

# DS 5220

## Supervised Machine Learning and Learning Theory

Alina Oprea  
Associate Professor, CCIS  
Northeastern University

September 25 2019

# Outline

- Brief review
- Gradient descent
  - Batch algorithm
  - Line search optimization
- Gradient descent for linear regression
- Regularization
  - Ridge and Lasso regression
  - Gradient descent for ridge regression

# Review

- |   |  |
|---|--|
| <ul style="list-style-type: none"><li>• Regression<ul style="list-style-type: none"><li>– Linear regression</li><li>– MSE loss</li><li>– Closed-form solution</li><li>– Simple and multiple regression</li><li>– Bias-variance tradeoff</li></ul></li></ul> | <ul style="list-style-type: none"><li>• Classification<ul style="list-style-type: none"><li>– Linear models</li><li>– Perceptron</li><li>– Generative models:<br/>Linear Discriminant Analysis (LDA)</li></ul></li></ul> |
|---|--|

# Training algorithms

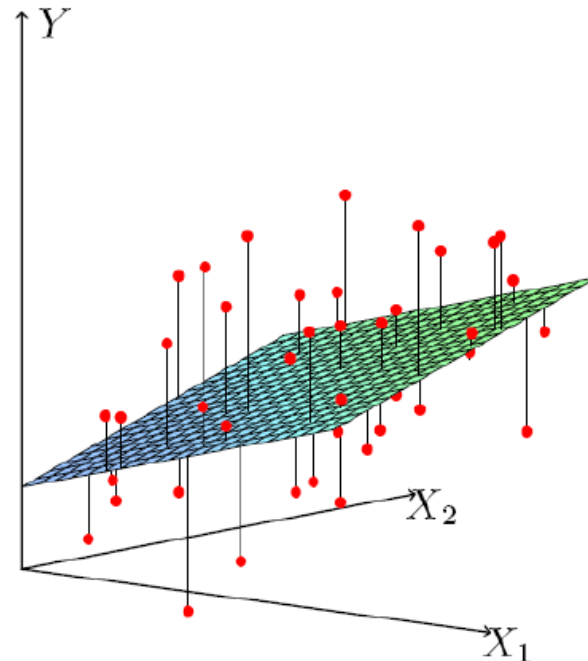
- Linear Regression
  - Minimize MSE
  - Compute closed-form solution (linear model that minimizes the MSE)
- LDA
  - Estimating probabilities of each class using Bayes Theorem
  - Learn normal distribution of data in each class
  - Estimate mean vector and covariance matrix

# Multiple Linear Regression

- Dataset:  $x_i \in R^d, y_i \in R$
- Hypothesis  $h_{\theta}(x) = \theta^T x$
- $MSE = \frac{1}{N} \sum (\theta^T x_i - y_i)^2$  **Loss / cost**

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

What are the drawbacks of computing the closed-form solution?



# How to optimize loss functions?

- Dataset:  $x_i \in R^d, y_i \in R$
- Hypothesis  $h_\theta(x) = \theta^T x$
- $J(\theta) = \frac{1}{N} \sum (\theta^T x_i - y_i)^2$  **Loss / cost**
  - Strictly convex function (unique minimum)
- **General method to optimize a multi-variate function**
  - Practical (low asymptotic complexity)
  - Convergence guarantees to global minimum

# What Strategy to Use?



# Follow the Slope



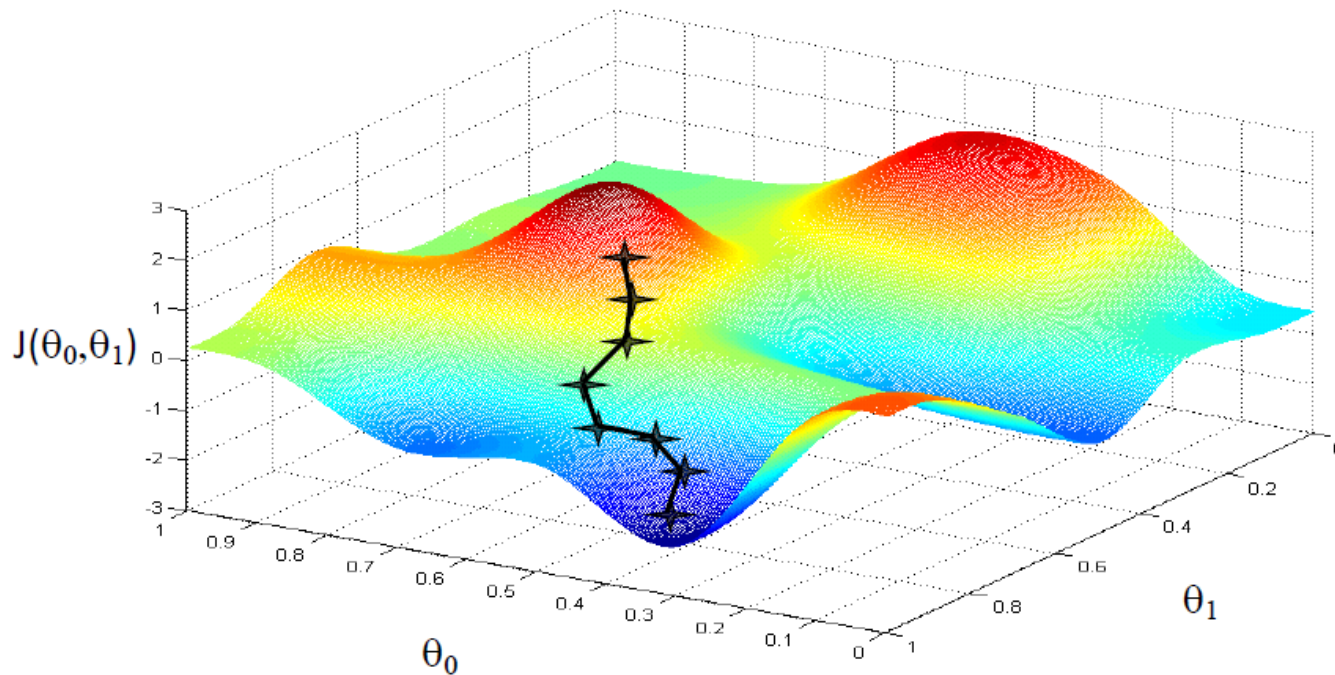
Follow the direction of steepest descent!



# How to optimize $J(\theta)$ ?

- Choose initial value for  $\theta$
- Until we reach a minimum:
  - Choose a new value for  $\theta$  to reduce  $J(\theta)$   $\longrightarrow$

Direction of  
steepest  
descent!



