

CS 360: Machine Learning

Prof. Sara Mathieson

Fall 2019



Admin

- **Office hour** modifications (all in Hilles 110)
 - Mon 2-3:45pm
 - Tues 12:30-1:30pm
 - Friday 4-5pm
- Lab 3 due **TONIGHT**
- Lab 4 due **Sunday night** (short Naïve Bayes exercise)
- Reading quiz Thursday (just **Duame 9.3** => 2 pages!)
- Midterm next Thursday (short in-lab part + take home)

Technology and Justice Keynote: Thursday!

Keynote: Vulnerabilities: How Social Media and Data Infrastructure are Exploited for Fun, Profit, and Politics

Thursday, September 26, 2019

4:30 p.m.

Sharpless Auditorium

danah boyd is a Principal Researcher at Microsoft Research, the founder and president of Data & Society, and a Visiting Professor at New York University. Her research is focused on addressing social and cultural inequities by understanding the relationship between technology and society.



Feedback forms

Understand well

- KNN 20
- Entropy 7
- Decision Trees 16
- Over/under fitting 1

Needs most work

- Linear regression 13
- SGD 4
- Decision trees 1
- Intuition behind entropy 4
- Bias/variance/loss 2
- Continuous features 2
- Tying topics together 2

Additional questions

Difficulty/workload
Count

1	2	3	4	5
		10	15	3

	Less	More	As Is
Slides	2	3	21
Board		7	19
Group	2	5	19
Handouts	1	10	15

- Random partners yes 19
- Random partners no 3
- Neutral 3

Computer policy

Students mentioned:

- noise of typing is distracting to other students
- temptation to look at non-class material is distracting for student and those around them
- laptops take up space on the small tables
- laptops can prevent students from seeing the board

Guidelines:

- Laptops strongly discouraged
- If a laptop will help your interaction with the class material, talk to me ahead of time!

Other thoughts

- Most take notes, posted notes are helpful
- Classroom setup is a bit awkward
- Labs have been long, after getting behind it can be hard to catch up
- More examples, more pseudocode during class
- Several mentioned implementation is a struggle

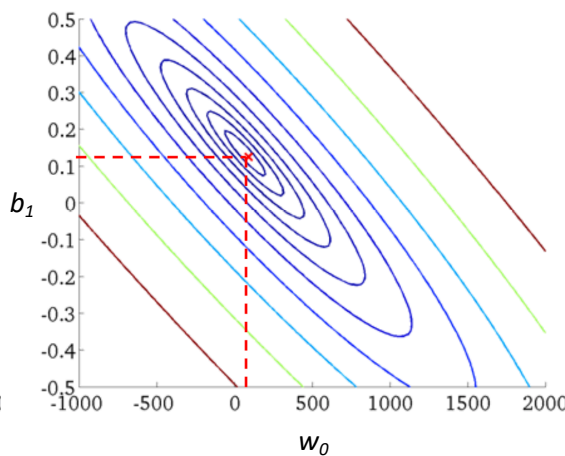
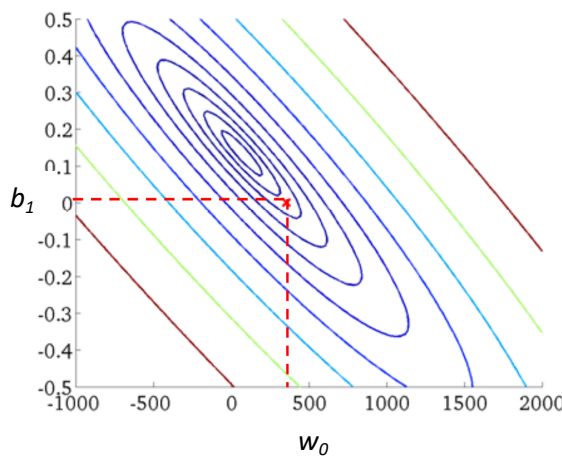
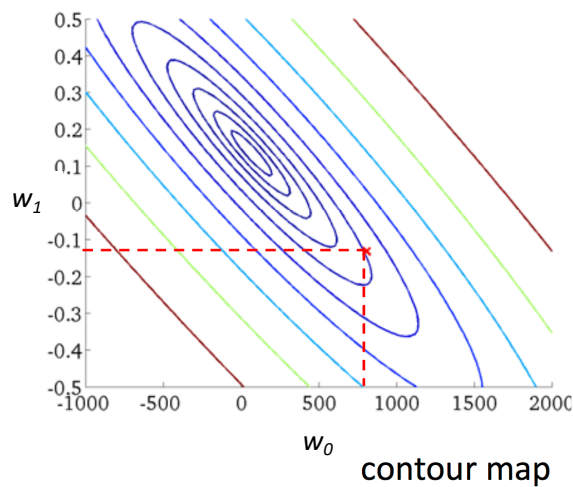
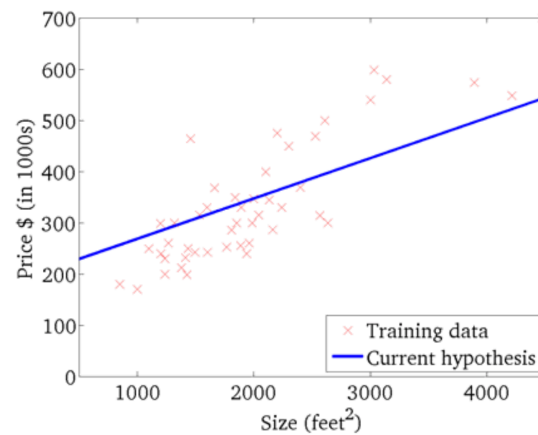
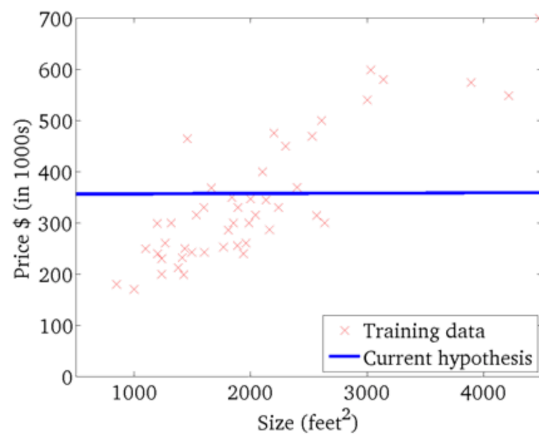
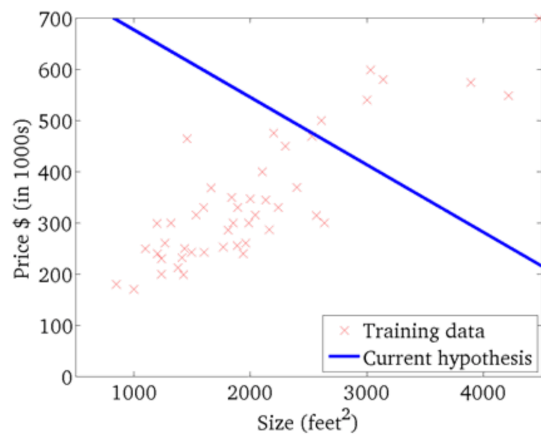
Outline for September 24

- Recap two ways of solving linear regression
 - SGD
 - Normal equations (closed-form solution)
- Regularization
- Introduction to conditional probability
 - Clinical trials example
- Begin: Naïve Bayes

