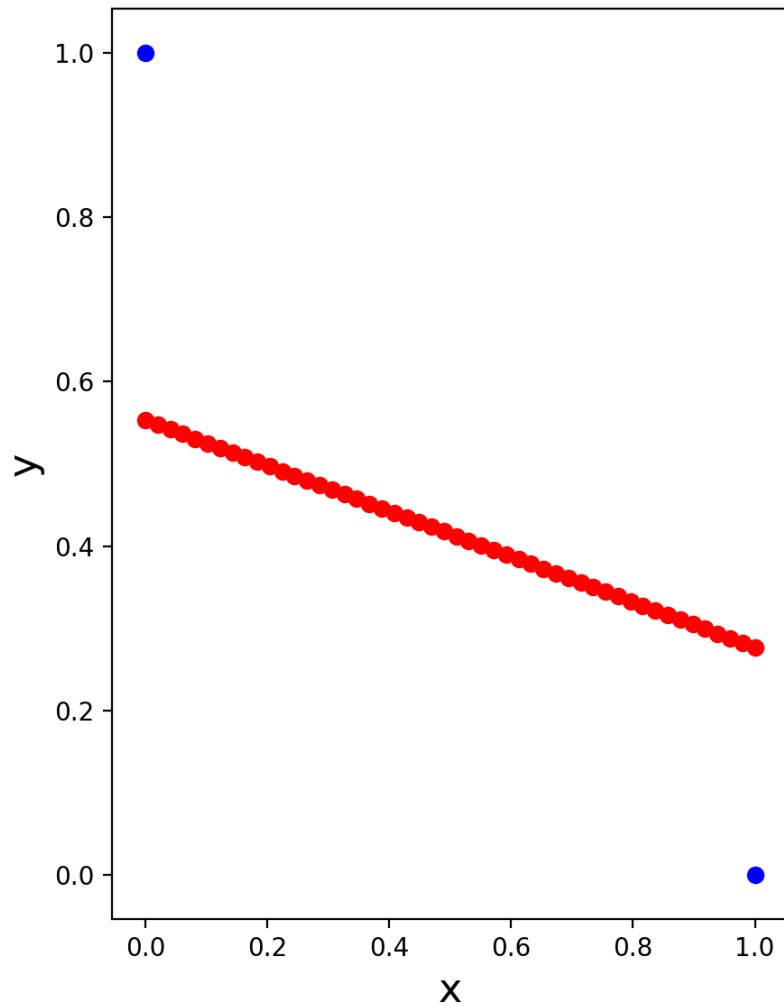
iteration: 1, cost: 0.410000