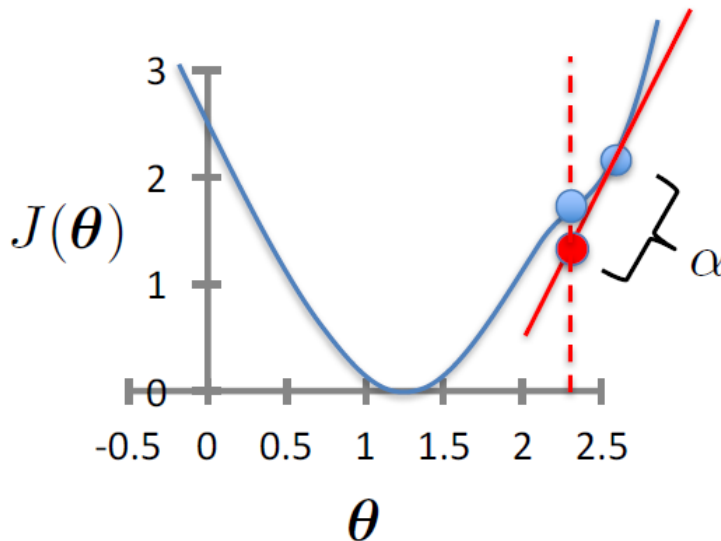
# Batch Gradient Descent

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

simultaneous update  
for  $j = 0 \dots d$

learning rate (small)  
e.g.,  $\alpha = 0.05$



- Gradient = slope of line tangent to curve
- Function decreases faster in negative direction of gradient
- Step is proportional to learning rate

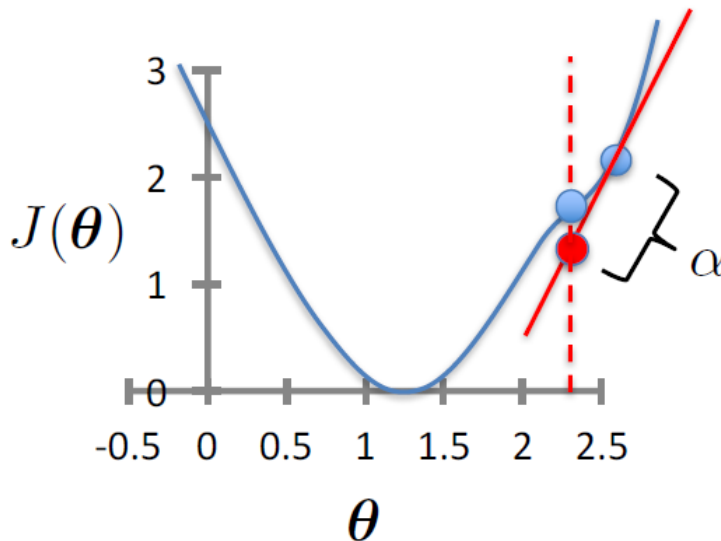
# Batch Gradient Descent

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

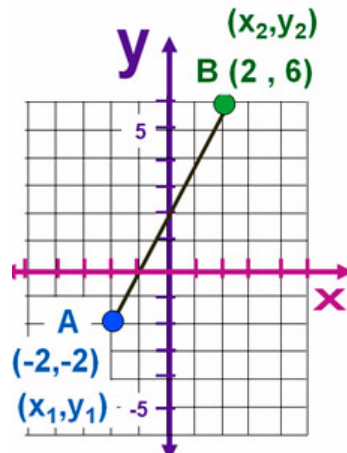
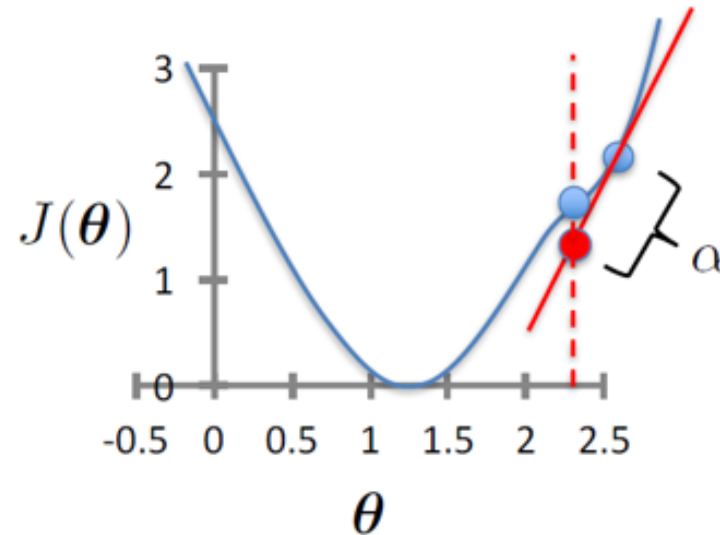
simultaneous update  
for  $j = 0 \dots d$

learning rate (small)  
e.g.,  $\alpha = 0.05$



$$\text{Vector update rule: } \theta \leftarrow \theta - \frac{\partial J(\theta)}{\partial \theta}$$

# Gradient Descent



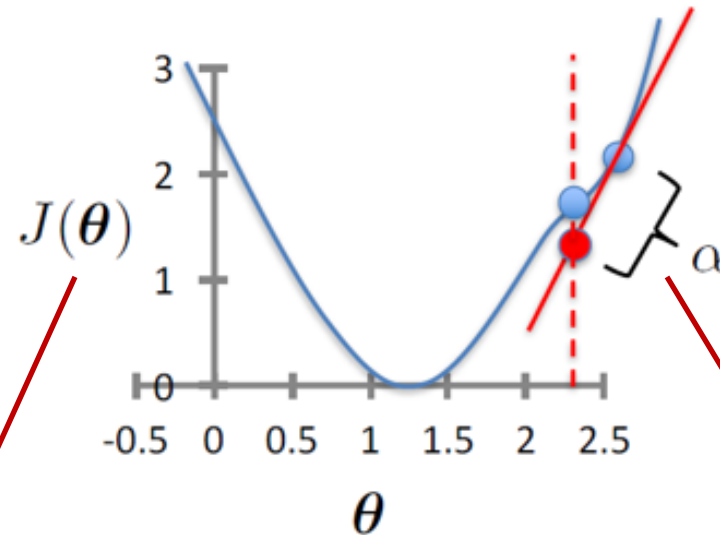
The Gradient "m" is:

$$m = \frac{y_2 - y_1}{x_2 - x_1} = \frac{\Delta Y}{\Delta X}$$

$$m = \frac{6 - -2}{2 - -2}$$

$$m = 8 / 4 = 2 \checkmark$$

# Gradient Descent



- If  $\theta$  is on the left of minimum, slope is negative
- Increase value of  $\theta$

- If  $\theta$  is on the right of minimum, slope is positive
- Decrease value of  $\theta$