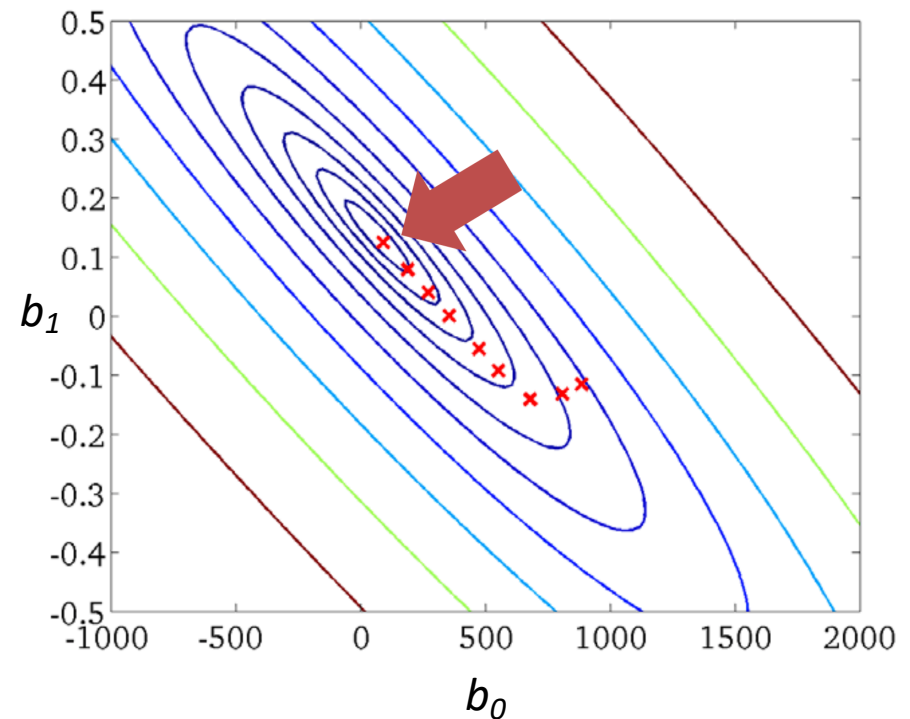
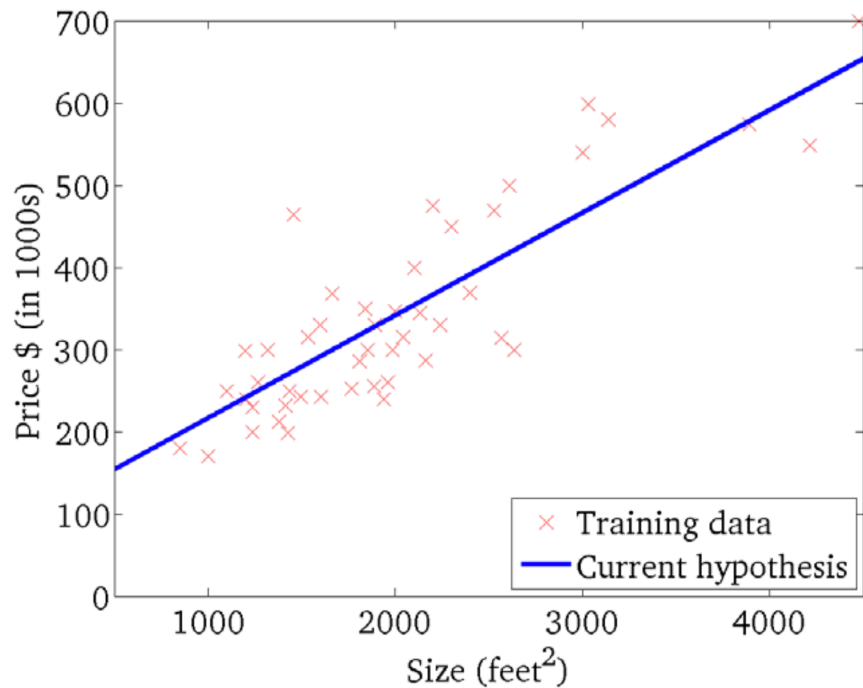
Outline for September 24

- Recap two ways of solving linear regression
 - SGD
 - Normal equations (closed-form solution)
- Regularization
- Introduction to conditional probability
 - Clinical trials example
- Begin: Naïve Bayes

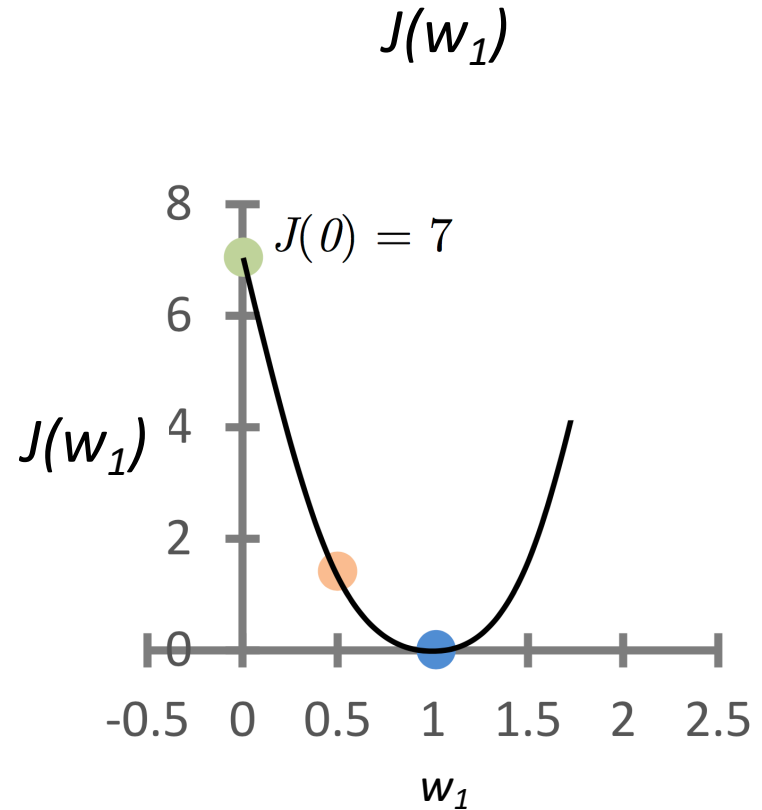
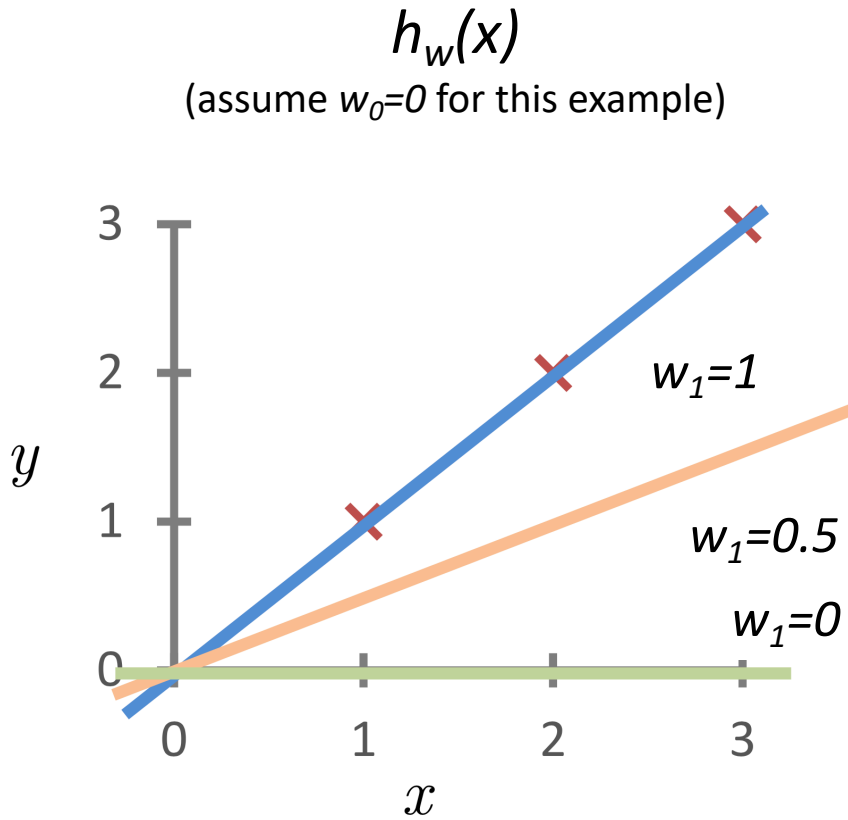
Cost Function



