

REGRESSION

INTRODUCTION TO DATA SCIENCE

TIM KRASKA

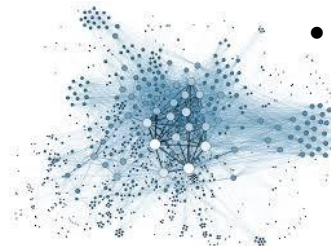


Recap



Data Management

- Relational Model
- SQL
- Star Schema
- ..



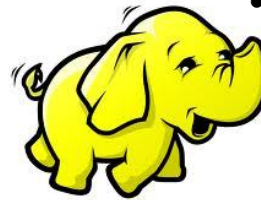
Visualization

- Types
- Senses
- Design goals
- ...



Data Clearning

- Schema Matching
- Entity Resolution
- Data Fusion
- Jaccard Similarity
- IDF Weight
- ...



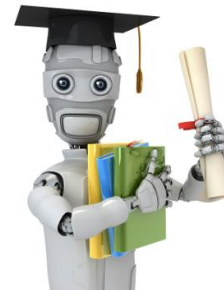
Map/Reduce

- Abstraction
- Implementation
- Spark
- ...



Information Retrieval

- Text Search
- Tokenization/Stemming
- Inverted Index
- Ranking
- Precision/Recall
- ...



Machine Learning

- Naïve Bayes
- Kmeans
- Page Rank
- Testing/Cross-Validation
- ..



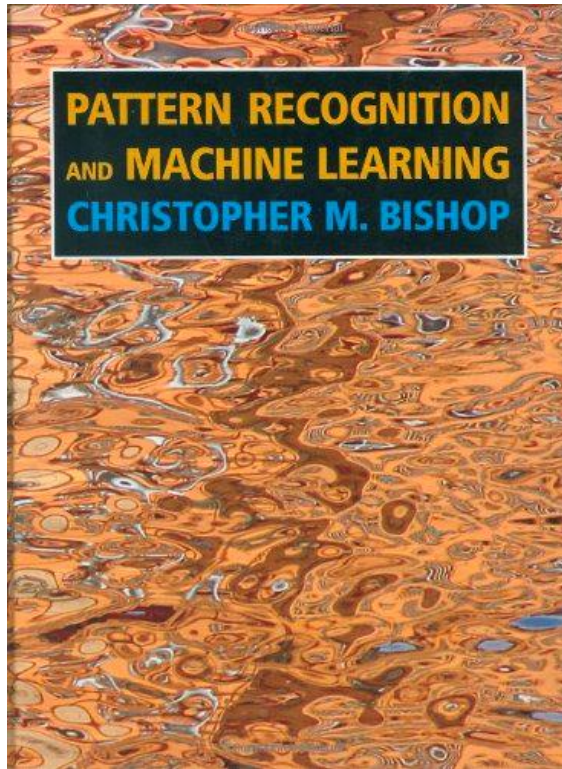
Probabilities and Statistics

- Law of Large Numbers
- Central Limit Theorem
- Testing
- ...

MATRIX

	<i>Supervised Learning</i>	<i>Unsupervised Learning</i>
<i>Discrete</i>	<div>classification or categorization</div>	clustering
<i>Continuous</i>	regression	dimensionality reduction

RECOMMENDED READING



Chapter 3 (math heavy, but complete)

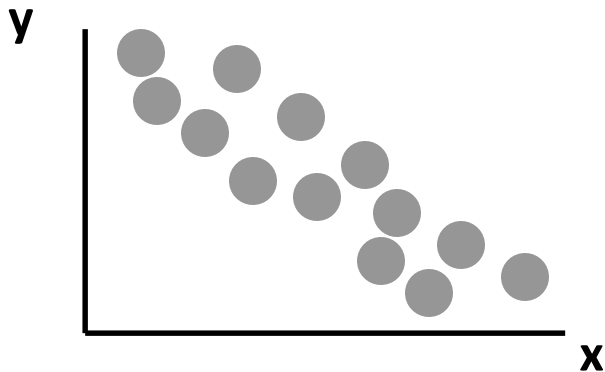
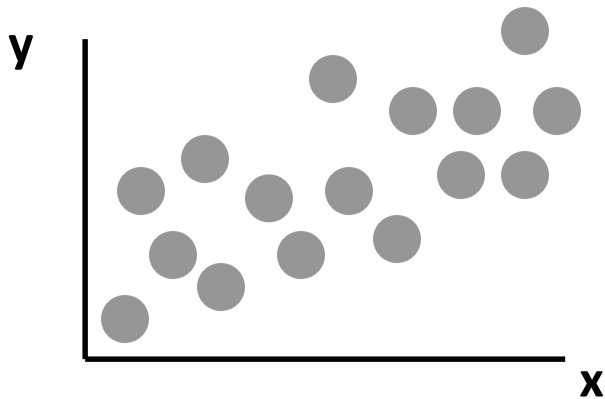
A bit easier to digest:

<http://cs229.stanford.edu/notes/cs229-notes1.pdf>

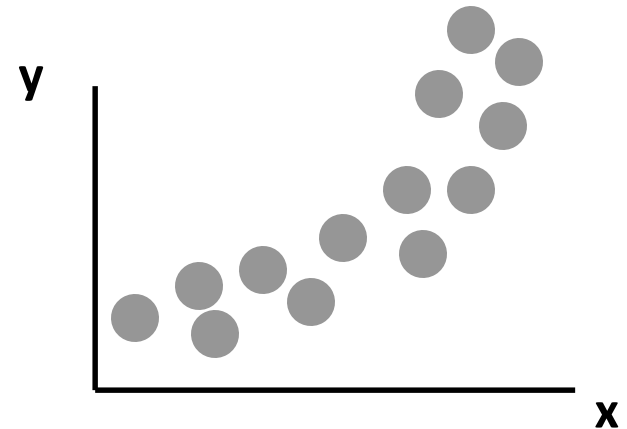
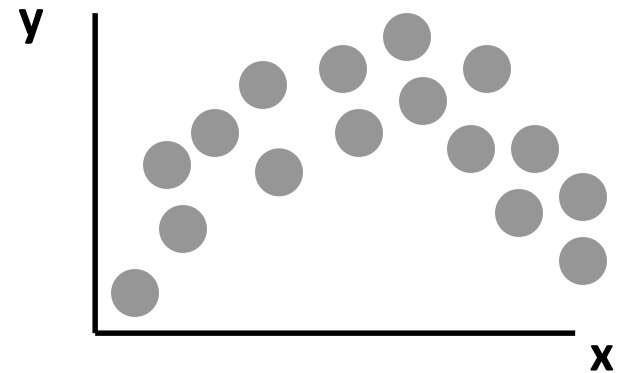
LINEAR REGRESSION

SCATTER PLOT EXAMPLES

Linear relationships



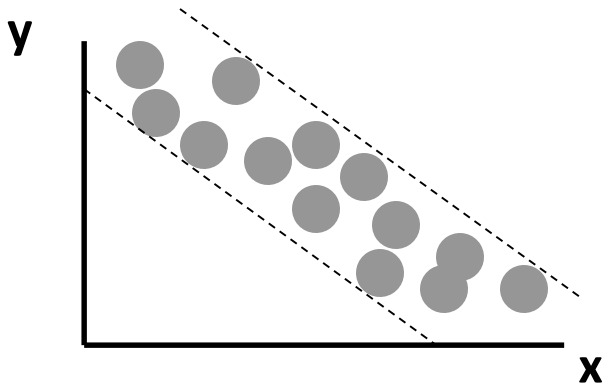
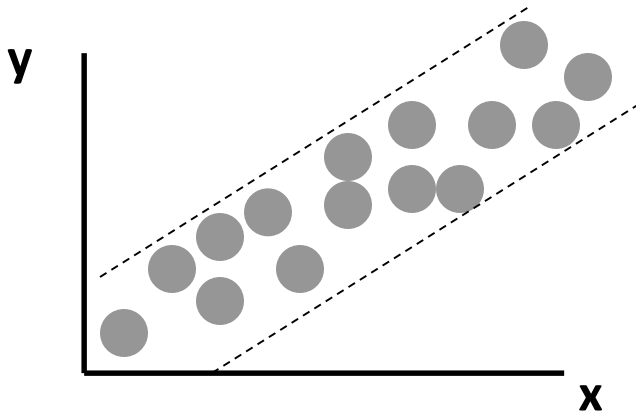
Curvilinear relationships



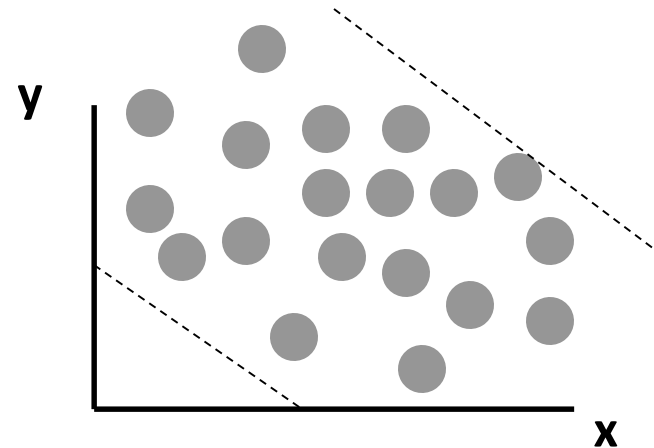
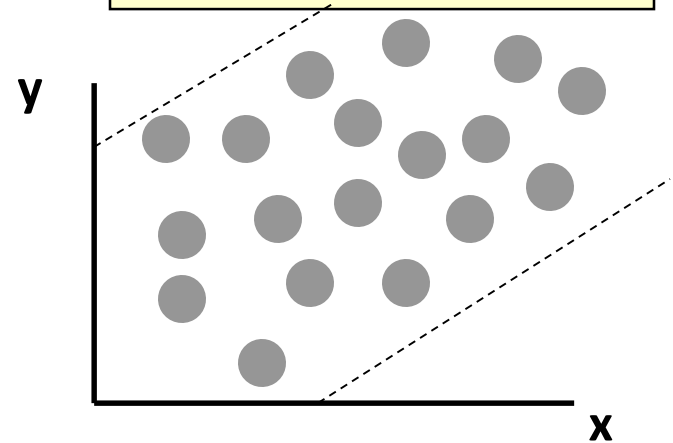
SCATTER PLOT EXAMPLES