In both cases  $\theta$  gets closer to minimum

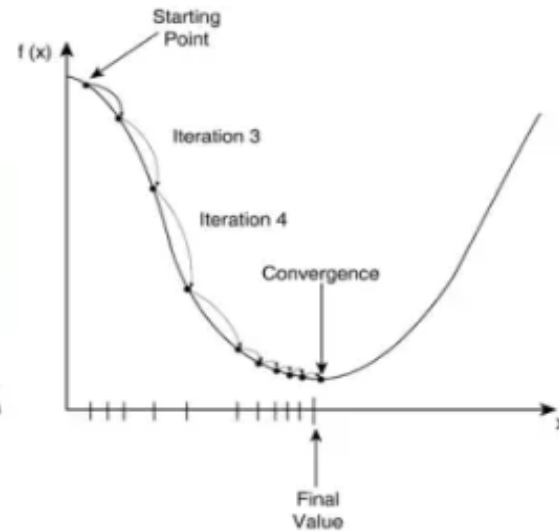
# Gradient Descent

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

simultaneous update  
for  $j = 0 \dots d$

learning rate (small)  
e.g.,  $\alpha = 0.05$



- As approach minimum, slope gets smaller (GD takes smaller steps)

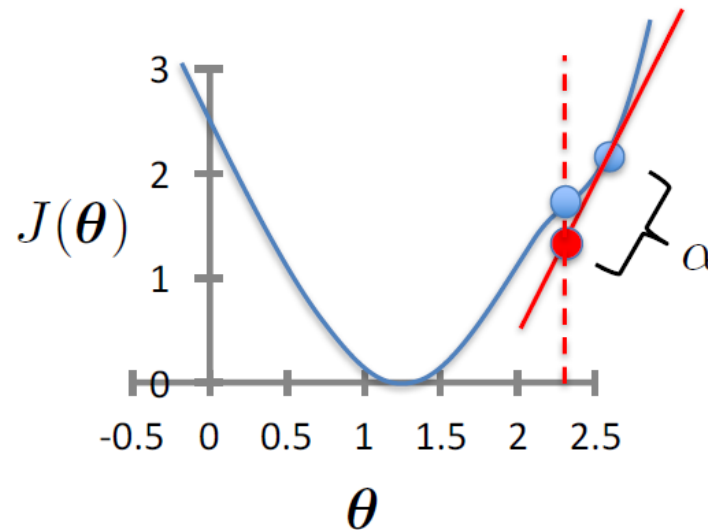
# Gradient Descent

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

simultaneous update  
for  $j = 0 \dots d$

learning rate (small)  
e.g.,  $\alpha = 0.05$



- What happens when  $\theta$  reaches a local minimum?
- The slope is 0, and gradient descent converges!
- Strictly convex functions only have global minimum

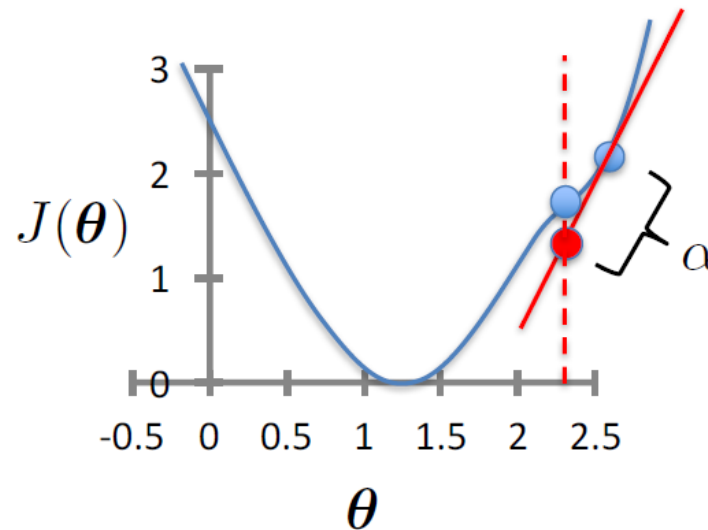
# Stopping Condition

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

simultaneous update  
for  $j = 0 \dots d$

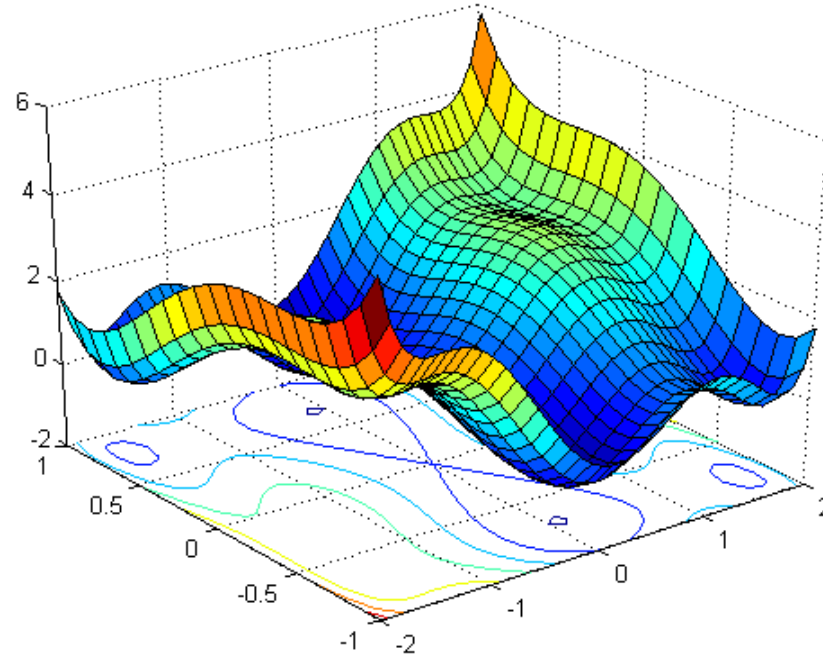
learning rate (small)  
e.g.,  $\alpha = 0.05$



- When should the algorithm stop?
- When the update in  $\theta$  is below some threshold