Gradient Descent: walking toward the minimum

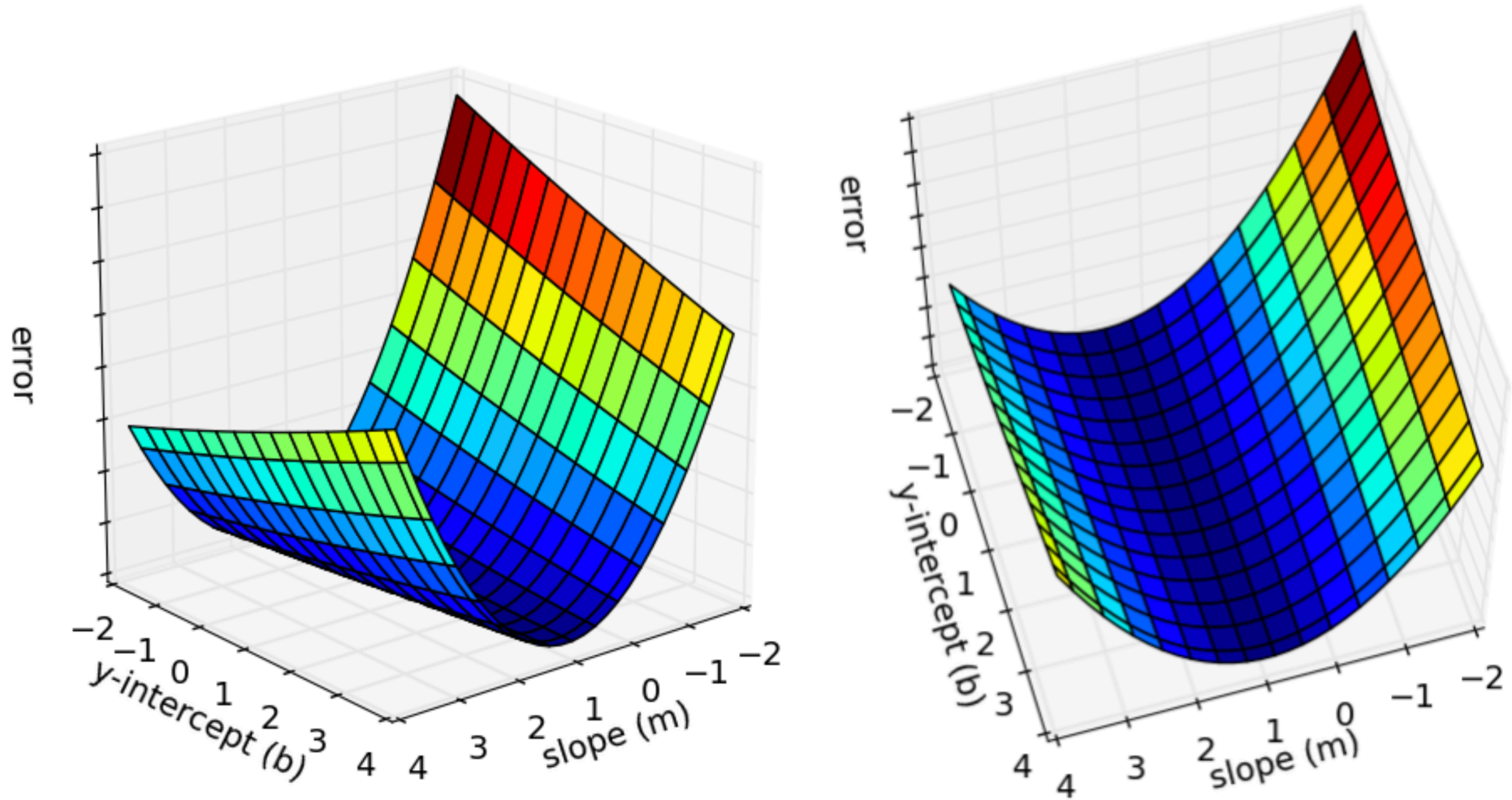


Cost Function (extra practice)



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

Error as a function of slope & y-intercept

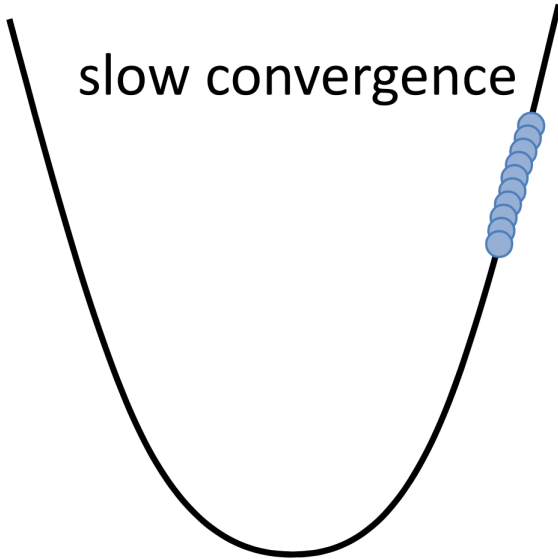


Convex \Rightarrow single global optimum (can prove it is a minimum with second derivative)!