(continued)

Strong relationships



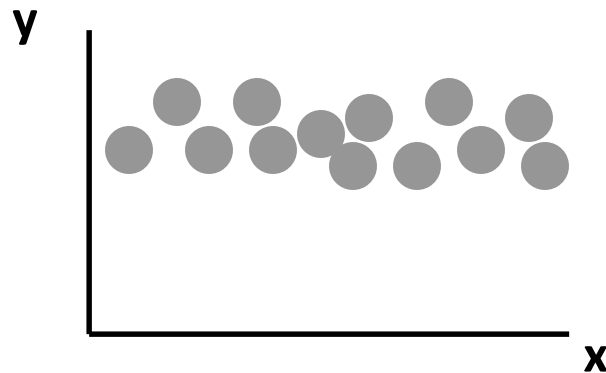
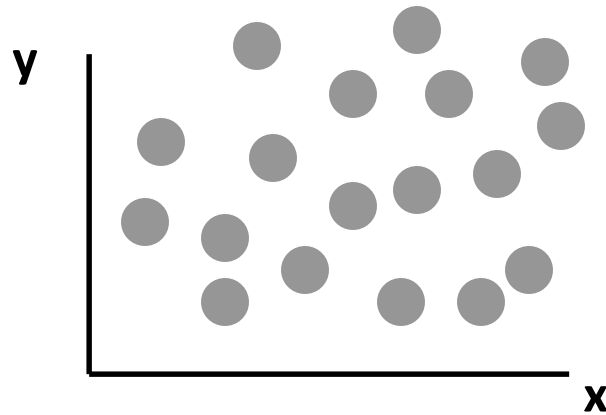
Weak relationships



SCATTER PLOT EXAMPLES

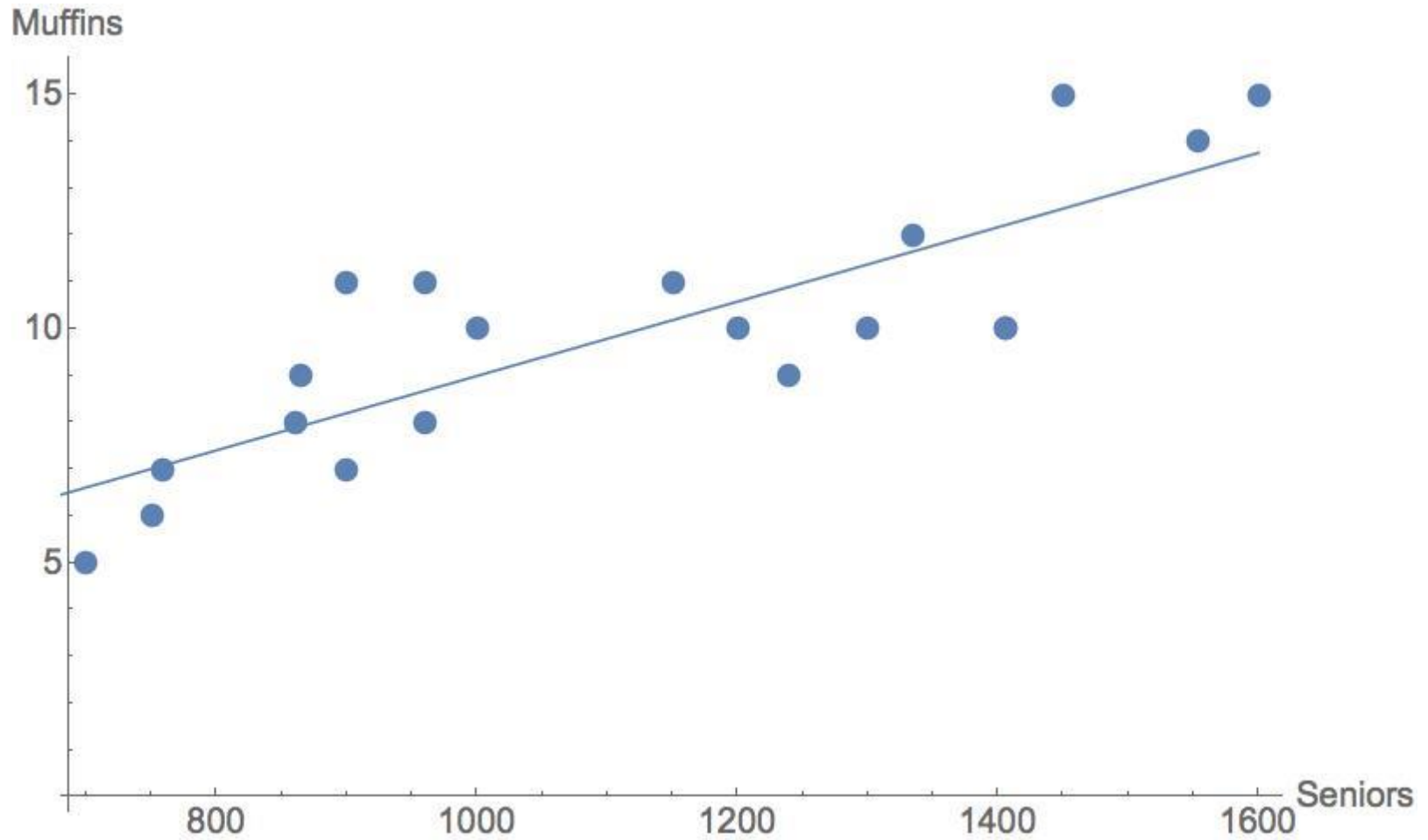
(continued)

No relationship



Naked Cup-Cake Run

Naked Cup-Cake Run



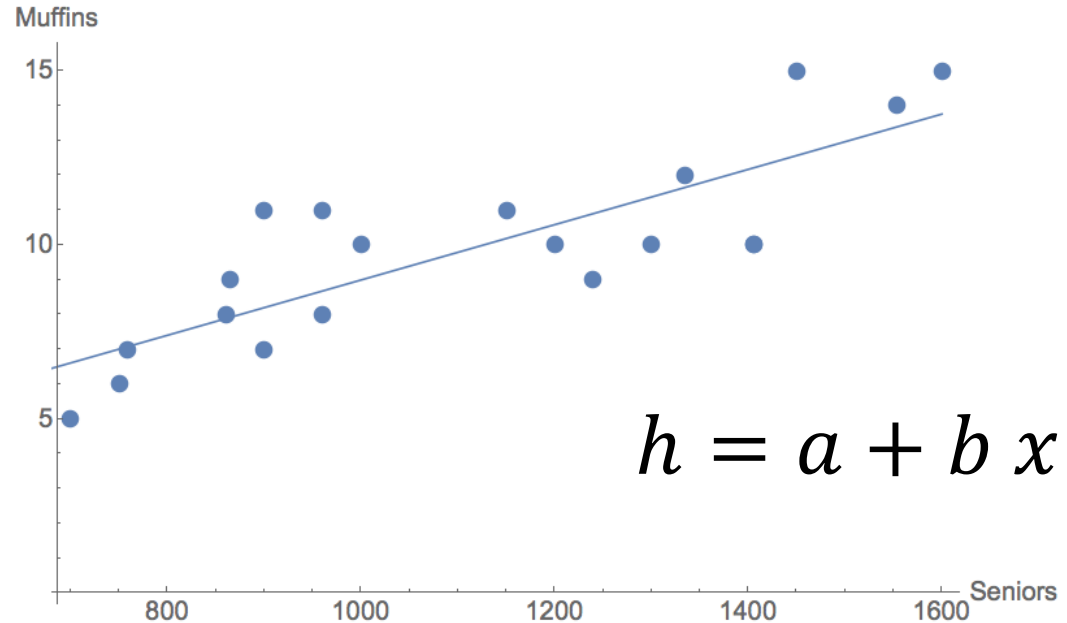
Naked Cup-Cake Run

How do we represent the problem:

$$q = \begin{pmatrix} q_0 \\ q_1 \end{pmatrix}$$

$$x = \begin{pmatrix} 1 \\ x_1 \end{pmatrix} = \begin{pmatrix} 1 \\ \text{Seniors} \end{pmatrix}$$

intercept term ($x_0=1$)