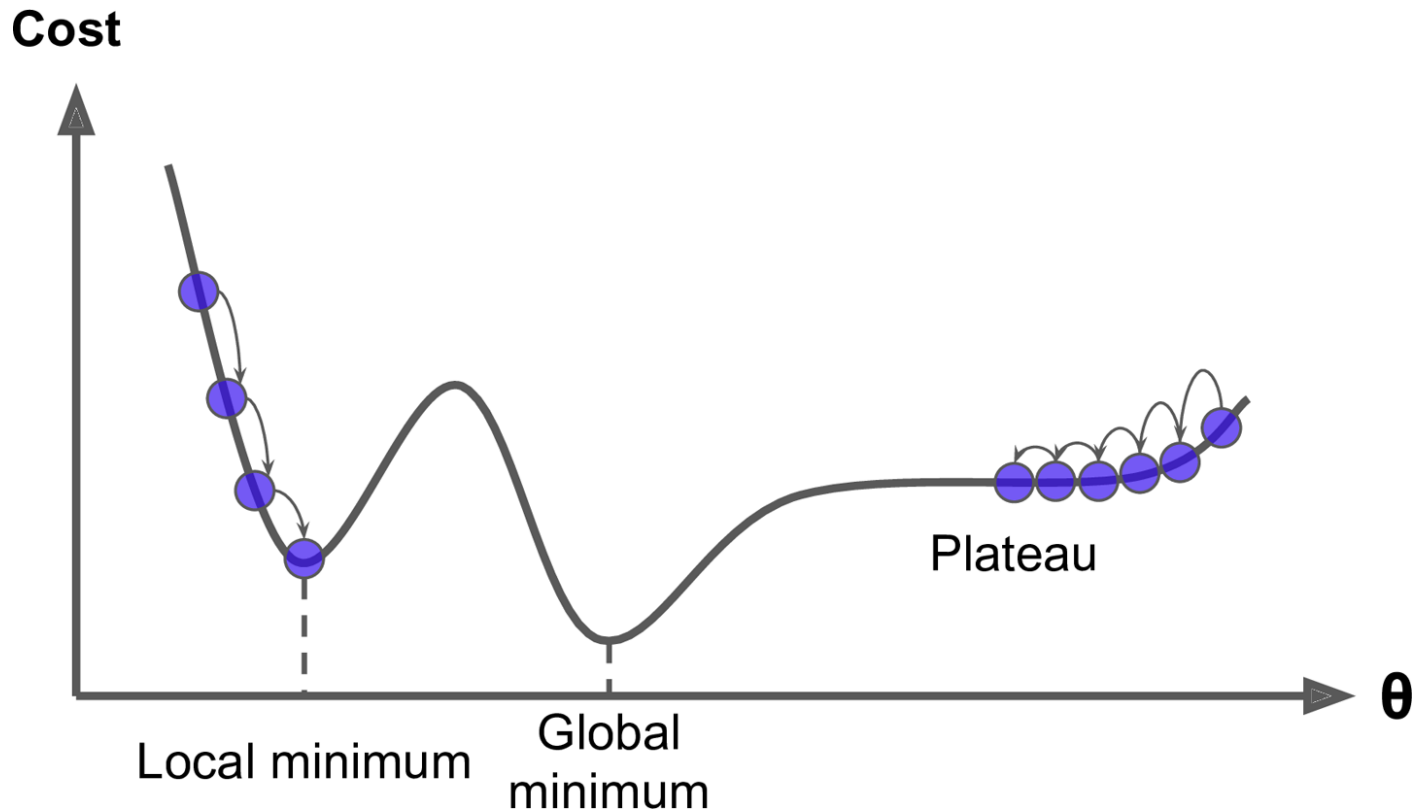


# Complex loss function



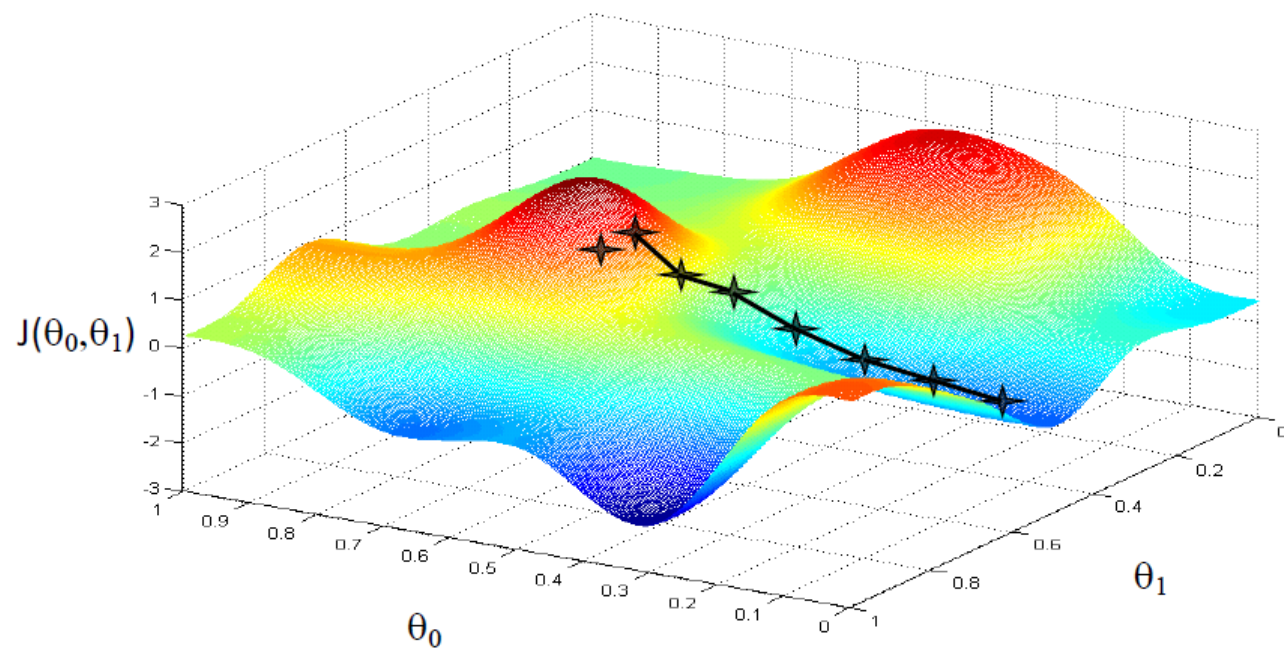
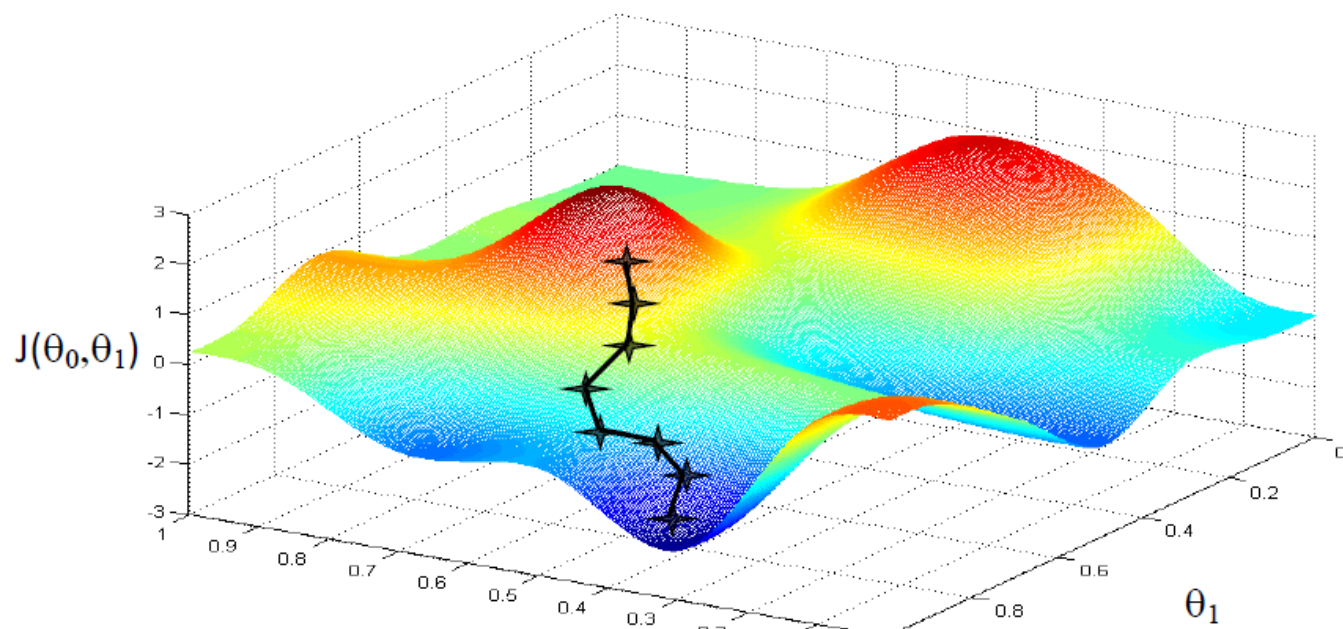
- Complex loss functions are more difficult to optimize