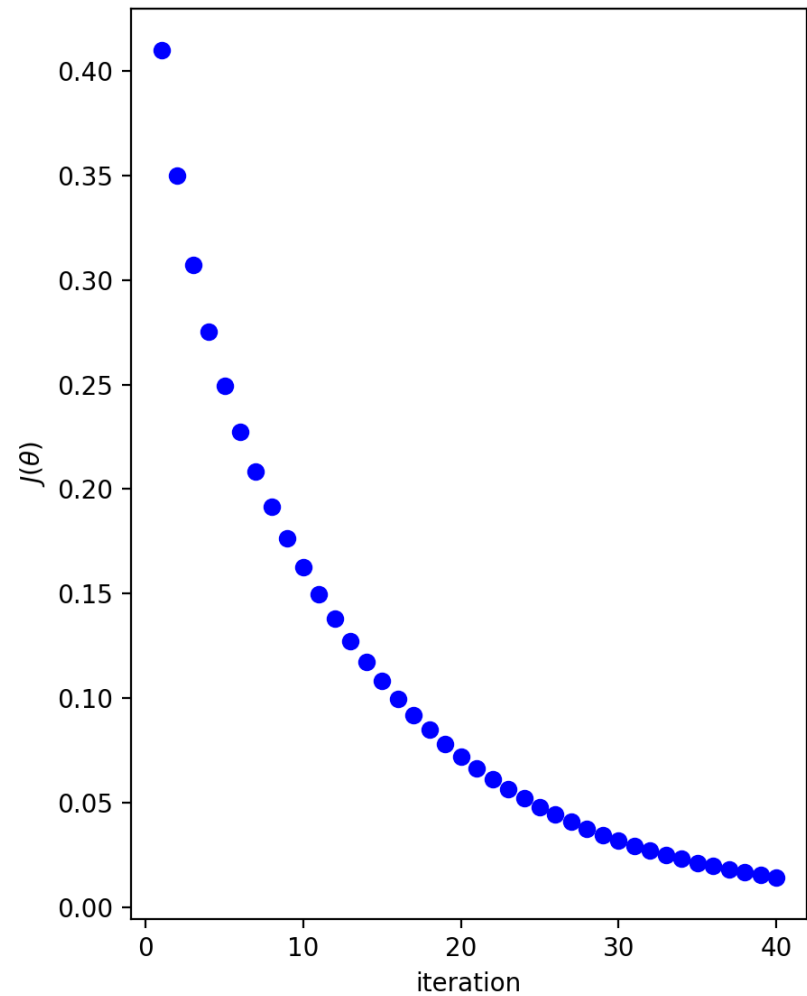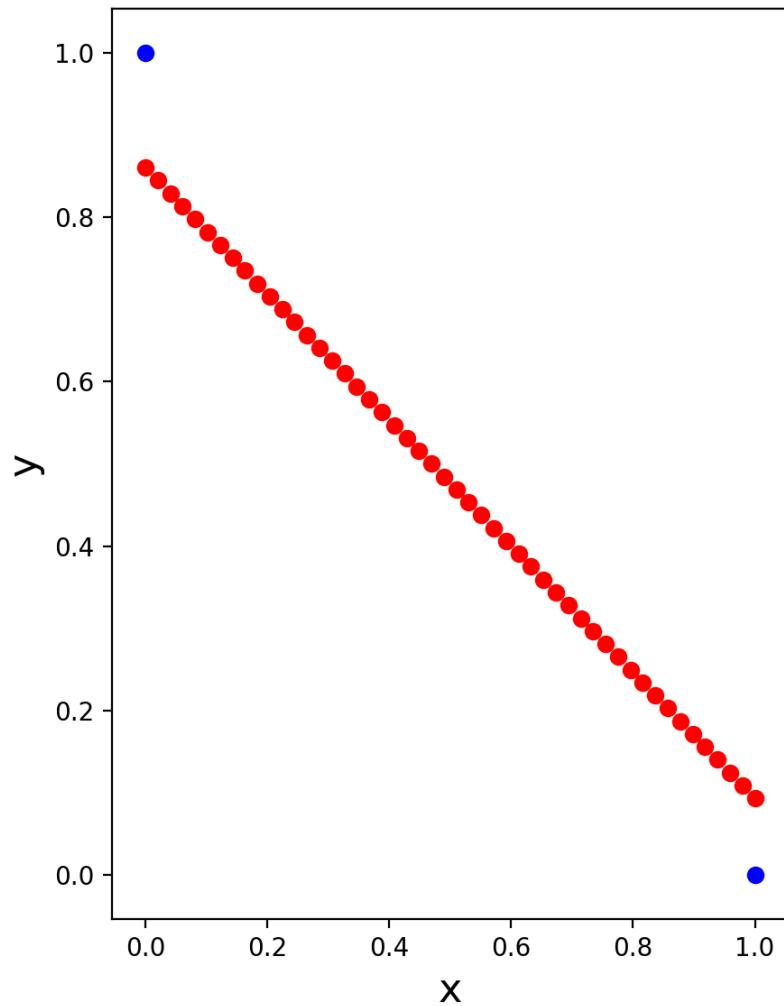# Toy example, iteration 2