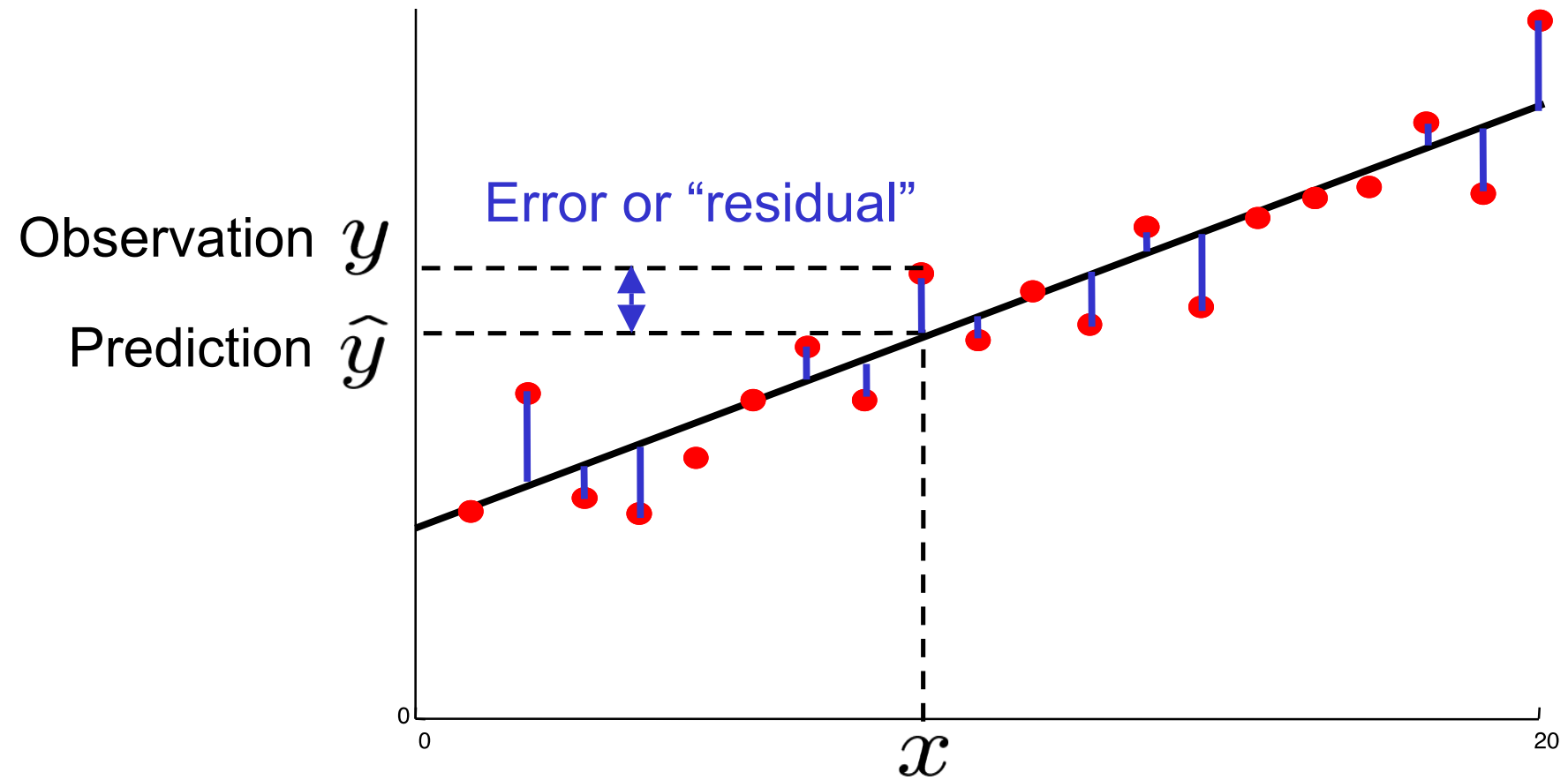
$$h = a + b x$$

Hypothesis: $h = a + b x$

$$h_{\Theta} = \Theta_0 x_0 + \Theta_1 x_1$$

$$h_{\Theta} = \Theta^T x$$

Idea: Chose θ so that $h_{\theta}(x)$ is close to h (muffins)
for training examples (Seniors, Muffins)



Naked Cup-Cake Run

Hypothesis:

$$h_{\Theta}(x) = \Theta^T x$$

$$h_{\Theta}(x) = \Theta_0 x_0 + \Theta_1 x_1$$

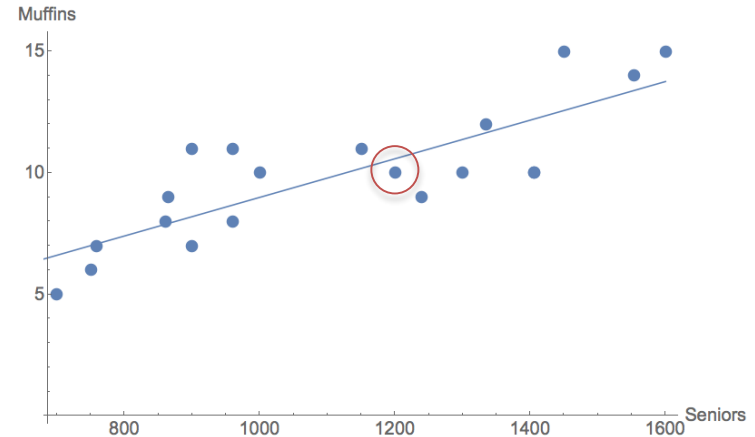
Notation:

$$x^i = \begin{pmatrix} 1 \\ x_1^i \\ \vdots \\ x_n^i \end{pmatrix}$$

Example

$$x^{20} = \begin{pmatrix} 1 \\ 1200 \end{pmatrix}$$

$$y^{20} = 9$$



Minimize: Squared Error (there are others)

$$\min_q J(q) = \frac{1}{m} \sum_{i=1}^m \left(h_q(x^{(i)}) - y^{(i)} \right)^2$$

Annotations for the equation above:

- Cost-Function** points to $J(q)$.
- #Data Points** points to m .
- Prediction** points to $h_q(x^{(i)})$.
- Observed** points to $y^{(i)}$.

Minimize: Squared Error (there are others)

$$\min_q J(q) = \frac{1}{m} \sum_{i=1}^m \left(h_q(x^{(i)}) - y^{(i)} \right)^2$$

#Data Points (points to m)
Prediction (points to $h_q(x^{(i)})$)
Observed (points to $y^{(i)}$)
Cost-Function (points to $J(q)$)

Alternative Versions of Squared Error Cost-Functions

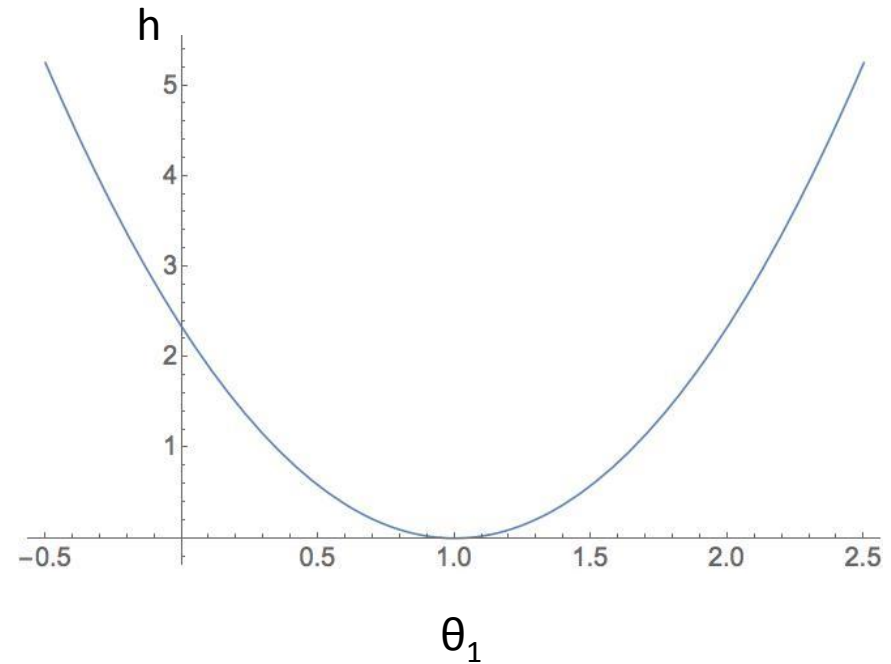
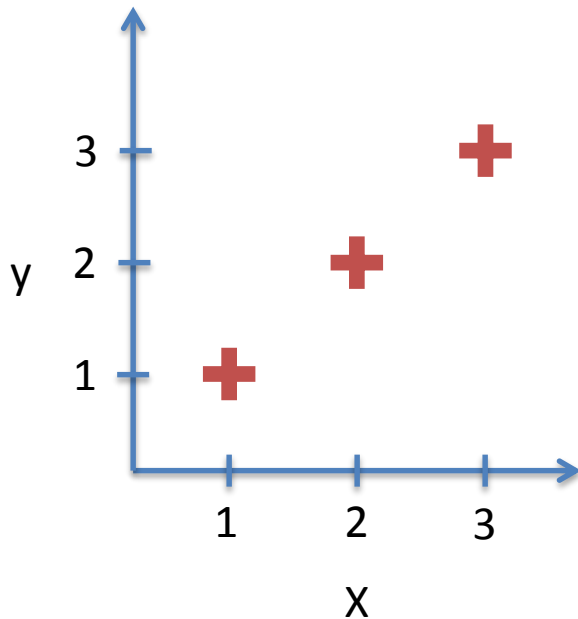
$$J(q) = \frac{1}{2m} \sum_{i=1}^m \left(h_q(x^{(i)}) - y^{(i)} \right)^2$$

$$J(q) = \frac{1}{2} \sum_{i=1}^m \left(h_q(x^{(i)}) - y^{(i)} \right)^2$$

1/2 is just for mathematical convenience

The Shape of the Cost Function

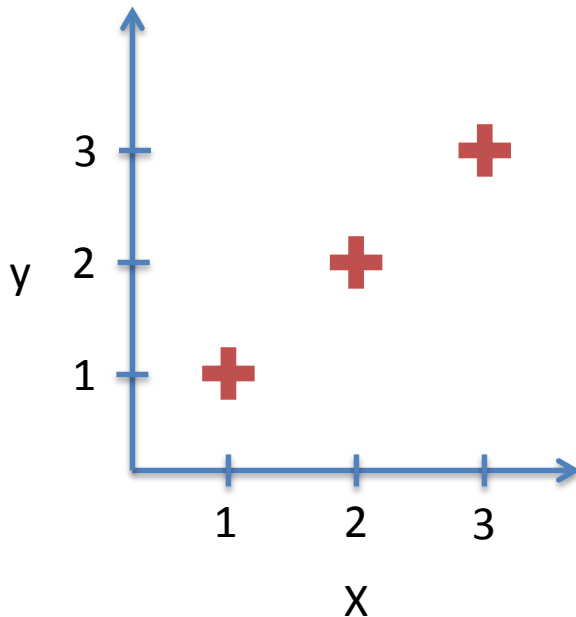
Lets assume $\theta_0=0$: $h_{q_1} = q_1 x_1$



$$J(q) = \frac{1}{2m} \sum_{i=1}^m \left(h_q(x^{(i)}) - y^{(i)} \right)^2$$

Clicker Question

Suppose our hypothesis is $h_{q_1} = q_1 x_1$
and we have the following training
set ($m=3$):



Our cost function is:

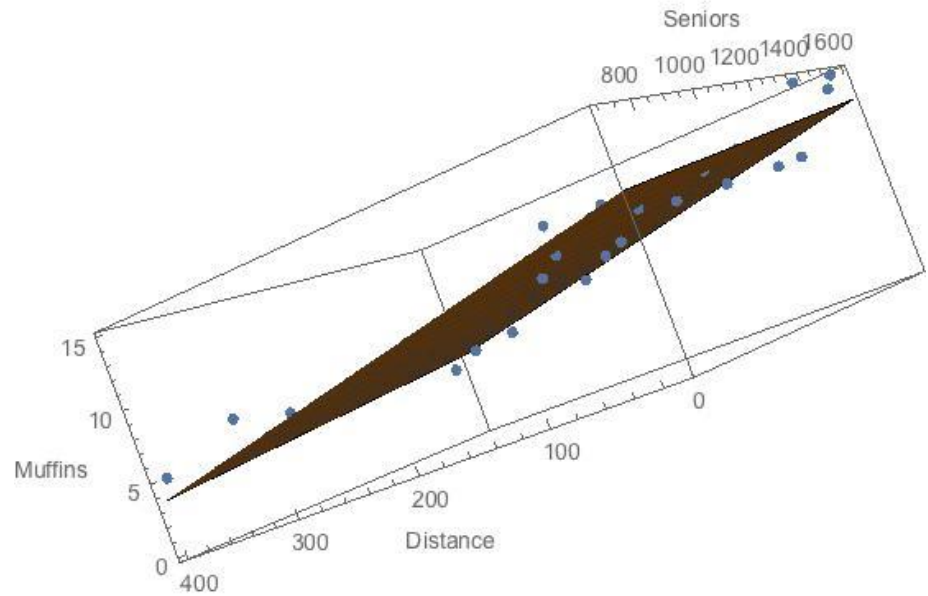
$$J(q_1) = \frac{1}{2m} \sum_{i=1}^m (h_{q_1}(x^{(i)}) - y^{(i)})^2$$

What is the cost of $\theta_1=0$

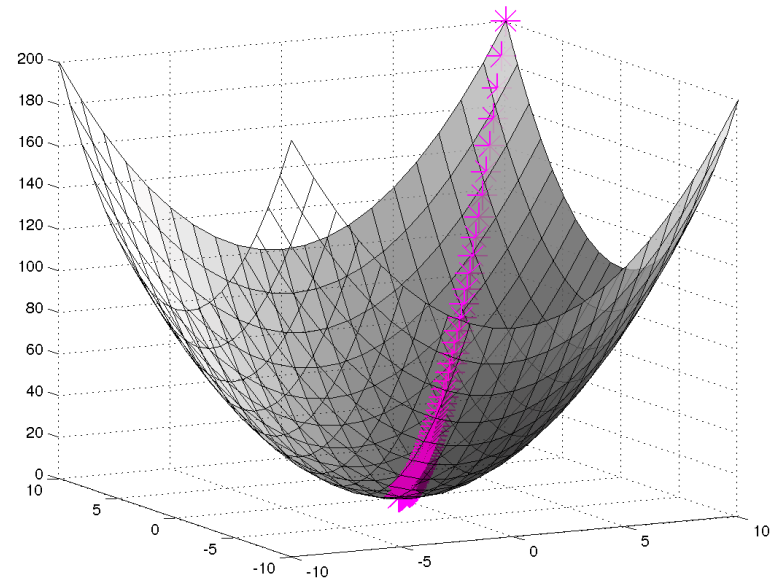
- a) 0
- b) 1/6
- c) 1
- d) 14/6

Works also in Higher Dimensions

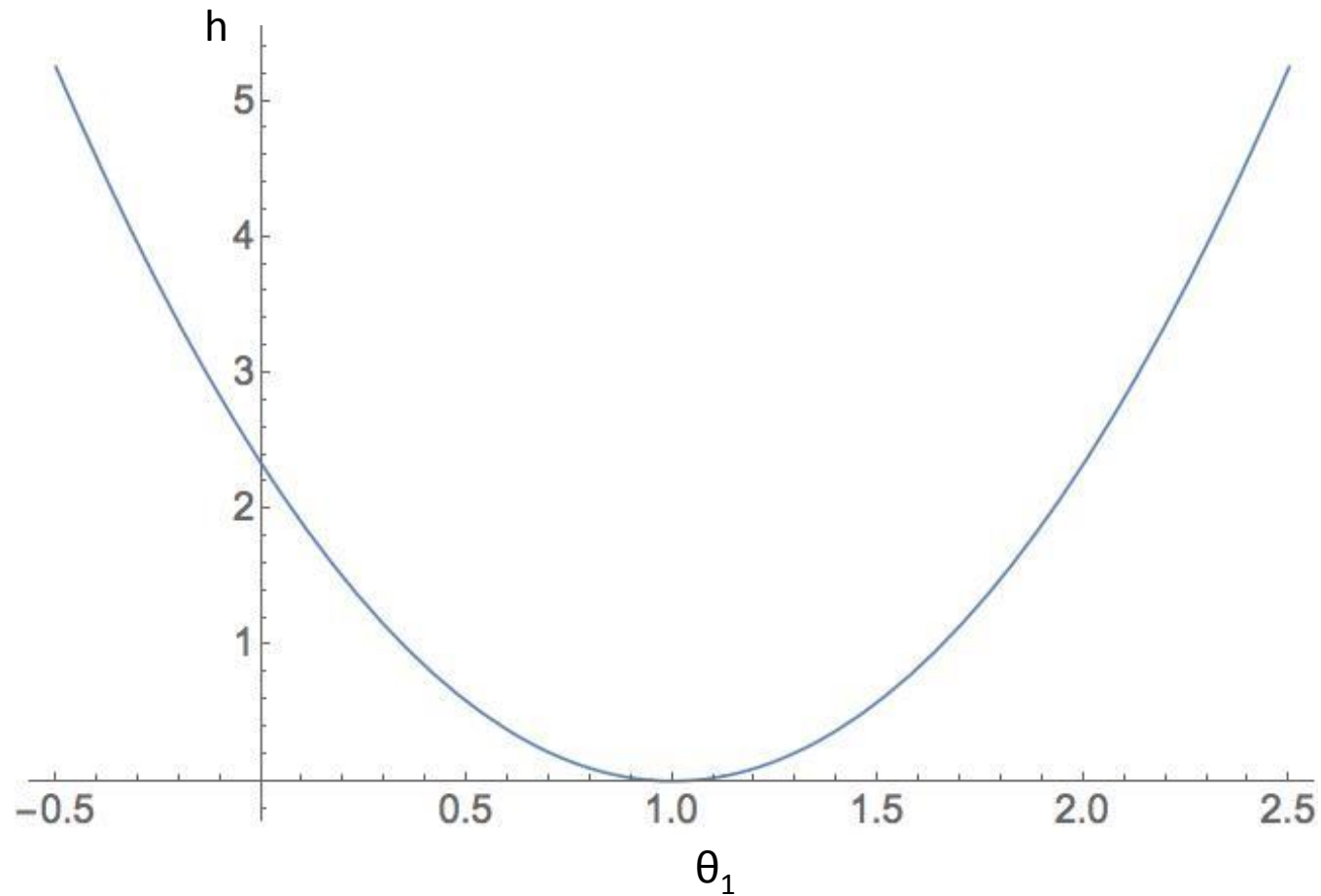
$$h_q(x) = q^T x$$



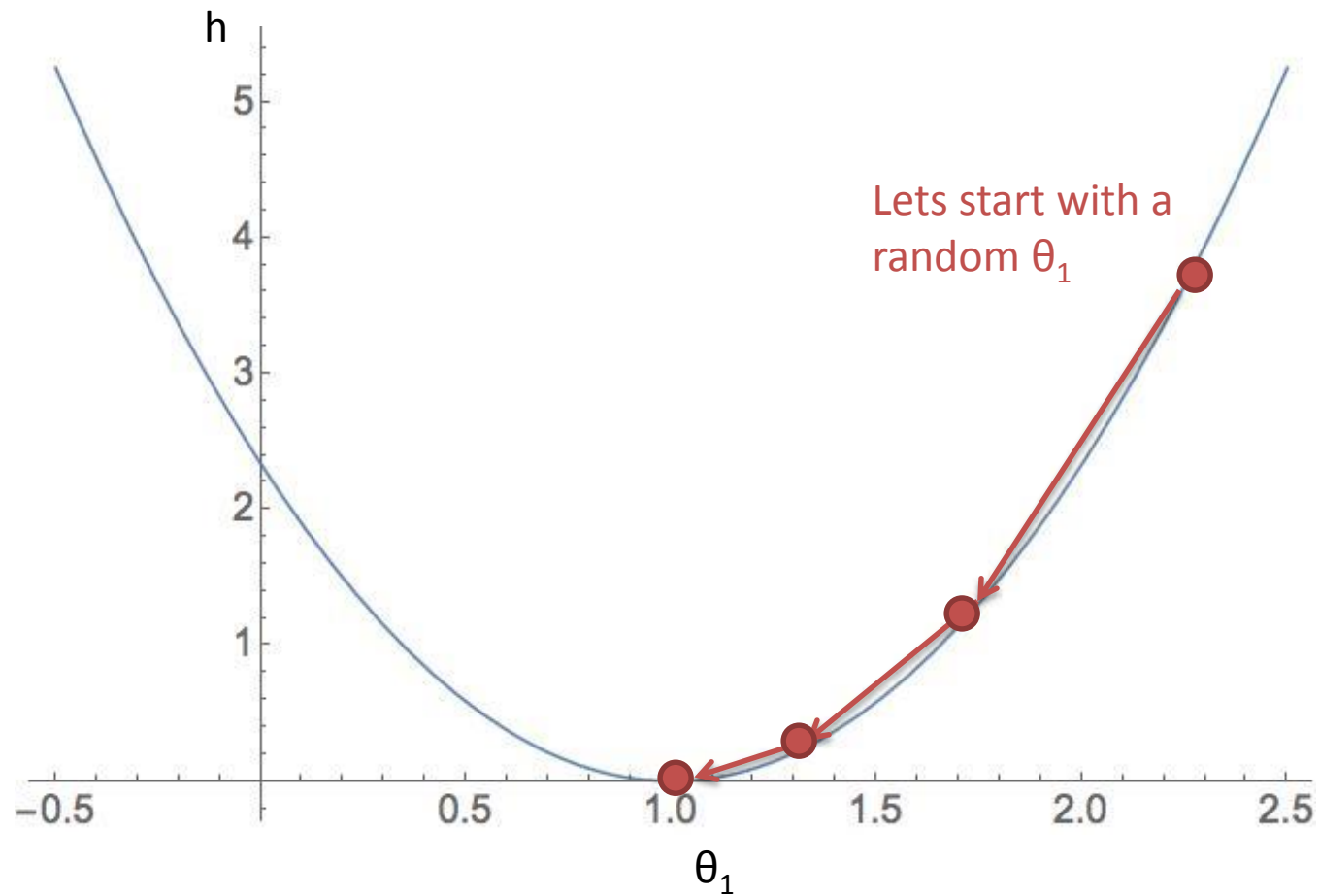
$$J(q_1) = \frac{1}{2m} \sum_{i=1}^m (h_q(x^{(i)}) - y^{(i)})^2$$



How Can We Find The Minimum?