# GD Convergence Issues



- Local minimum: Gradient descent stops
- Plateau: Almost flat region where slope is small

**Solution: start from multiple random locations**



# Simple Linear Regression

- Dataset  $x_i \in R, y_i \in R, h_{\theta}(x) = \theta_0 + \theta_1 x$
- $J(\theta) = \frac{1}{N} \sum_{i=1}^n (\theta_0 + \theta_1 x_i - y_i)^2$

**MSE / Loss**

- Solution of min loss

$$\begin{aligned} -\theta_0 &= \bar{y} - \theta_1 \bar{x} \\ -\theta_1 &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \end{aligned}$$

**Variance of x**

**Co-variance of x and y**

$$\begin{aligned} \bar{x} &= \frac{\sum_{i=1}^n x_i}{n} \\ \bar{y} &= \frac{\sum_{i=1}^n y_i}{n} \end{aligned}$$

# GD for Simple Linear Regression

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

simultaneous update  
for  $j = 0 \dots d$

- $J(\theta) = \frac{1}{N} \sum_{i=1}^N (\theta_0 + \theta_1 x_i - y_i)^2$
- $\frac{\partial J(\theta)}{\partial \theta_0} = \frac{2}{N} \sum_{i=1}^N (\theta_0 + \theta_1 x_i - y_i) = \frac{2}{N} \sum_{i=1}^N (h_{\theta}(x_i) - y_i)$
- $\frac{\partial J(\theta)}{\partial \theta_1} = \frac{2}{N} \sum_{i=1}^N (h_{\theta}(x_i) - y_i) x_i$

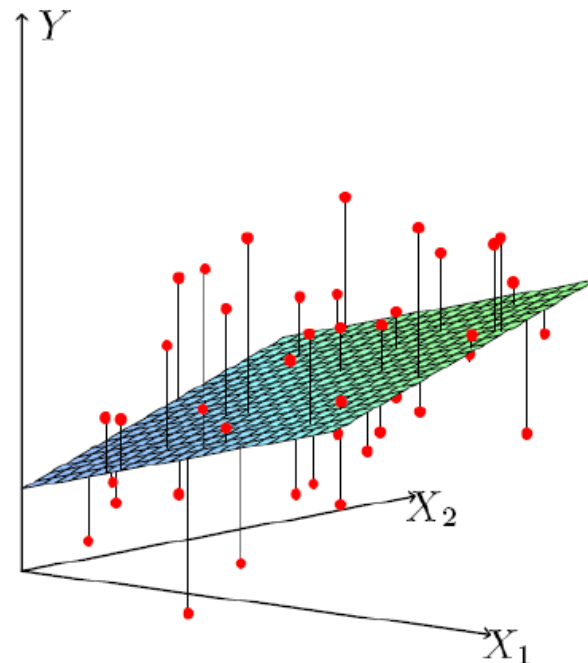
**Batch: Update of each parameter  
component depends on all training data**

# Multiple Linear Regression

- Dataset:  $x_i \in R^d, y_i \in R$
- Hypothesis  $h_{\theta}(x) = \theta^T x$
- $MSE = \frac{1}{N} \sum (\theta^T x_i - y_i)^2$  Loss / cost

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

MSE is a strictly convex function  
and has unique minimum



# GD for Multiple Linear Regression

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

simultaneous update  
for  $j = 0 \dots d$

- $J(\theta) = \frac{1}{N} \sum_{i=1}^N (\sum_k \theta_k x_{ik} - y_i)^2$
- $$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta_j} &= \frac{2}{N} \sum_{i=1}^N (\sum_k \theta_k x_{ik} - y_i) \frac{\partial (\sum_k \theta_k x_{ik} - y_i)}{\partial \theta_j} \\ &= \frac{2}{N} \sum_{i=1}^N (h_{\theta}(x_i) - y_i) x_{ij} \end{aligned}$$

# GD for Linear Regression

- Initialize  $\theta$
- Repeat until convergence  $\|\theta_{new} - \theta_{old}\| < \epsilon$  or  $iterations == MAX\_ITER$

$$\theta_j \leftarrow \theta_j - \alpha \frac{2}{N} \sum_{i=1}^N (h_{\theta}(x_i) - y_i) x_{ij}$$

simultaneous  
update  
for  $j = 0 \dots d$

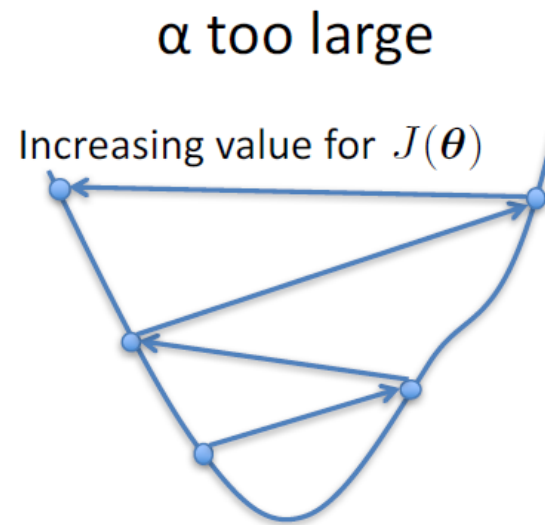
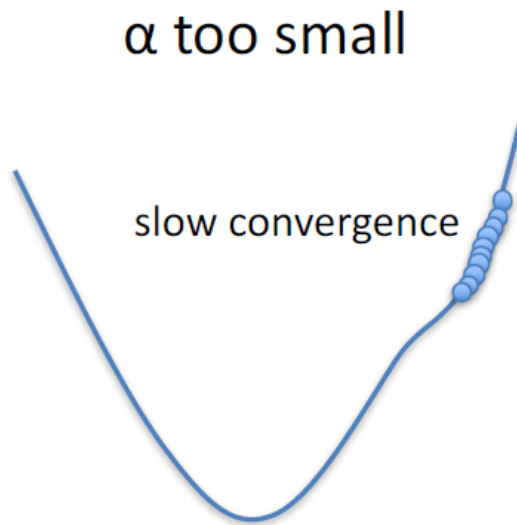
- Assume convergence when  $\|\theta_{new} - \theta_{old}\|_2 < \epsilon$

$$\text{L}_2 \text{ norm: } \|v\|_2 = \sqrt{\sum_i v_i^2} = \sqrt{v_1^2 + v_2^2 + \dots + v_{|v|}^2}$$