iteration: 2, cost: 0.350001

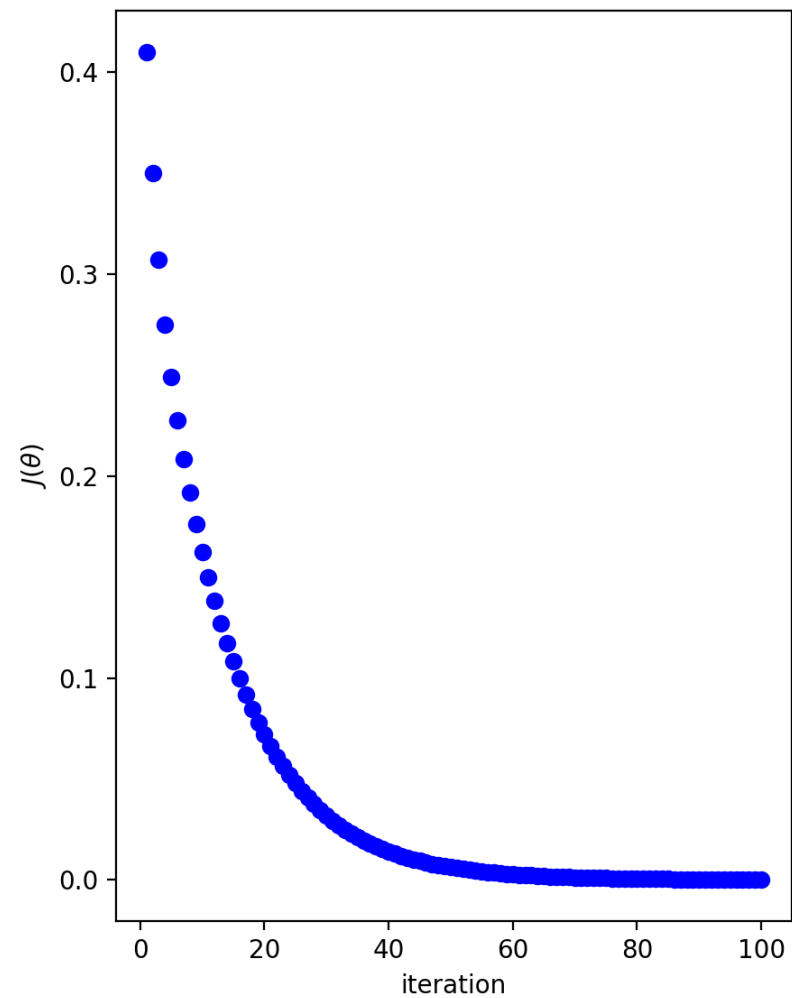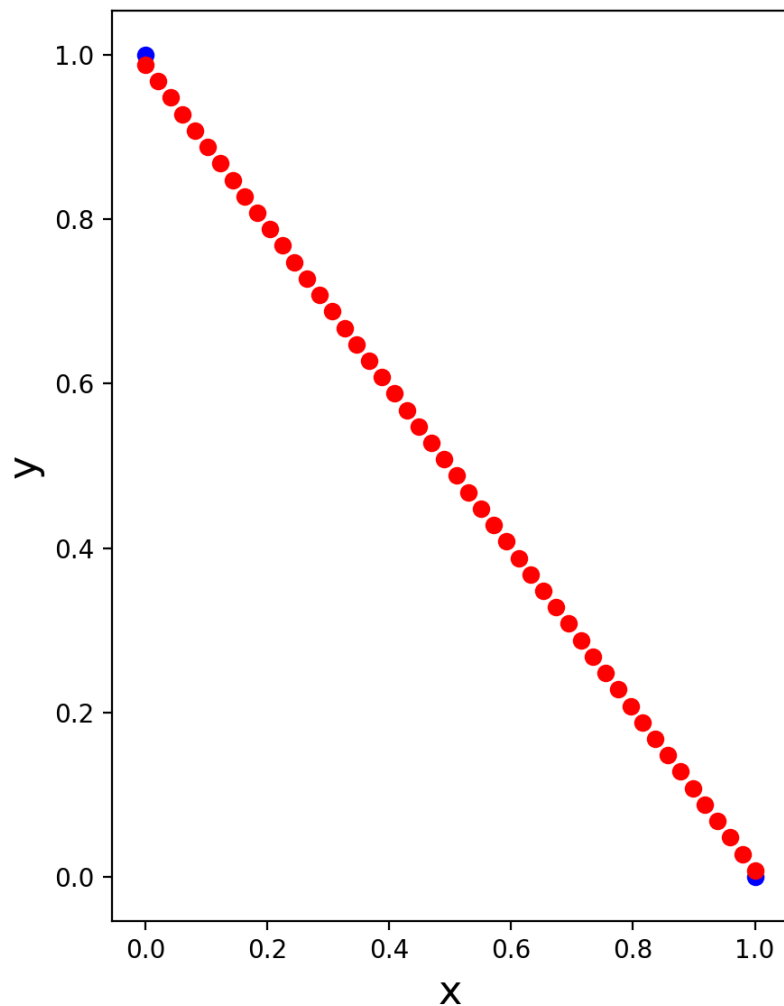# Toy example, iteration 12

iteration: 12, cost: 0.138047

# Toy example, iteration 40

iteration: 40, cost: 0.014064

# Toy example, iteration 100

iteration: 100, cost: 0.000105

# Outline for September 25

- Recap bias/variance tradeoff

- Simple linear regression

- SGD (Stochastic Gradient Descent)

- Normal equations solution

# Pros and Cons

## 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 $O(p^3)$