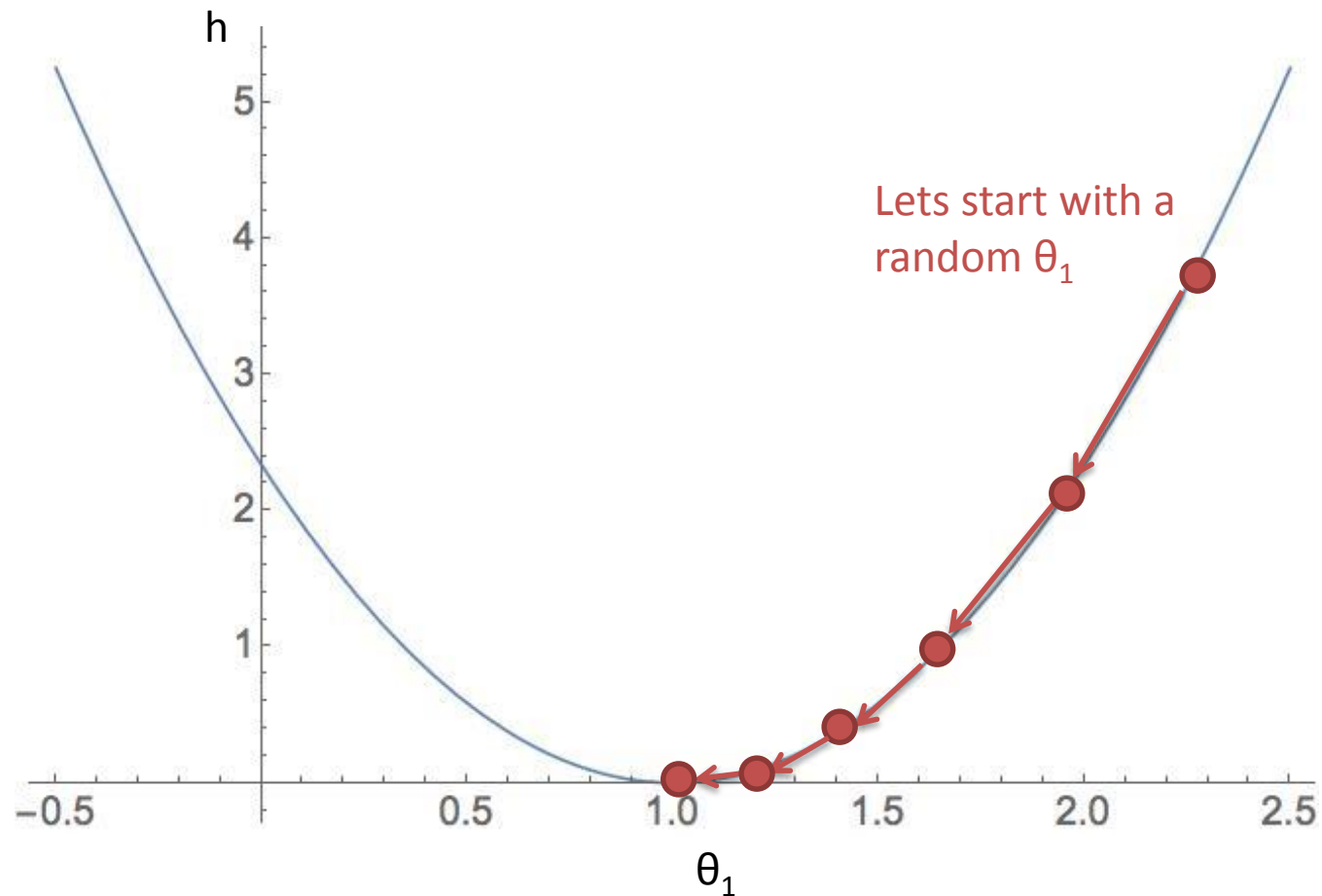
How Can We Find The Minimum?



$$q_1 = q_1 - a \frac{\partial J(q_1)}{\partial q_1}$$

Step Size / Learning Rate

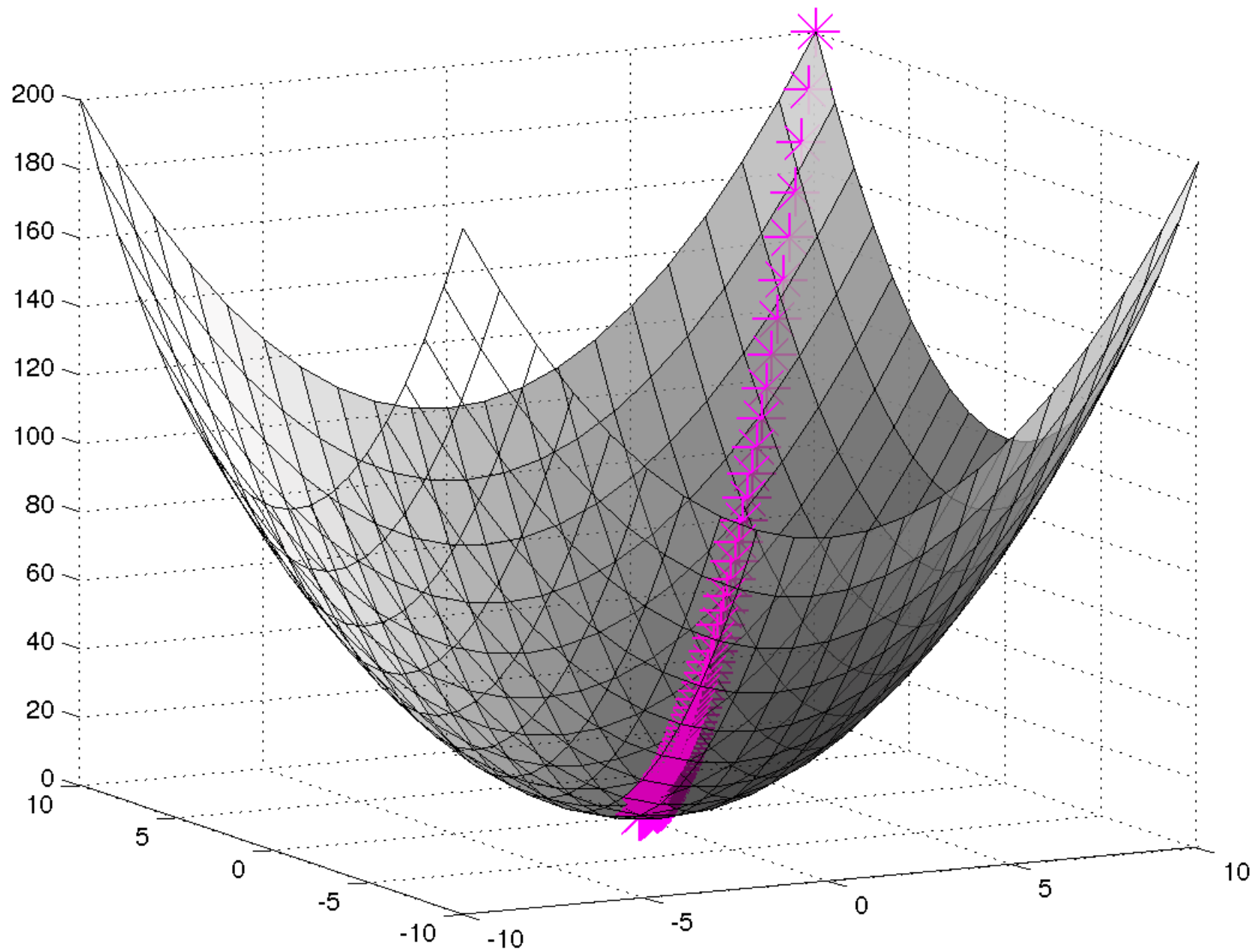
How Can We Find The Minimum?



Steps are automatically smaller the closer they get to the minimum

$$q_1 = q_1 - a \frac{\partial}{\partial q_1} J(q_1)$$

Step Size / Learning Rate



Gradient Descent

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



**Step Size /
Learning Rate**

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$h(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x$$

Gradient Descent

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

Lets assume single data point (m=1)

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_{\theta}(x) - y)^2$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$h(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x$$

Chain Rule

$$\frac{d}{dx} [f(g(x))] = f'(g(x))g'(x)$$

Gradient Descent

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

Lets assume single data point (m=1)

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_{\theta}(x) - y)^2 \\ &= 2 \cdot \frac{1}{2} (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_{\theta}(x) - y) \end{aligned}$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$h(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x$$

Chain Rule

$$\frac{d}{dx} [f(g(x))] = f'(g(x))g'(x)$$

Gradient Descent

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

Lets assume single data point (m=1)

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_{\theta}(x) - y)^2 \\ &= 2 \cdot \frac{1}{2} (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_{\theta}(x) - y) \\ &= (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left(\sum_{i=0}^n \theta_i x_i - y \right) \end{aligned}$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$h(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x$$

Chain Rule

$$\frac{d}{dx} [f(g(x))] = f'(g(x))g'(x)$$

Gradient Descent

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

Lets assume single data point (m=1)

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_{\theta}(x) - y)^2 \\ &= 2 \cdot \frac{1}{2} (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_{\theta}(x) - y) \\ &= (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left(\sum_{i=0}^n \theta_i x_i - y \right) \\ &= (h_{\theta}(x) - y) x_j \end{aligned}$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$h(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x$$

Chain Rule

$$\frac{d}{dx} [f(g(x))] = f'(g(x))g'(x)$$

LMS Update Rule

(Widrow-Hoff Learning)