Choosing step size α

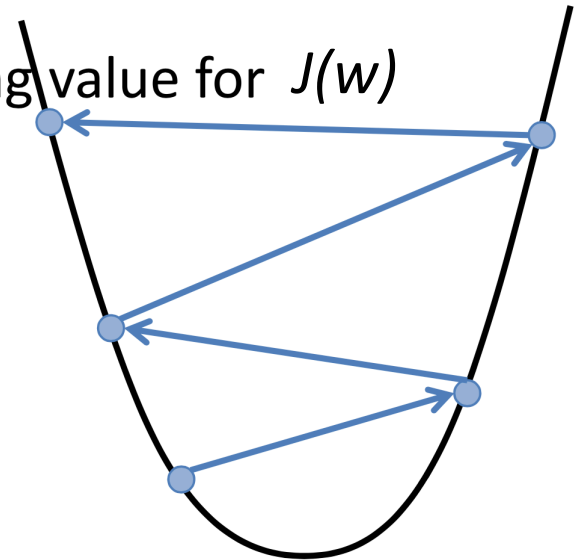
α too small

slow convergence



α too large

increasing value for $J(w)$



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

Pros and Cons

Gradient Descent

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

Normal Equations

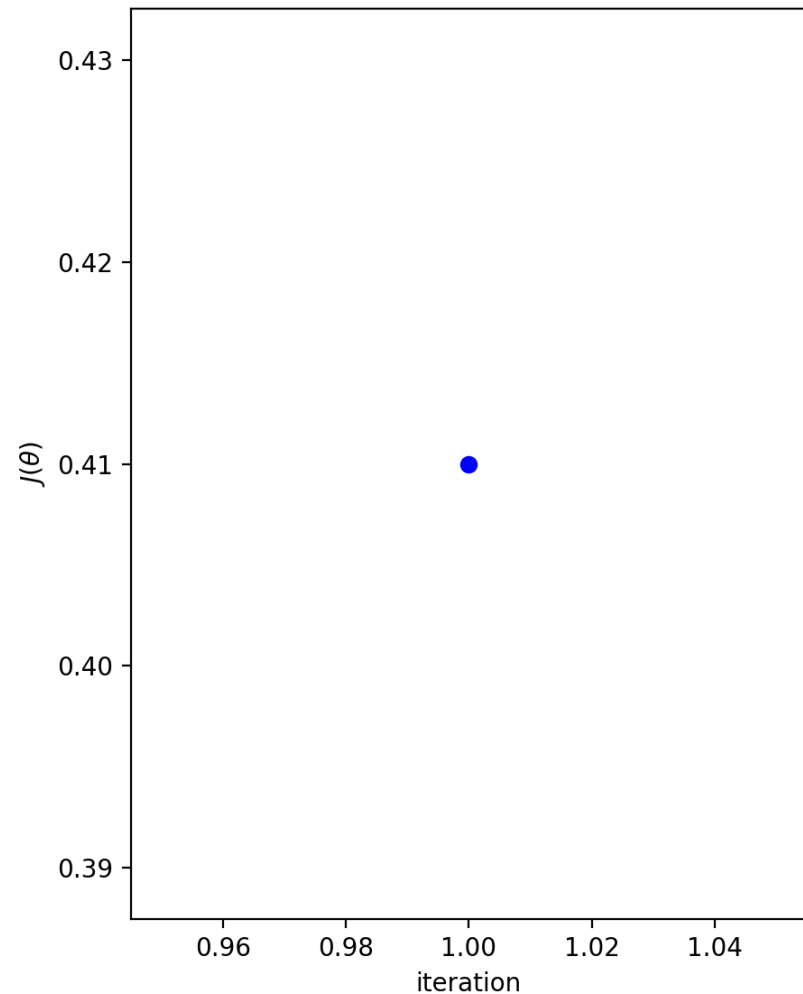
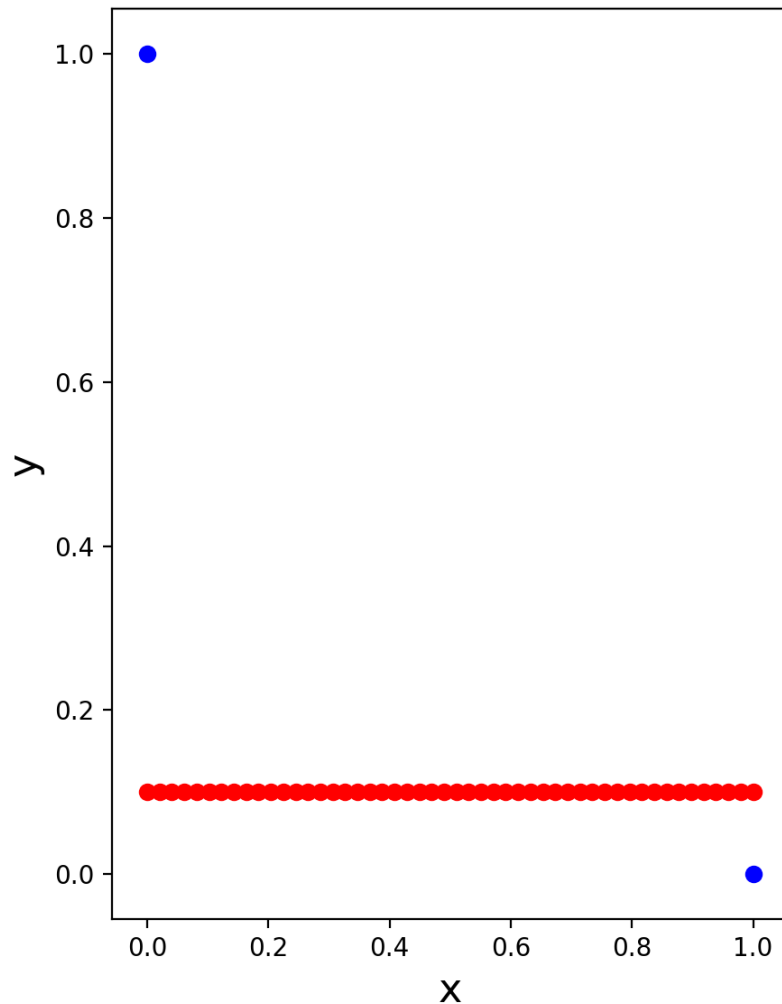
- non-iterative
- no need for α
- slow if p is large
 - matrix inversion is $O(p^3)$

Lab 3 applied to Handout 5

Toy example, iteration 1

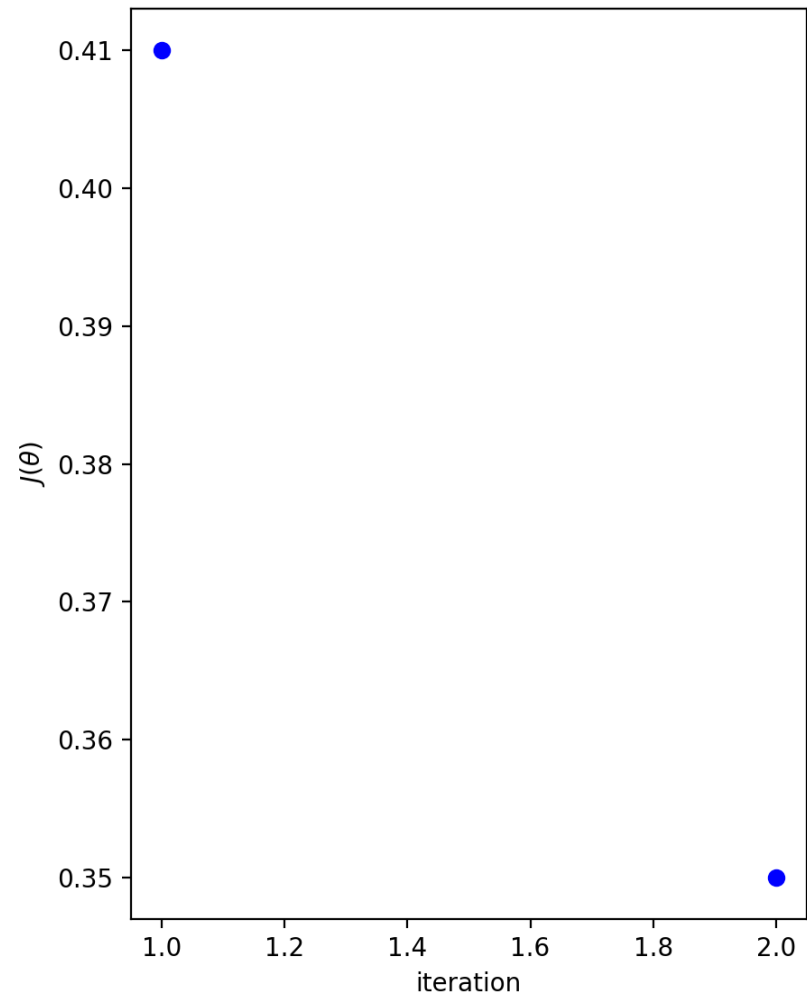
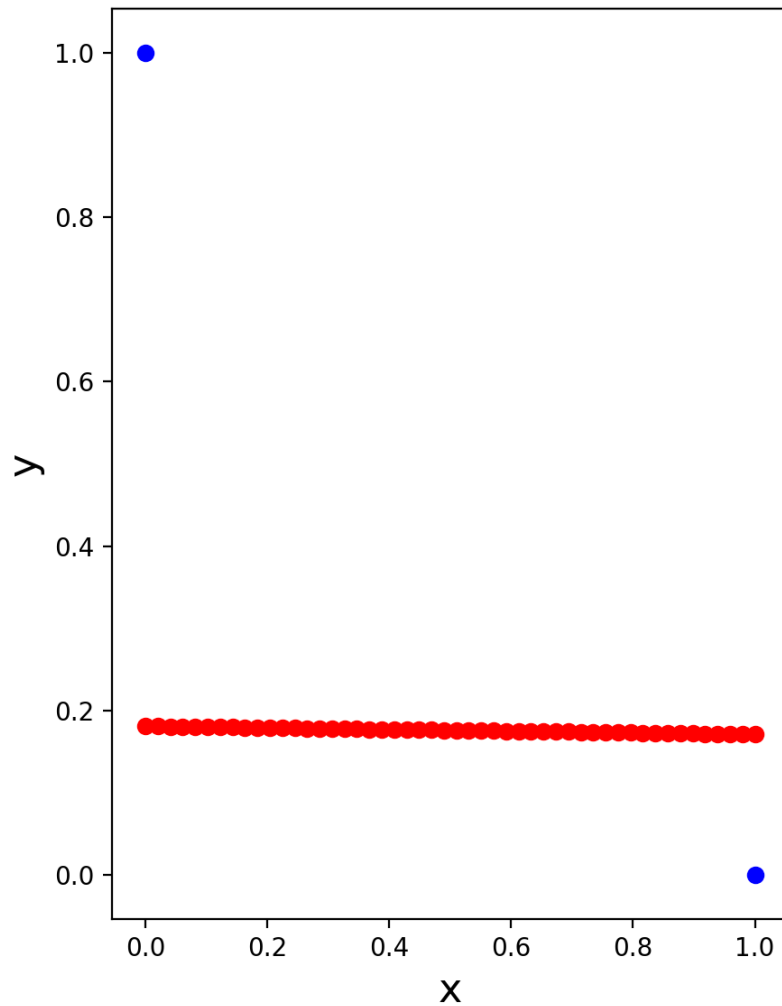
This is what you
should have obtained
in Handout 5!