Can also bound number of iterations



# Choosing learning rate

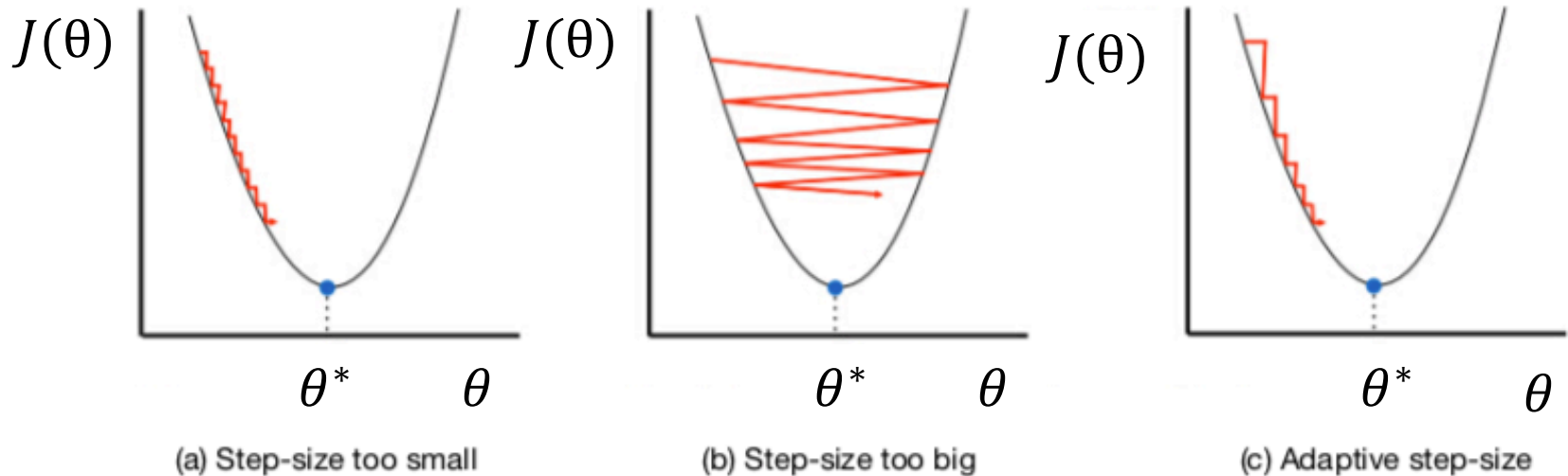


- May overshoot the minimum
- May fail to converge
- May even diverge

To see if gradient descent is working, print out  $J(\theta)$  each iteration

- The value should decrease at each iteration
- If it doesn't, adjust  $\alpha$

# Adaptive step size



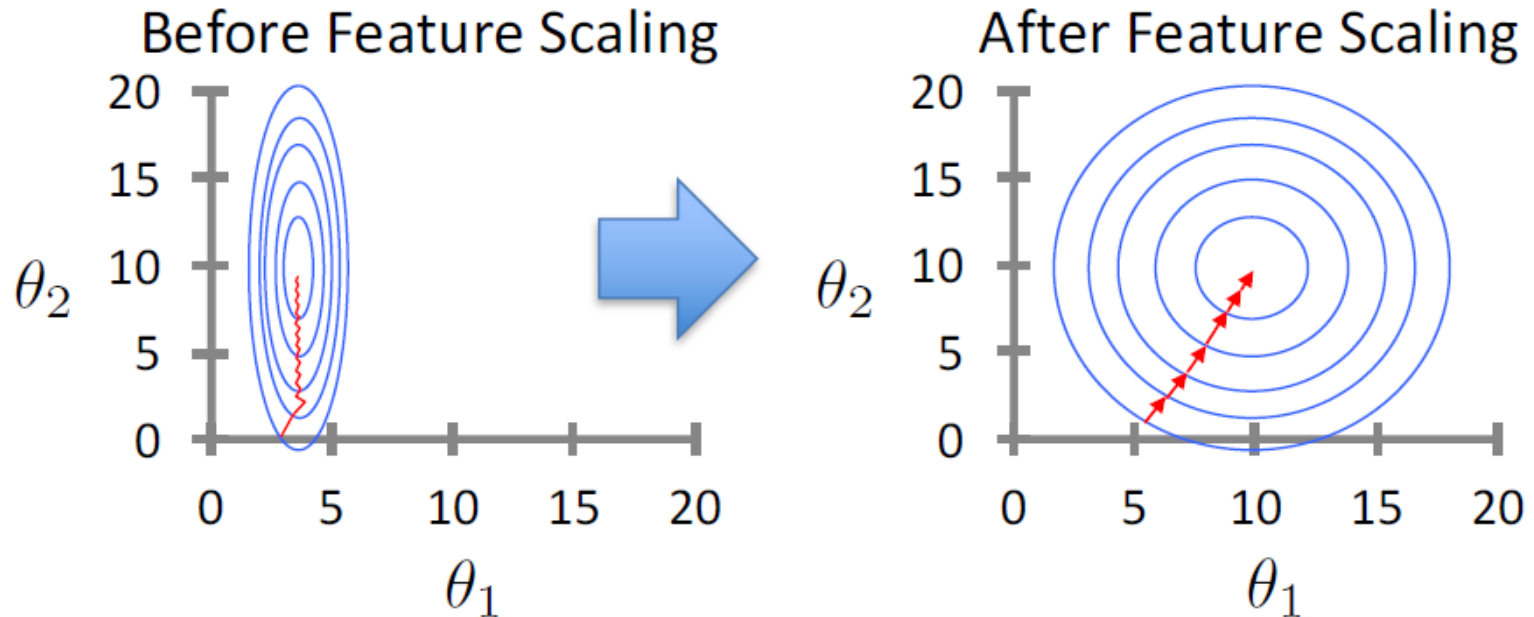
- Start with large step size and reduce over time, adaptively
- Measure how objective decreases

# Gradient Descent with Line Search

- Input:  $\alpha_{max}$  max. learning rate  
 $\tau \in [0.5, 0.9]$ ,  $\epsilon$  (tolerance),  $T$  (backtrack)
- Initialize  $\theta$  with random value
- $\alpha \leftarrow \alpha_{max}$  // Set at maximum learning rate
- Repeat until convergence
  - $\theta_{try} \leftarrow \theta$
  - Repeat max  $T$  times // Maximum number backtrack
  - $\theta_j^{try} \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$  for all  $j$ 
    - If  $J(\theta) - J(\theta_{try}) > \epsilon$ : then  $\theta \leftarrow \theta_{try}$ ; break // Improved objective
    - Else  $\alpha \leftarrow \tau \alpha$  (reduce step size). // Backtrack to smaller rate
  - If  $T$  times is reached, break and start over