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

Batch Gradient Descent

Repeat until convergence {

$$\theta_j := \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)} \quad (\text{for every } j)$$

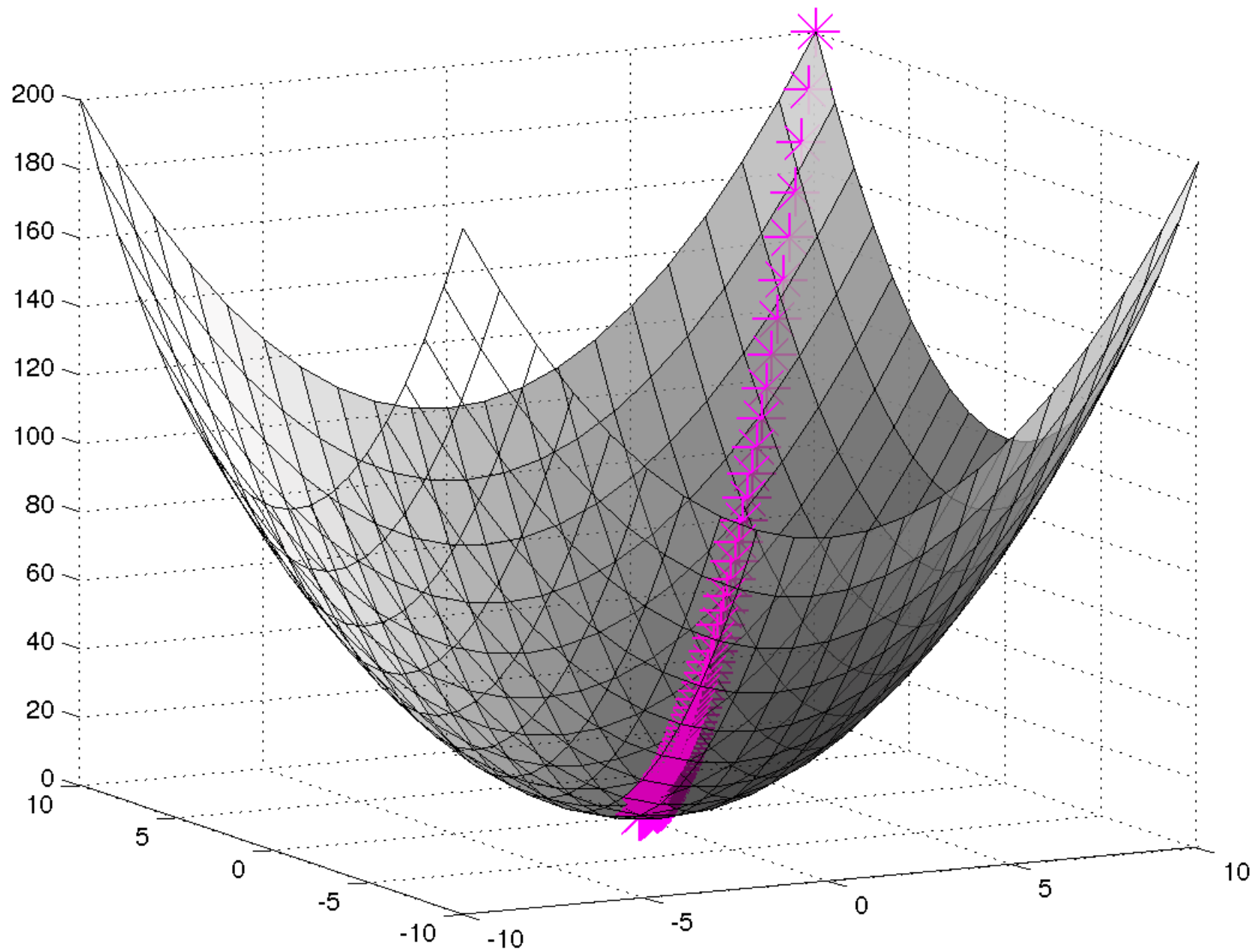
}

Clicker:

Can we get “stuck” in a local minima?

a) Yes

b) No



Clicker Question

Batch Gradient Descent

Repeat until convergence {

Data (x,y)

(1,1)

(2,2)

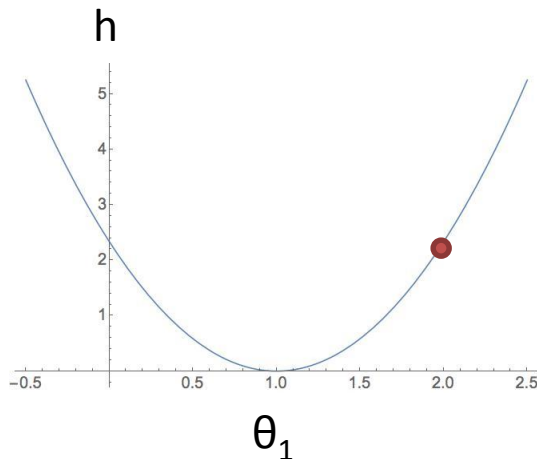
(3,3)

$$\theta_j := \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)} \quad (\text{for every } j)$$

}

Our hypothesis $h_{\theta}(x) = q_0 + q_1 x_1$

Current values: $\alpha = 0.01$ $q_0 = 0$ $q_1 = 2$



What is the next θ_1

a) 1

b) 1.54

c) 1.86

Batch Gradient Descent

Repeat until convergence {

Data (x,y)

(1,1)

(2,2)

(3,3)

$$\theta_j := \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)} \quad (\text{for every } j)$$

}

$$q_1 = q_1 + a \sum_{i=1}^m (y^{(i)} - h_q(x^{(i)})) x_j^{(i)}$$

Our hypothesis

$$h_q(x) = q_0 + q_1 x_1$$

$$q_1 = q_1 + a \left[\begin{aligned} & \left(y^{(1)} - (q_0 + q_1 x^{(1)}) \right) x^{(1)} + \\ & \left(y^{(2)} - (q_0 + q_1 x^{(2)}) \right) x^{(2)} + \\ & \left(y^{(3)} - (q_0 + q_1 x^{(3)}) \right) x^{(3)} \end{aligned} \right]$$

Current values:

$$a = 0.01$$

$$q_0 = 0$$

$$q_1 = 2$$

$$q_1 = 2 + 0.01 \left[\begin{aligned} & \left((1 - 2)1 + \right) \\ & \left((2 - 4)2 + \right) \\ & \left((3 - 6)3 \right) \end{aligned} \right]$$

$$q_1 = 2 + 0.01(-1 - 4 - 9) = 2 - 0.14 = 1.86$$

How To Implement It

Repeat Until Convergence {

$$Temp_0 = q_0 + a \sum_{i=0}^m (y^{(i)} - h_q(x^{(i)})) x_0^{(i)}$$

$$Temp_1 = q_1 + a \sum_{i=1}^m (y^{(i)} - h_q(x^{(i)})) x_1^{(i)}$$

\vdots \vdots

$$Temp_n = q_n + a \sum_{i=1}^m (y^{(i)} - h_q(x^{(i)})) x_n^{(i)}$$

$$\theta_0 = Temp_0$$

$$\theta_1 = Temp_1$$

\vdots

$$\theta_n = Temp_n$$

}

What is the problem with Batch Gradient Descent for Big Data?

Repeat until convergence {

$$\theta_j := \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)} \quad (\text{for every } j)$$

}

Stochastic Gradient Descent (Incremental Gradient Descent)

```
Loop {  
    for i= 1 to m, {  
         $\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$       (for every j ).  
    }  
}
```

Clicker: Does this algorithm always converge?

a) Yes

b) No

Compromise

- **Batch gradient descent:**
Use all m examples in each step
- **Stochastic gradient descent:**
Use 1 example in each step
- **Mini-batch gradient descent:**
Use b (e.g., 10) examples in each step
 - Easier to parallelize (e.g., SIMD)
 - Bit more robust against local minima (depending on b)

Important!

For linear regression there exists a closed form

- Good for smaller data sizes

Practical Considerations: Features

X_1 = Age (e.g., 24; 32; 100)

X_2 = Salary (e.g., \$120000; \$200000; \$210 000)

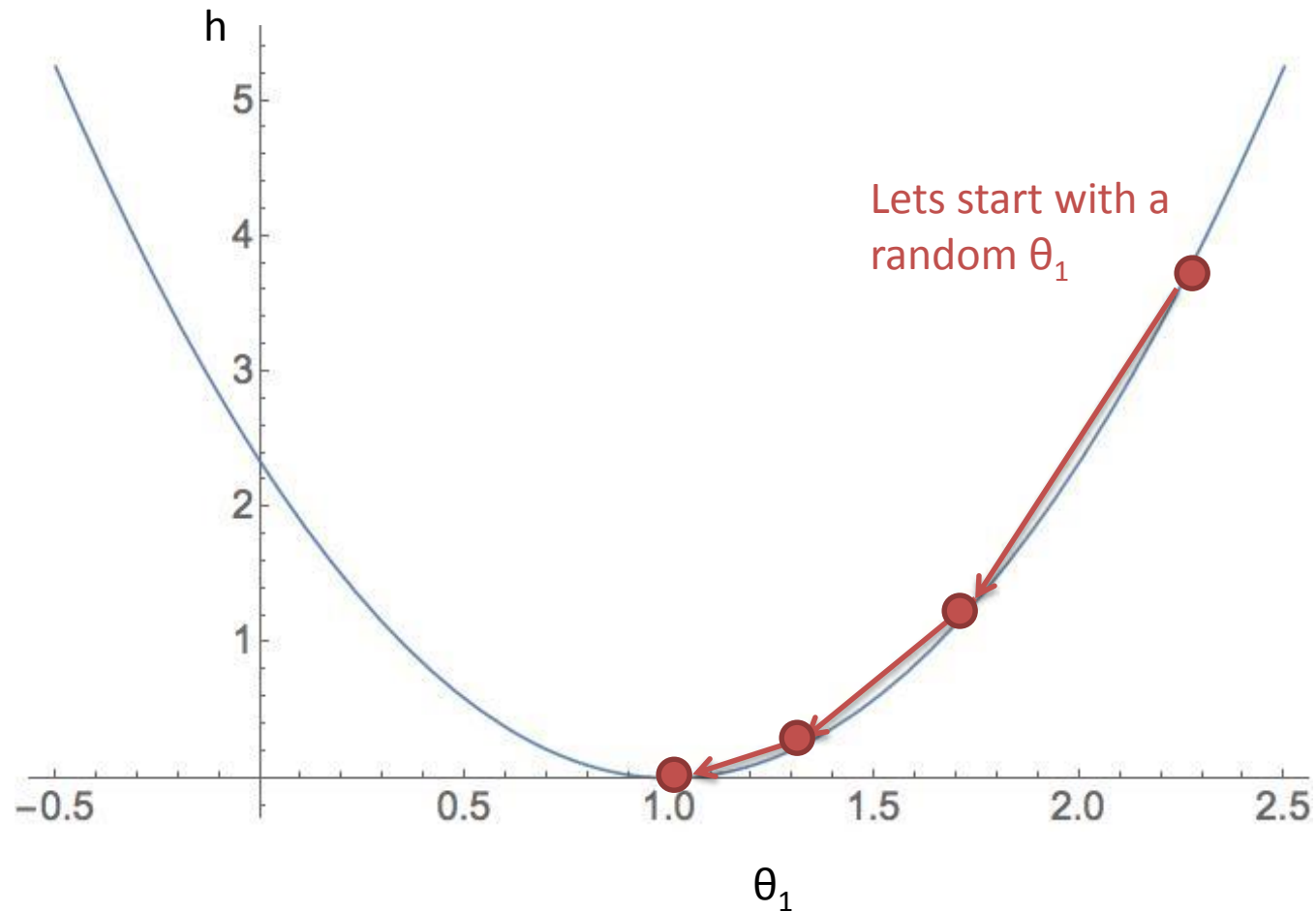
Solution: Feature Scaling $-1 \leq x \leq 1$