iteration: 1, cost: 0.410000



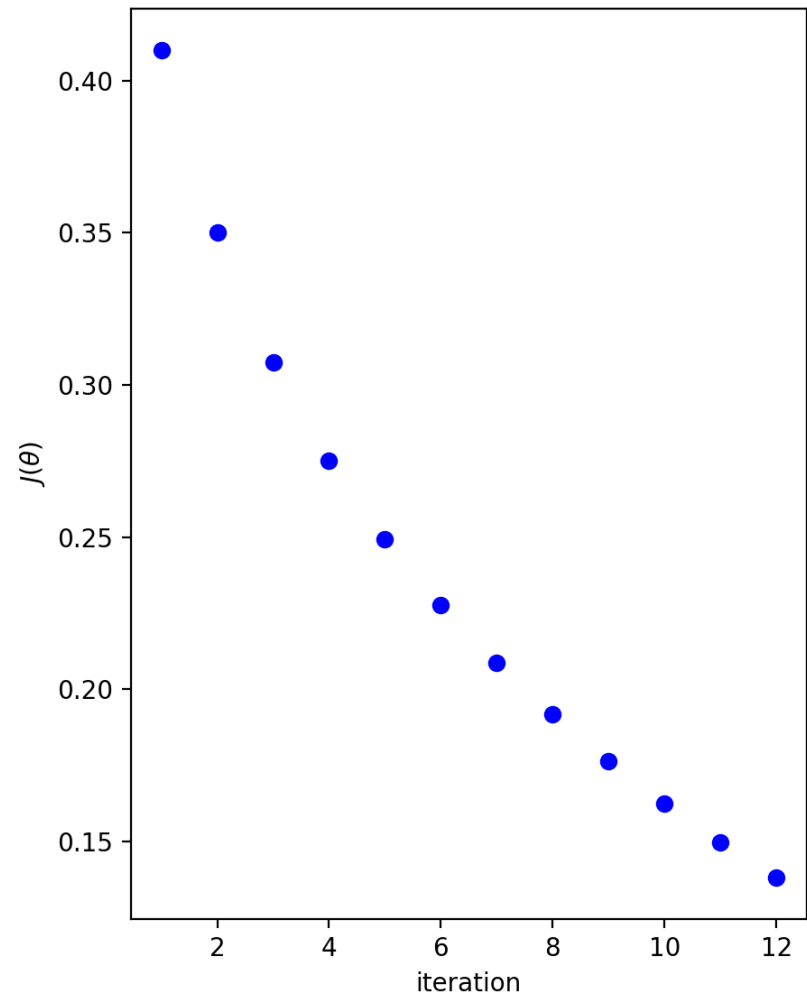
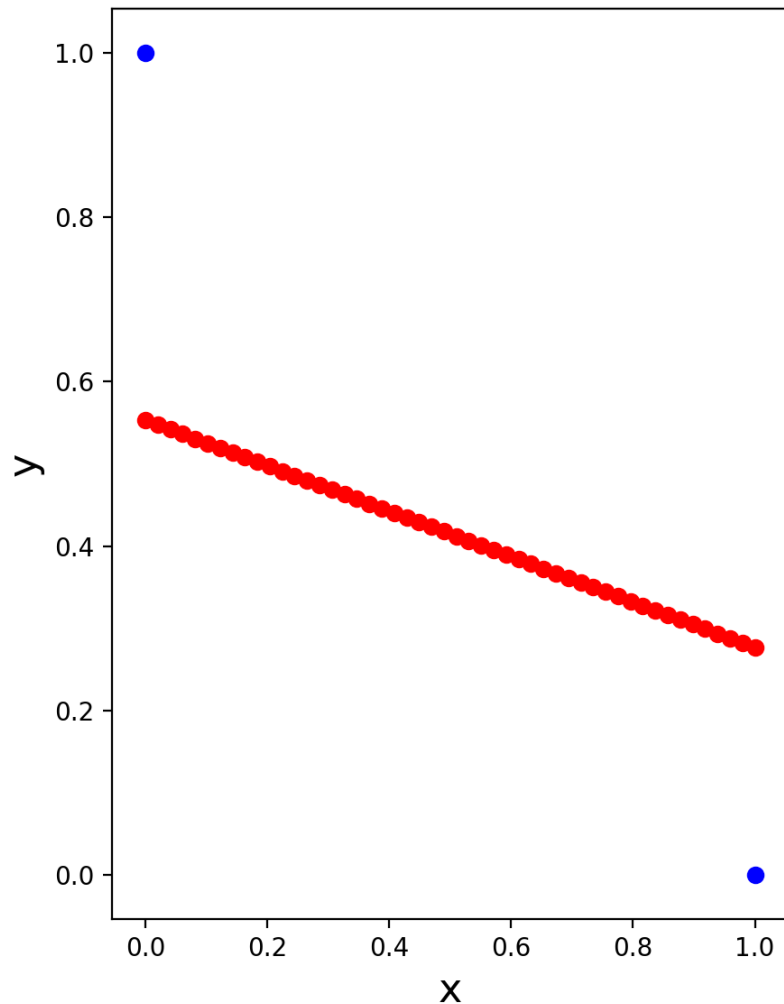
Toy example, iteration 2

iteration: 2, cost: 0.350001



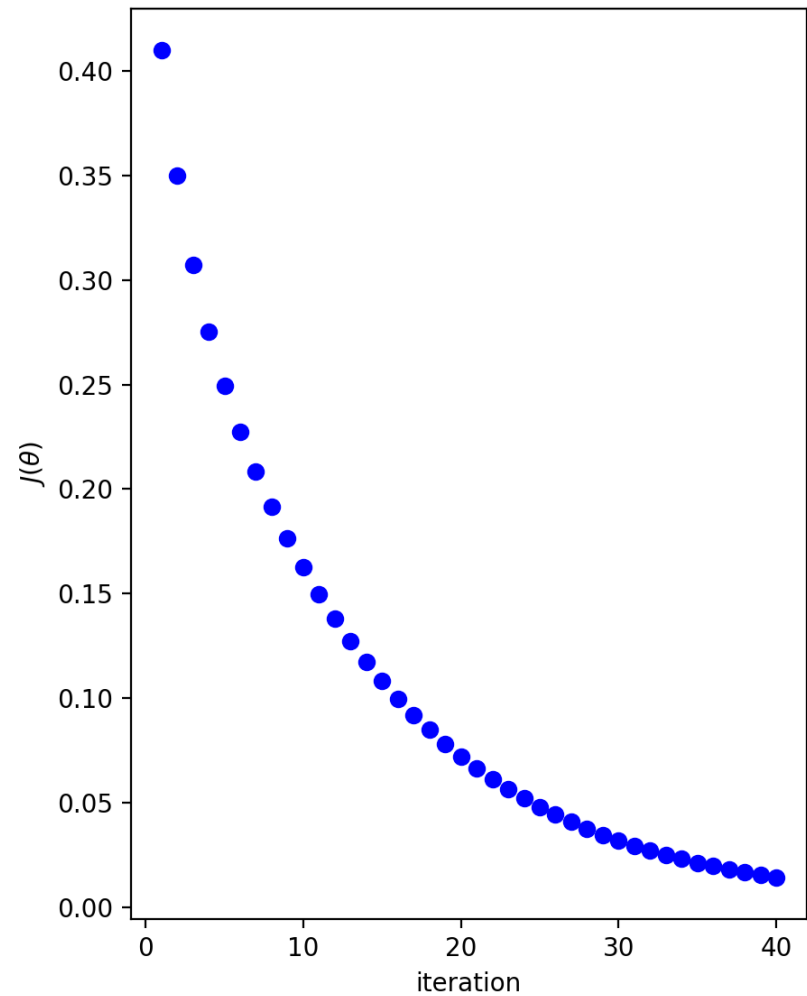
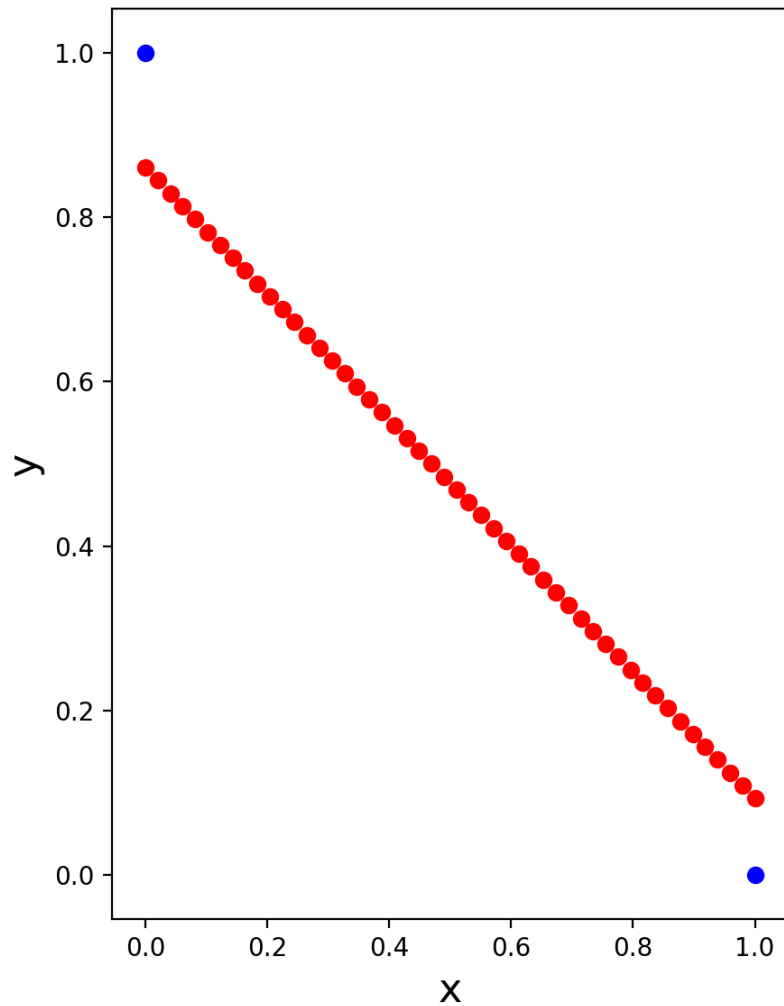
Toy example, iteration 12