# Feature Scaling

- **Idea:** Ensure that feature have similar scales



- Makes gradient descent converge *much* faster

# Gradient Descent in Practice

- Asymptotic complexity
  - $O(NTd)$ ,  $N$  is size of training data,  $d$  is feature dimension, and  $T$  is number of iterations
- Most popular optimization algorithm in use today
- At the basis of training
  - Linear Regression
  - Logistic regression
  - SVM
  - Neural networks and Deep learning
  - Stochastic Gradient Descent variants

# Gradient Descent vs Closed Form

## Gradient Descent

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

simultaneous update  
for  $j = 0 \dots d$

## Closed form

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

### • Gradient Descent

- + Linear increase in  $d$  and  $N$
- + Generally applicable
- Need to choose  $\alpha$  and stopping conditions
- Might get stuck in local optima

### • Closed Form

- + No parameter tuning
- + Gives the global optimum
- Not generally applicable
- Slow computation:  $O(d^3)$

# Issues with Gradient Descent

- Might get stuck in local optimum and not converge to global optimum
  - Restart from multiple initial points
- Only works with differentiable loss functions
- Small or large gradients
  - Feature scaling helps
- Tune learning rate
  - Can use line search for determining optimal learning rate

# Review Gradient Descent

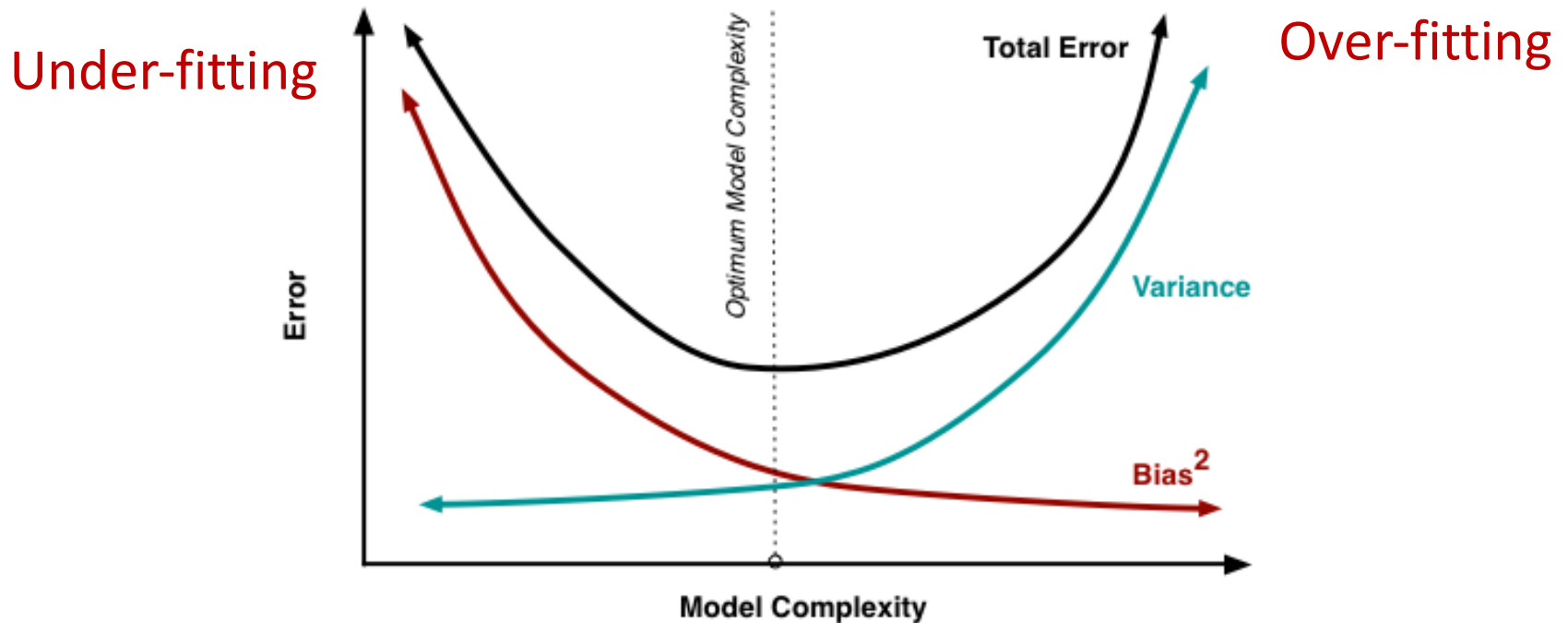
- Gradient descent is an efficient algorithm for optimization and training ML models
  - The most widely used algorithm in ML!
  - Much faster than using closed-form solution for linear regression
  - Main issues with Gradient Descent is convergence and getting stuck in local optima (for neural networks)
- Gradient descent is guaranteed to converge to optimum for strictly convex functions if run long enough



# Outline

- Brief review
- Gradient descent
  - Batch algorithm
  - Line search optimization
- Gradient descent for linear regression
- Regularization
  - Ridge and Lasso regression
  - Gradient descent for ridge regression

# Bias-Variance Tradeoff



- Bias = Difference between estimated and true models
- Variance = Model difference on different training sets

**MSE is proportional to Bias + Variance**

# Regularization

- A method for automatically controlling the complexity of the learned hypothesis
- **Idea:** penalize for large values of  $\theta_j$ 
  - Can incorporate into the cost function
  - Works well when we have a lot of features, each that contributes a bit to predicting the label
- Can also address overfitting by eliminating features (either manually or via model selection)