Mean Normalization $-0.5 \leq x \leq 0.5$

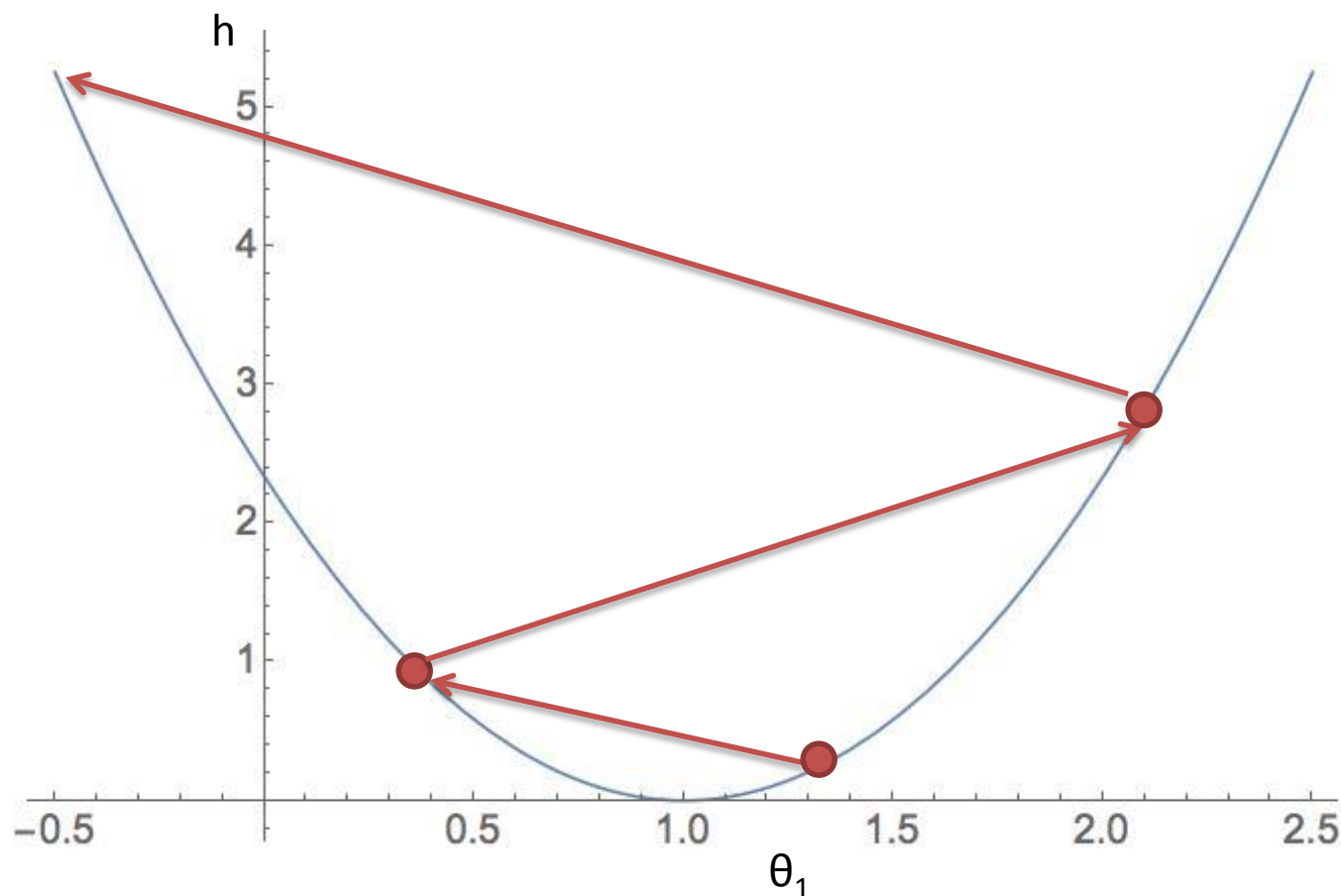
$$x_j = \frac{x_j - m_j}{\max(x_j) - \min(x_j)}$$

$$x_j^i = \frac{x_j^i - m_j}{S_j}$$

Practical Considerations: Learning Rate



Practical Considerations: Learning Rate



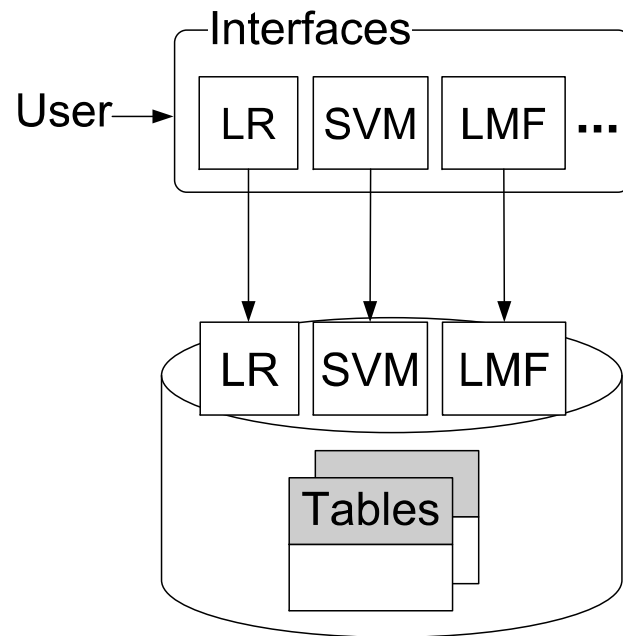
For sufficiently small α , $J(\theta)$ should decrease on every iteration

Bug if α is too small, it will take a long time to converge

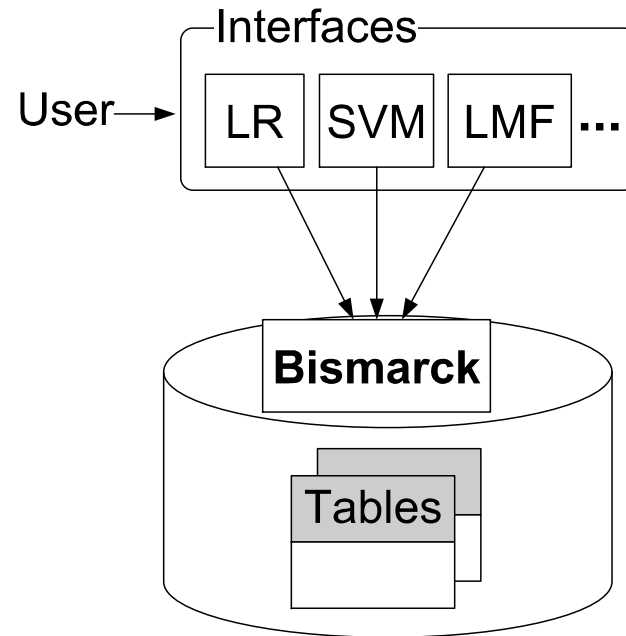
Try 2x or 3x steps. For example: 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1.0