iteration: 12, cost: 0.138047



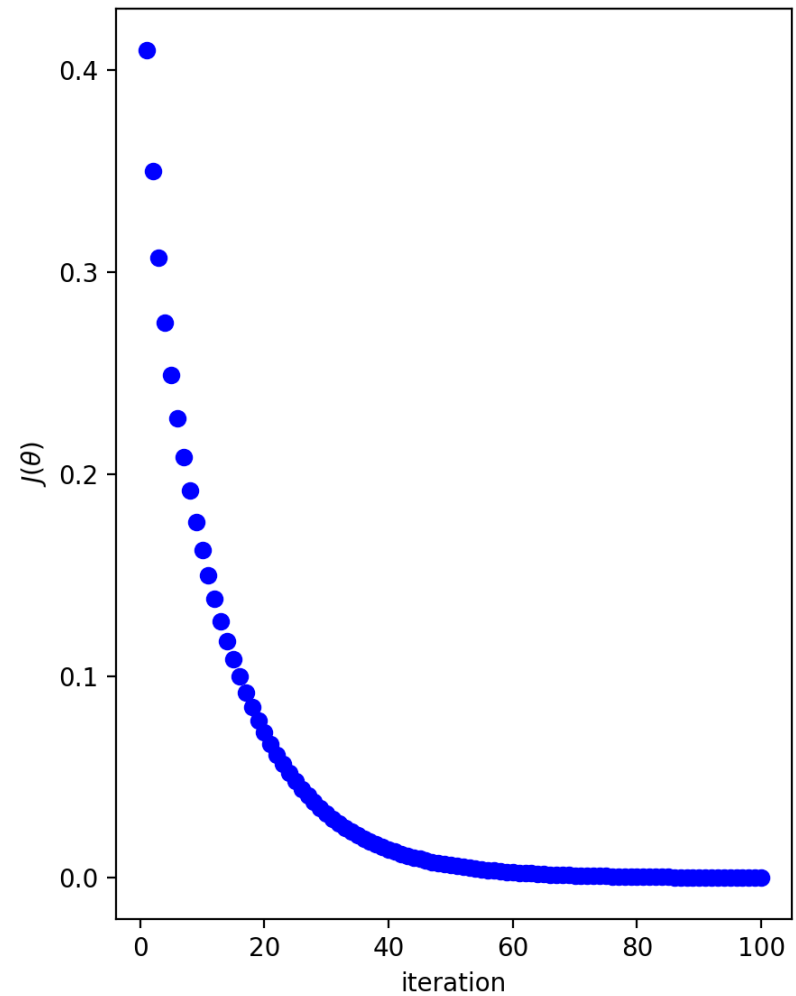
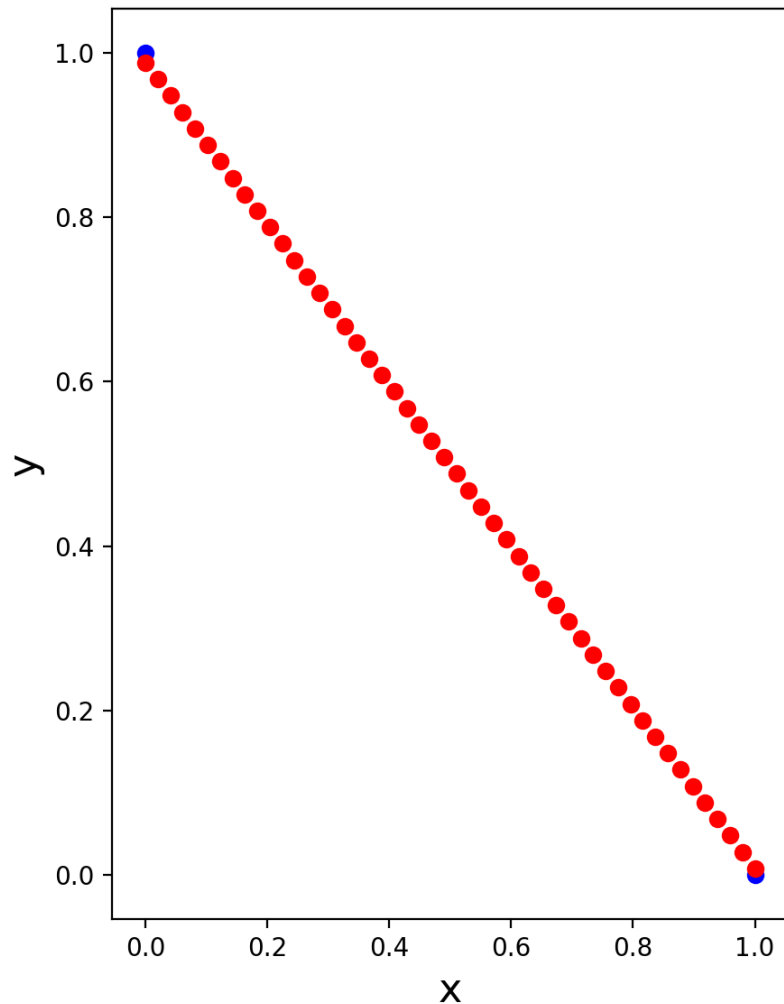
Toy example, iteration 40