Reduce model complexity

Reduce model variance

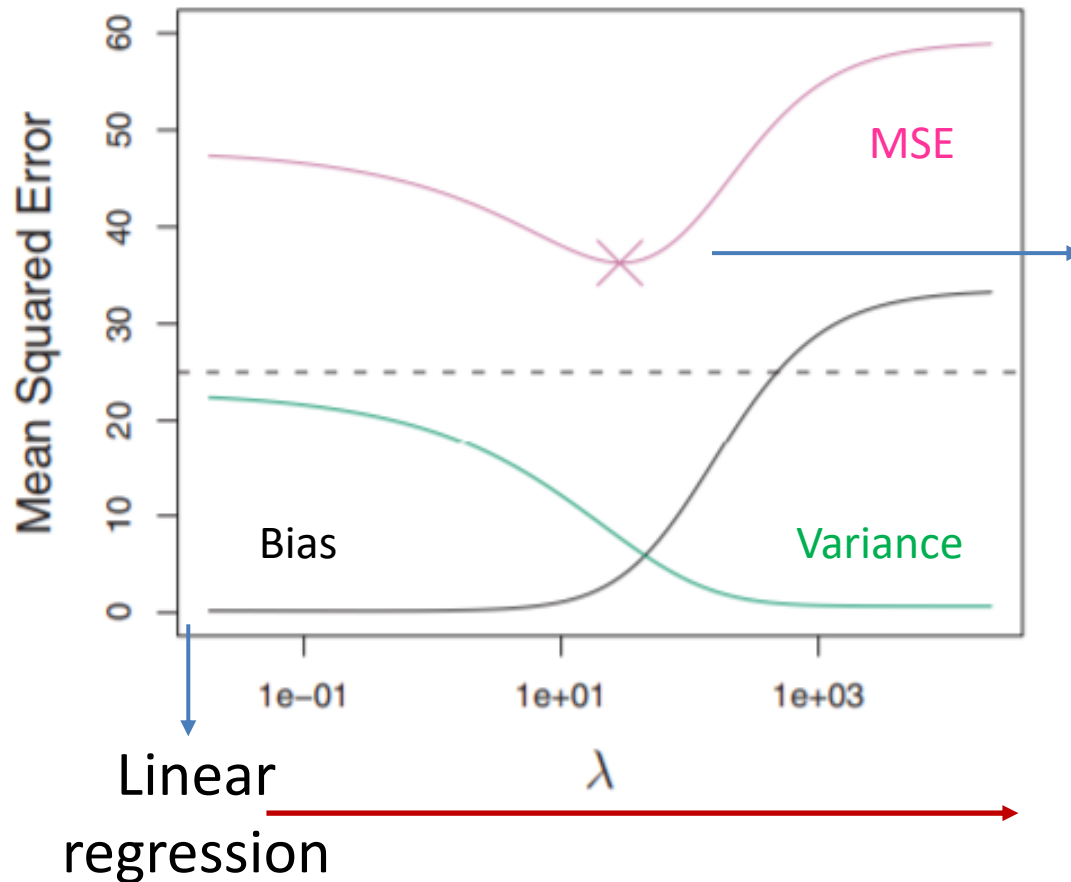
# Ridge regression

- Linear regression objective function

$$J(\theta) = \underbrace{\frac{1}{2} \sum_{i=1}^N (h_{\theta}(x_i) - y_i)^2}_{\text{model fit to data}} + \underbrace{\frac{\lambda}{2} \sum_{j=1}^d \theta_j^2}_{\text{regularization}}$$

- $\lambda$  is the regularization parameter (  $\lambda \geq 0$  )
  - No regularization on  $\theta_0$ !
- If  $\lambda = 0$ , we train linear regression
  - If  $\lambda$  is large, the coefficients will shrink close to 0

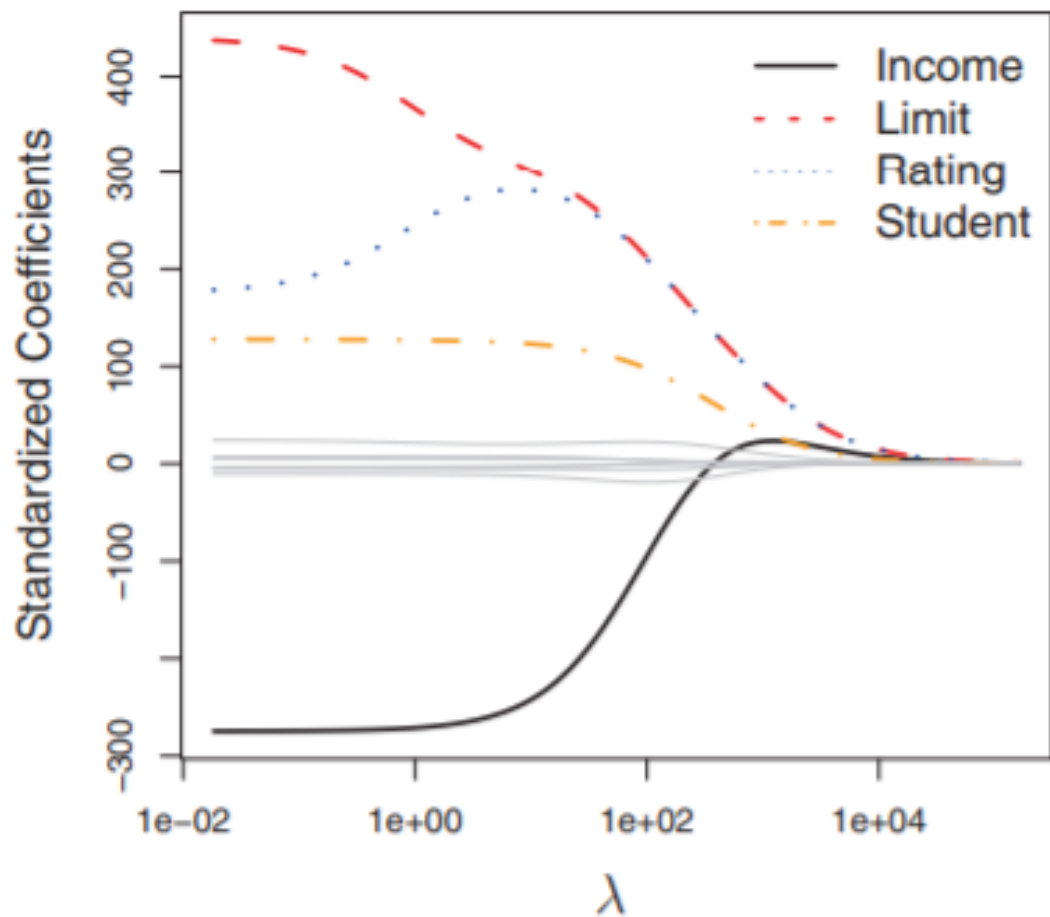
# Bias-Variance Tradeoff



Optimal  
Ridge regression

Reduced model  
complexity

# Coefficient shrinkage



Predict credit card balance

# GD for Ridge Regression

Min MSE

$$J(\theta) = \frac{1}{2} \sum_{i=1}^N (h_{\theta}(x_i) - y_i)^2 + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$

Gradient update:  $\theta_0 \leftarrow \theta_0 - \alpha \sum_{i=1}^N (h_{\theta}(x_i) - y_i)$

$$\theta_j \leftarrow \theta_j - \alpha \sum_{i=1}^N (h_{\theta}(x_i) - y_i) x_{ij} - \underbrace{\alpha \lambda \theta_j}_{\text{Regularization}}$$

$$\theta_j \leftarrow \theta_j (1 - \alpha \lambda) - \alpha \sum_{i=1}^N (h_{\theta}(x_i) - y_i) x_{ij}$$

# Review

- Gradient descent is a general optimization algorithm that can be applied in training
  - Minimize loss function (error metric)
- Gradient descent converges to local minima
- Requires selection of learning rate
- Bias-Variance tradeoff shows that both bias and variance need to be minimized
  - Regularization is a method to reduce model complexity and avoid overfitting
  - Ridge and Lasso regularization can be added to any loss function
  - Regularized model trained with Gradient Descent



# Acknowledgements

- Slides made using resources from:
  - Andrew Ng
  - Eric Eaton
  - David Sontag
- Thanks!