Observe J and pick the one which has a good convergence rate

Towards a Unified Architecture for in-RDBMS Analytics



Current In-RDBMS Analytics

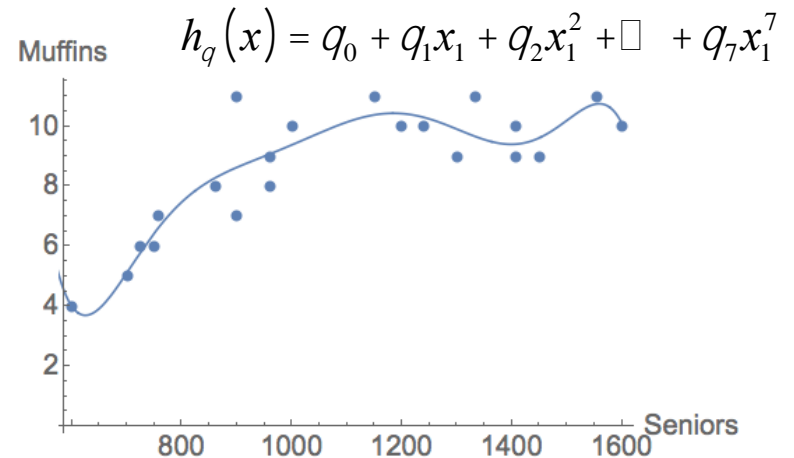
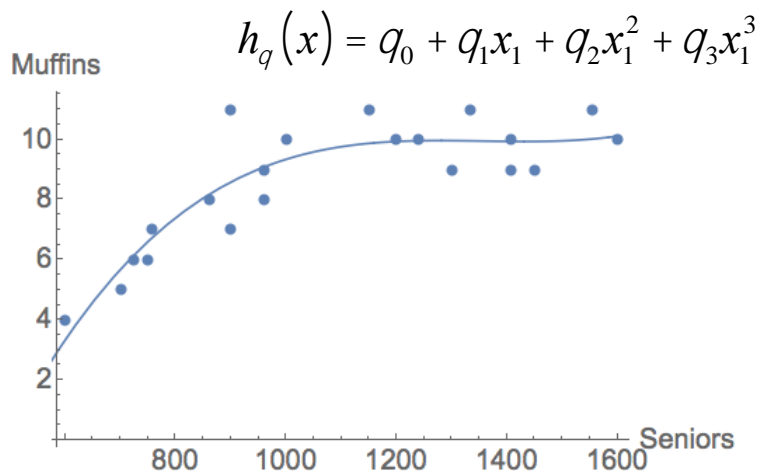
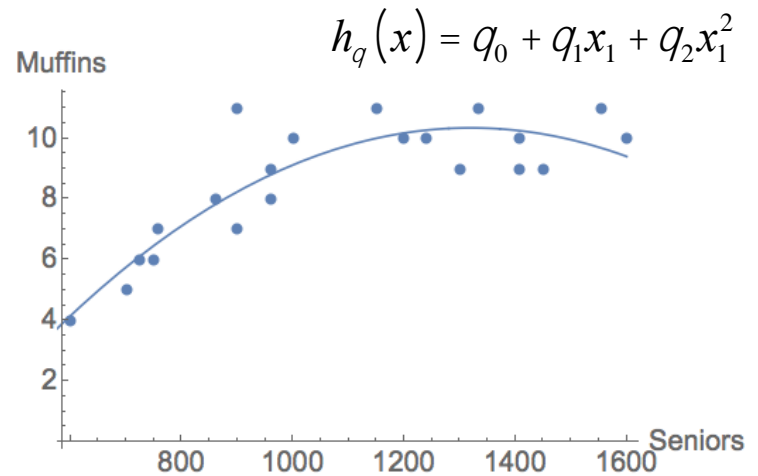
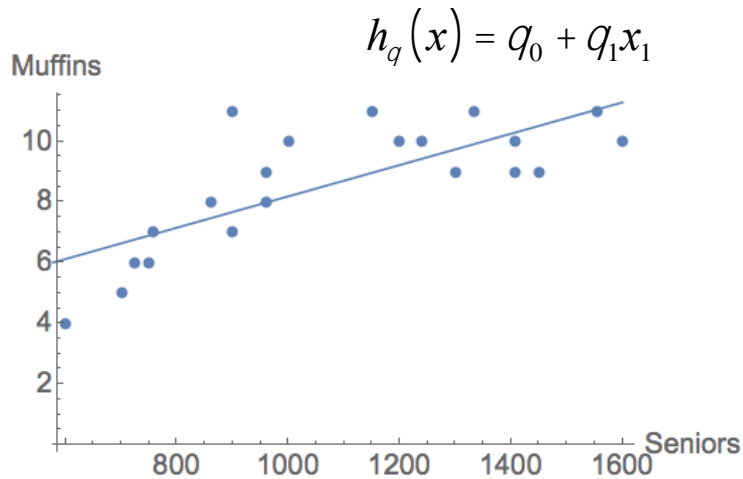


Bismarck In-RDBMS Analytics

Towards a Unified Architecture for in-RDBMS Analytics

Analytics Task	Objective
Logistic Regression (LR)	$\sum_i \log(1 + \exp(-y_i w^T x_i)) + \mu \ \vec{w}\ _1$
Classification (SVM)	$\sum_i (1 - y_i w^T x_i)_+ + \mu \ \vec{w}\ _1$
Recommendation (LMF)	$\sum_{(i,j) \in \Omega} (L_i^T R_j - M_{ij})^2 + \mu \ L, R\ _F^2$
Labeling (CRF) [48]	$\sum_k \left[\sum_j w_j F_j(y_k, x_k) - \log Z(x_k) \right]$
Kalman Filters	$\sum_{t=1}^T \ Cw_t - f(y_t)\ _2^2 + \ w_t - Aw_{t-1}\ _2^2$
Portfolio Optimization	$p^T w + w^T \Sigma w \quad \text{s.t.} \quad w \in \Delta$

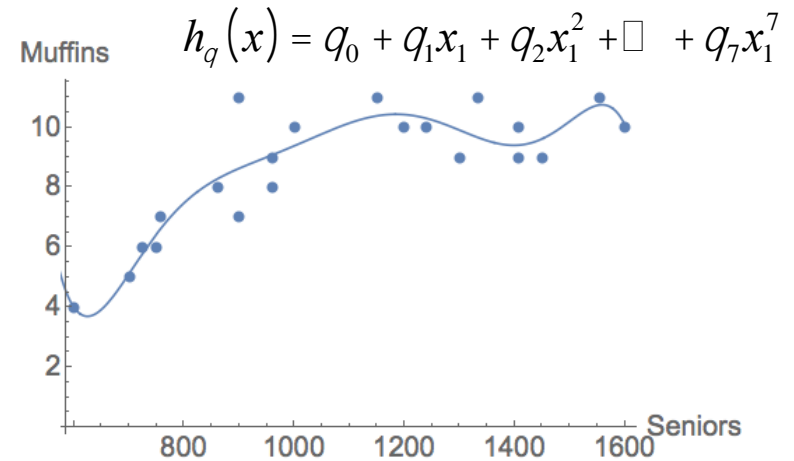
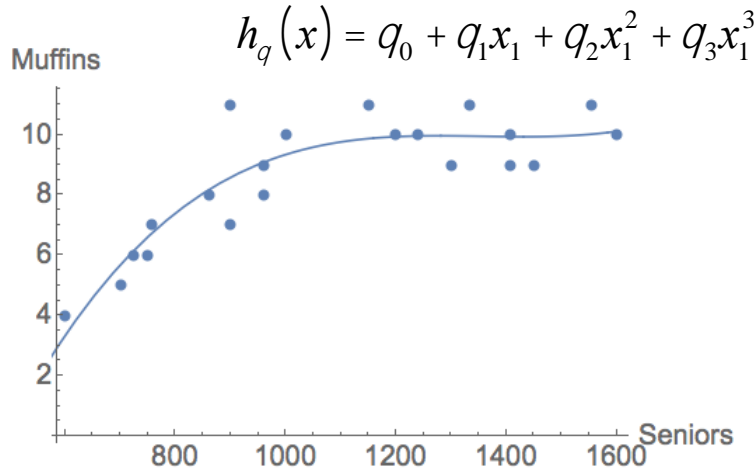
Polynomial Regression



How To Prevent Overfitting

- Adjust features
 - Reduce number of polynomials
 - Reduce number of features
- Regularization
 - Keep all the features, but reduce their impact
 - Works well when we have a lot of features, each of which contributes a little

Regularization: Intuition

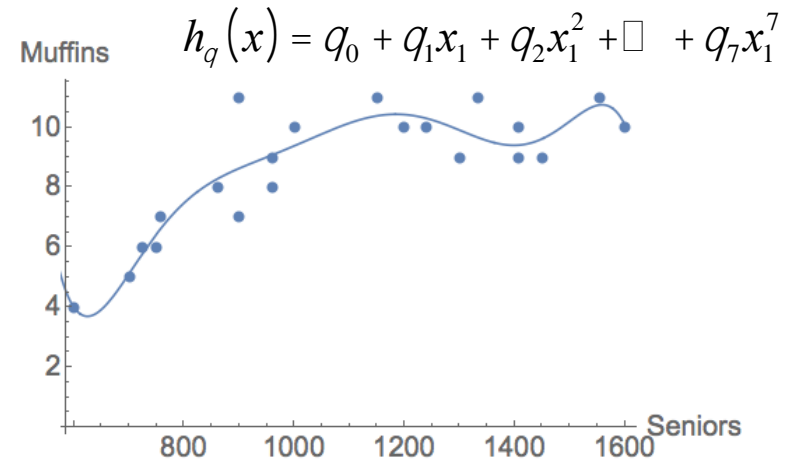
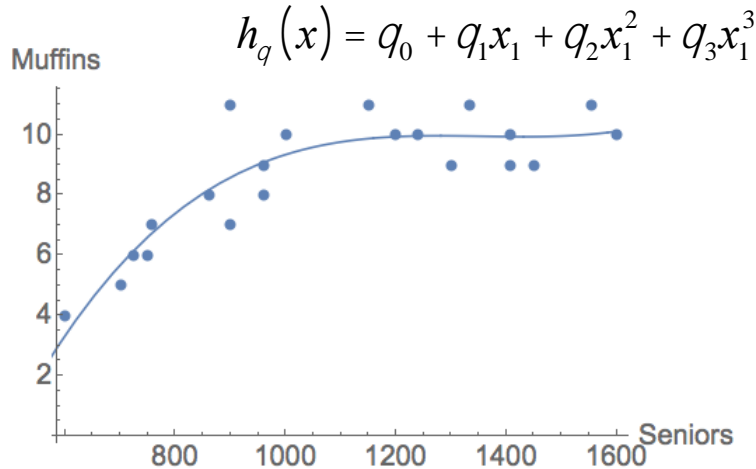


$$\min_q J(q) = \frac{1}{2m} \sum_{i=1}^m \left(h_q(x^{(i)}) - y^{(i)} \right)^2$$

$$+ 1000q_4 + 1000q_5 + 1000q_6 + 1000q_7$$

$$q_4, q_5, q_6, q_7 \rightarrow 0$$

Regularization: Intuition



$$\min_q J(q) = \frac{1}{2m} \sum_{i=1}^m \left(h_q(x^{(i)}) - y^{(i)} \right)^2 + \underbrace{\frac{\lambda}{2} \sum_{j=1}^n q_j^2}_{\text{Regularization term}}$$

Works ok for reducing the impact of polynomials (better to reduce nb of polynomials directly)
 Works great for a bag of equally important features.