iteration: 40, cost: 0.014064



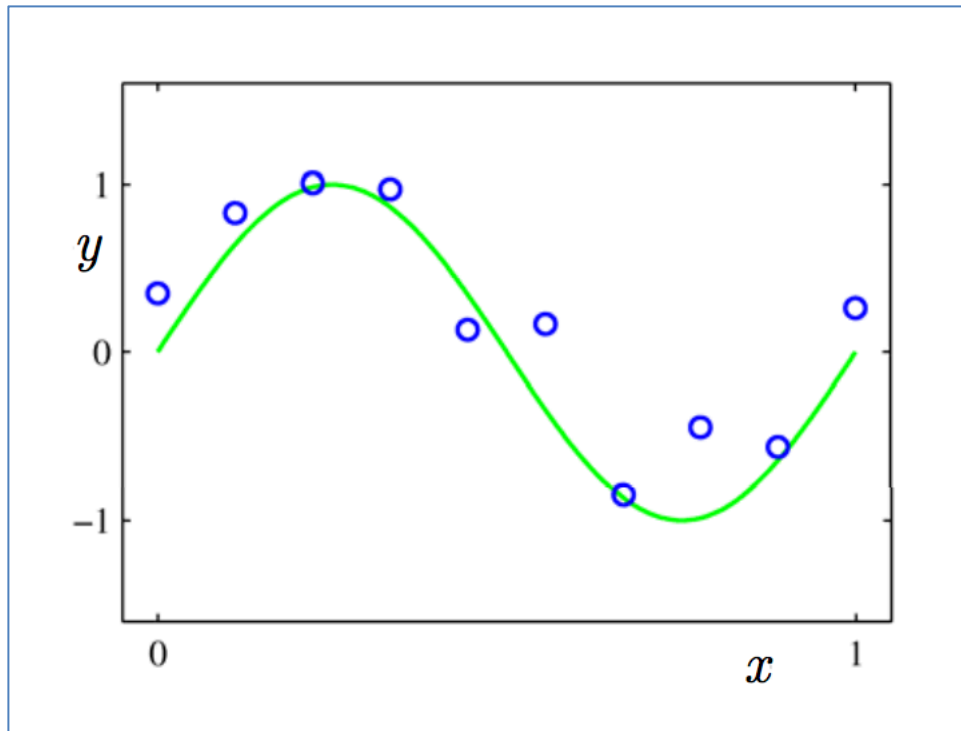
Toy example, iteration 100

iteration: 100, cost: 0.000105



Polynomial Regression

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



Recap
Normal
Equations

Closed
form

$$J(\vec{w}) = \frac{1}{2} \sum_{i=1}^n (\underbrace{\vec{w}^T \vec{x}_i}_{\hat{y}_i} - y_i)^2 \quad \left. \vphantom{\sum_{i=1}^n} \right\} \text{cost function, want to minimize!}$$

$$= \frac{1}{2} (\underbrace{X \vec{w} - \vec{y}}_{1 \times n})^T (\underbrace{X \vec{w} - \vec{y}}_{n \times 1})$$

$$= \frac{1}{2} (\underbrace{\vec{w}^T X^T X \vec{w}}_{(X^T X) "w^2" \text{ const}} - \underbrace{\vec{w}^T X^T \vec{y}}_{\text{const}} - \underbrace{\vec{y}^T X \vec{w}}_{\text{const}} + \underbrace{\vec{y}^T \vec{y}}_{\text{const}})$$

equal

$$\begin{matrix} (A B)^T & = & B^T A^T \\ \uparrow & \uparrow & \uparrow & \uparrow \\ n \times k & k \times m & m \times k & k \times n \\ & & m \times n & \end{matrix}$$

$$(A+B)^T = A^T + B^T$$

$$\left[\begin{array}{c} \bullet \\ \bullet \\ \bullet \end{array} \right]_{1 \times n} \left[\begin{array}{c} \bullet \\ \bullet \\ \bullet \end{array} \right]_{n \times 1} \quad \left. \vphantom{\left[\begin{array}{c} \bullet \\ \bullet \\ \bullet \end{array} \right]} \right\} \text{sum of squares!}$$

$$\left[\begin{array}{c} \bullet \\ \bullet \\ \bullet \end{array} \right] \left[\begin{array}{c} \bullet \\ \bullet \\ \bullet \end{array} \right]^T$$

$$\Rightarrow \hat{\beta} = \underbrace{(X^T X)^{-1}}_{\text{minimum var}(x)} \underbrace{X^T \vec{y}}_{\text{cov}(x, y)}$$

minimum $\nabla^2 f(x, y)$
(check with second derivative)

$$w_1 = \frac{\text{Cov}(x, y)}{\text{Var}(x)}$$

$$A^{-1}A \cdot x = A^{-1}B$$

$$X = A^{-1}B$$

SGD

SGD usually set $\vec{w} = \vec{0}$ \uparrow zero vector.

shuffle range(n): up

for i in range(n):

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

gradient of

$\underbrace{(\vec{w}^T \vec{x}_i - y_i) \vec{x}_i}_{\text{gradient of } J \text{ wrt one } \vec{x}_i}$

$$\begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \leftarrow \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} - \underbrace{\varphi(\bar{w}^T \bar{x}_i - y_i)}_{\text{Scalar}} \begin{bmatrix} 1 \\ x_i \end{bmatrix}$$

w_j high (pos) \Rightarrow feature positively correlated w/
 w_j large (neg) \Rightarrow " output
 $w_j \approx 0 \Rightarrow$ feature is uninformative.

Runtime

T = # of SGD iterations

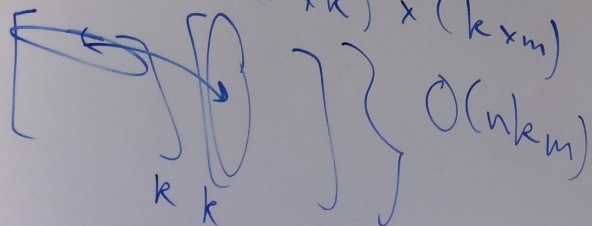
n = # of examples

p = # features

inverse : k dim matrix:

$$\Rightarrow O(k^3)$$

matrix mult : $(n \times k) \times (k \times m)$



$$\text{SGD} = O(Tnp)$$

Normal
Eqs =

$$X^T X \left\{ \begin{matrix} O(np^2) \\ (p \times n)(n \times p) \end{matrix} \right.$$

$$(X^T X)^{-1} \left\{ O(p^3) \right.$$

$$(X^T X)^{-1} X^T \left\{ O(p^2 n) \right.$$

$$\begin{matrix} \nearrow & \nwarrow \\ (p \times p) & (p \times n) \end{matrix}$$

$$\cdot \vec{y} \left\{ O(pn) \right.$$

p very high \Rightarrow SGD!

o.w. \Downarrow

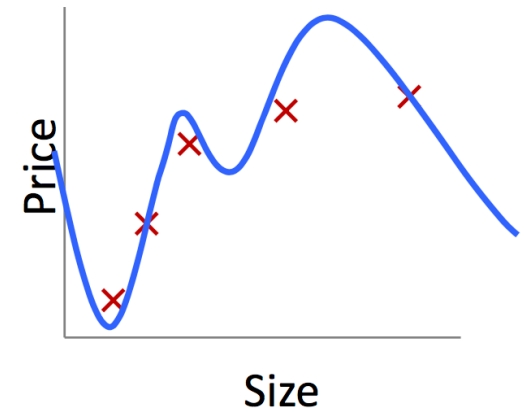
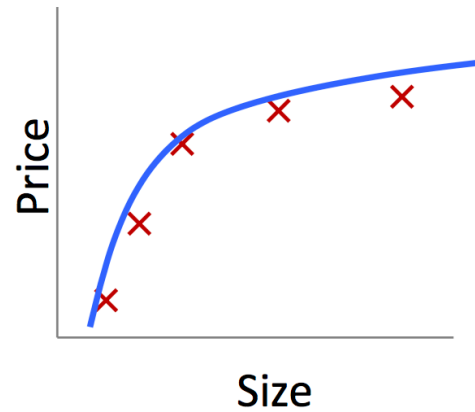
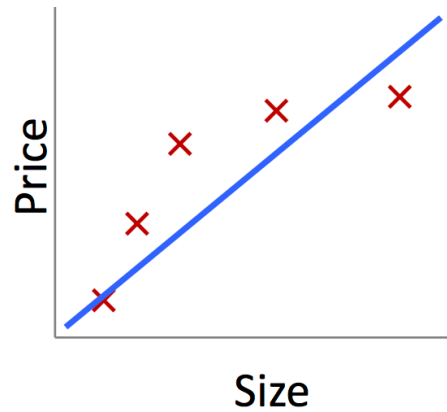
$$O(np^2 + p^3)$$

Outline for September 24

- Recap two ways of solving linear regression
 - SGD
 - Normal equations (closed-form solution)
- **Regularization**
- Introduction to conditional probability
 - Clinical trials example
- Begin: Naïve Bayes

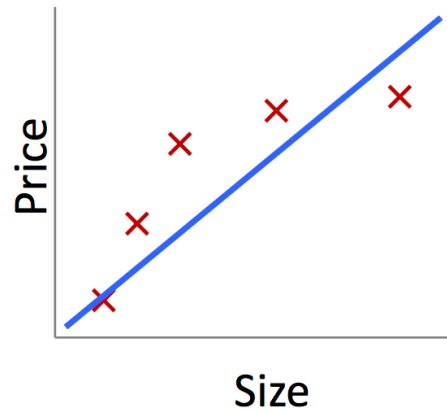
Generalization error

- Example: price vs. size (i.e. of a house or car)

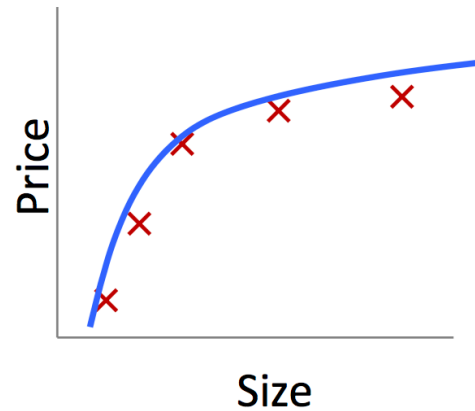


Generalization error

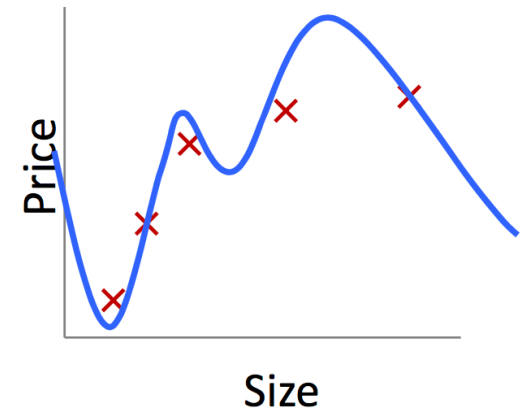
- Example: price vs. size (i.e. of a house or car)



underfitting
(high bias)



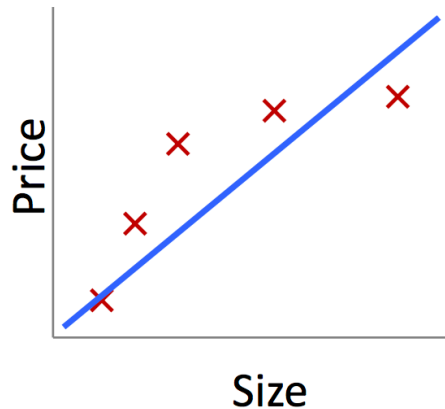
correct fit



overfitting
(high variance)

Generalization error

- Example: price vs. size (i.e. of a house or car)

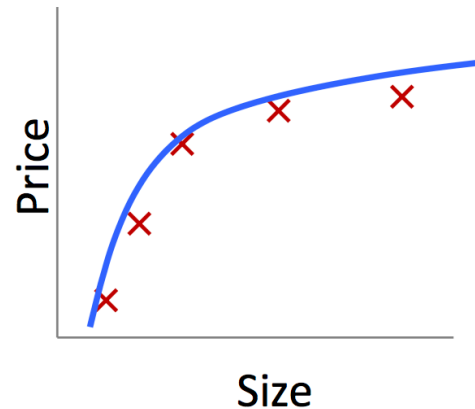


underfitting
(high bias)

Structural error:

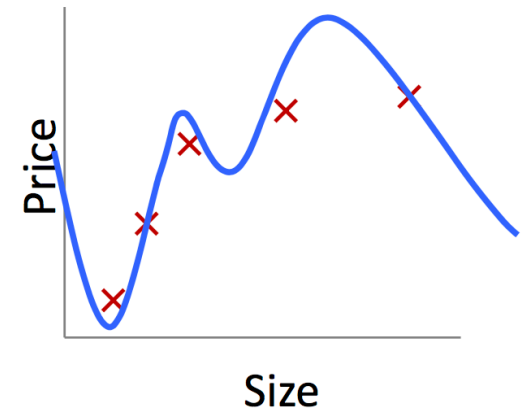
Hypothesis space cannot model true relationship

- More data doesn't help
- Need a more flexible model



correct fit

balance
↔



overfitting
(high variance)

Estimation (approximation) error:

Hypothesis space *can* model true relationship, BUT hard to identify correct model due to large hypothesis space, small n , or noise

- Reduce hypothesis space
- Add more data

Regularization

What if ...

- we have a limited # of training examples ($n < p$), or
- we want to automatically control the complexity of the learned hypothesis?

Regularization

What if ...

- we have a limited # of training examples ($n < p$), or
- we want to automatically control the complexity of the learned hypothesis?

Idea: penalize large values of w_j

Why prefer small weights?

- if large weights, small change in feature can result in large change in prediction
- prevent giving too much weight to any one feature
- might prefer zero weight for useless features

Common Regularizers

$$||\vec{w}||_0 = \sum_{j:w_j \neq 0} 1$$

L_0 norm

- Number of non-zero entries
- Minimizing L_0 norm is NP hard

Common Regularizers

$$||\vec{w}||_0 = \sum_{j:w_j \neq 0} 1$$

L_0 norm

- Number of non-zero entries
- Minimizing L_0 norm is NP hard

$$||\vec{w}||_1 = \sum_{j=1}^p |w_j|$$

L_1 norm

- Sum of magnitude of weights
- Not differentiable

Common Regularizers

$$||\vec{w}||_0 = \sum_{j:w_j \neq 0} 1$$

L_0 norm

- Number of non-zero entries
- Minimizing L_0 norm is NP hard

$$||\vec{w}||_1 = \sum_{j=1}^p |w_j|$$

L_1 norm

- Sum of magnitude of weights
- Not differentiable

$$||\vec{w}||_2 = \sqrt{\sum_{j=1}^p w_j^2}$$

L_2 norm

- Sum of squared weights
- Differentiable

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

Adding regularization to our cost function

SGD

$$\vec{w} \leftarrow \vec{w} - \alpha \left[(\vec{w}^T \vec{x}_i - y_i) \vec{x}_i + \lambda \vec{w} \right]$$

$$\vec{w} \leftarrow \vec{w} (1 - \alpha \lambda) \dots$$

Normal

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

identity

pulling weights down

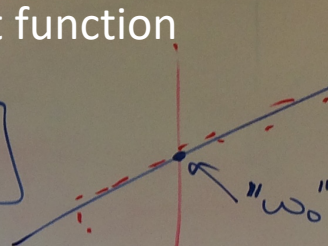
don't penalize w_0

$$\lambda \begin{bmatrix} \times & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

λ is small and positive

I^*

don't regularize "bias" shift away from 0



Outline for September 24

- Recap two ways of solving linear regression
 - SGD
 - Normal equations (closed-form solution)
- Regularization
- Introduction to conditional probability
 - Clinical trials example
- Begin: Naïve Bayes

Probability

$$\boxed{P(A, B) = P(A)P(B|A)}$$

↑
Bayes Rule.

$= P(A \text{ and } B)$

$$= P(B)P(A|B)$$

$$\boxed{P(A|B) = \frac{P(A)P(B|A)}{P(B)}}$$

$$P(\text{rain}) = 20\%$$

$$P(\text{umbrella} | \text{rain}) = 15\%$$

$$P(\text{U} | \text{R}) = ?$$

↑ "given"
↑ "given it is raining"

$$P(\text{U} | \text{R}) = \frac{P(\text{U}, \text{R})}{P(\text{R})}$$

$$= \frac{15\%}{20\%}$$

$$= 75\%$$