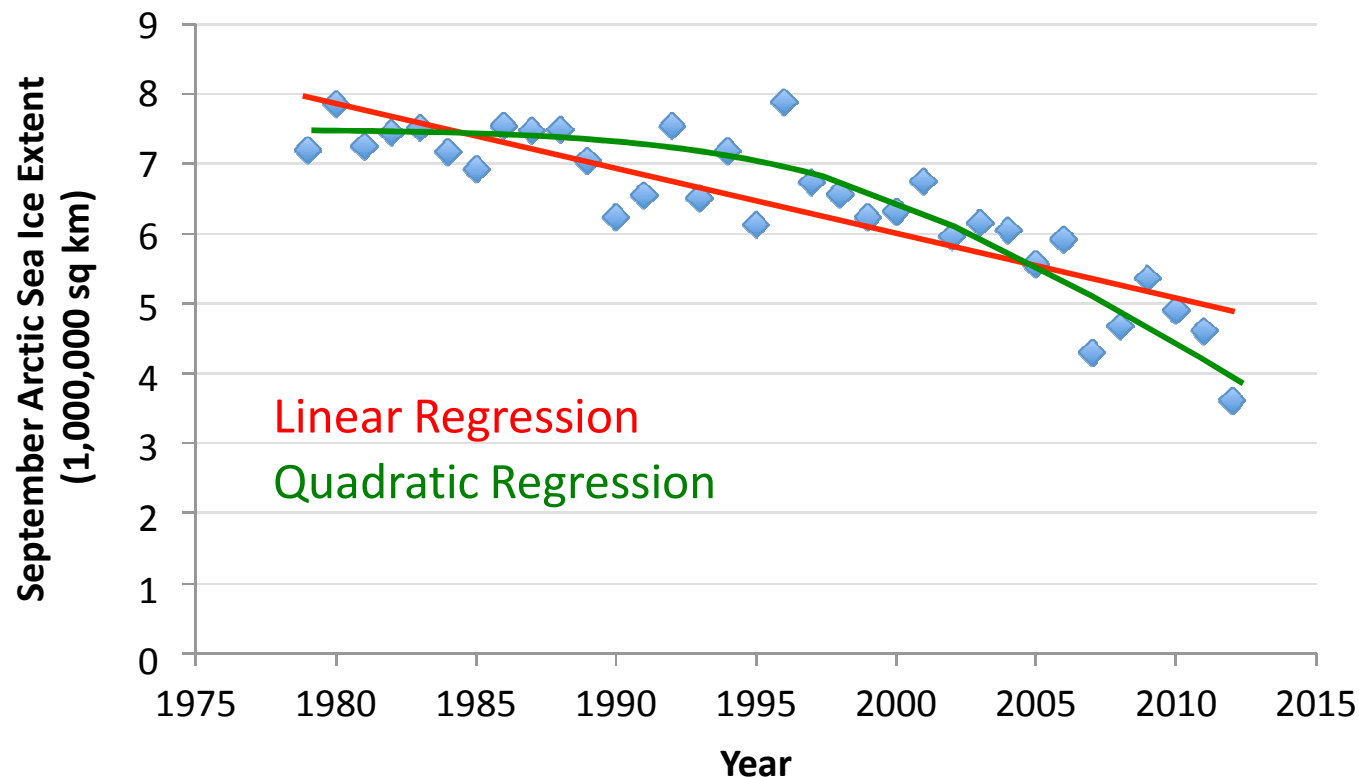# BBM406: Fundamentals of Machine Learning

Linear Regression, Cost Function, Gradient Descent

# Regression

Given:

- Data $X = \{x^{(1)}, x^{(2)}, \cdots, x^{(n)}\}$ , where $x^{(i)} \in R$

- Corresponding labels $y = \{y^{(1)}, y^{(2)}, \cdots, y^{(n)}\}$, where $x^{(i)} \in R$
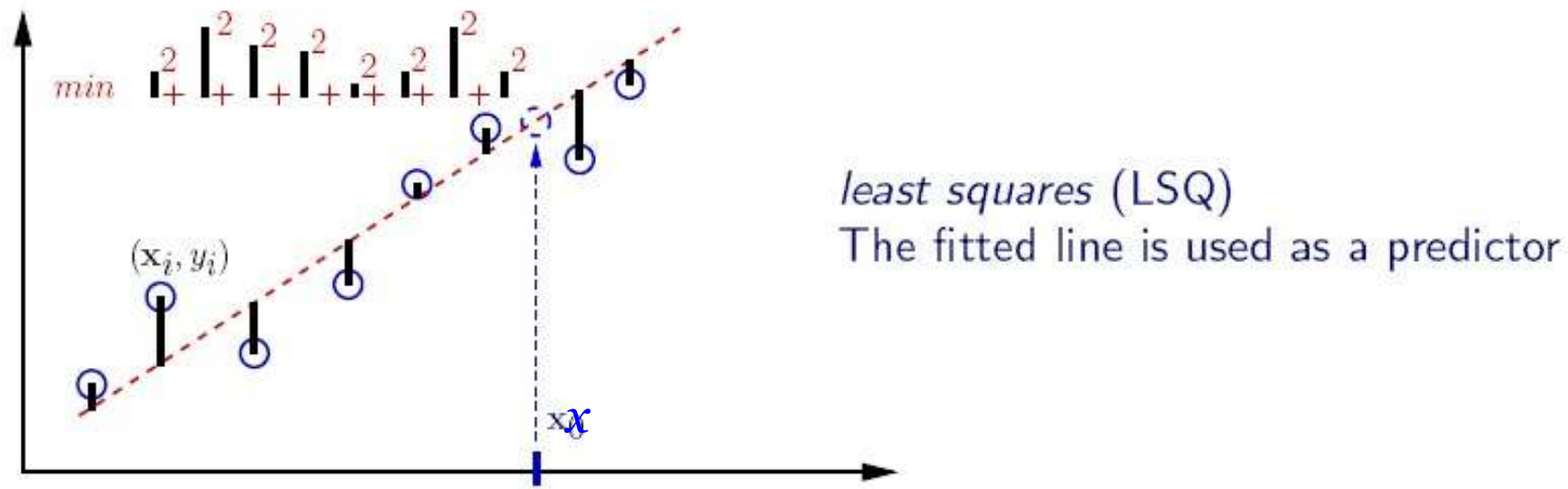
# Linear Regression

- Hypothesis:

Assume $x_0 = 1$

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_d x_d = \sum_{j=0}^{d} \theta_j x_j = h_\theta(X)$$
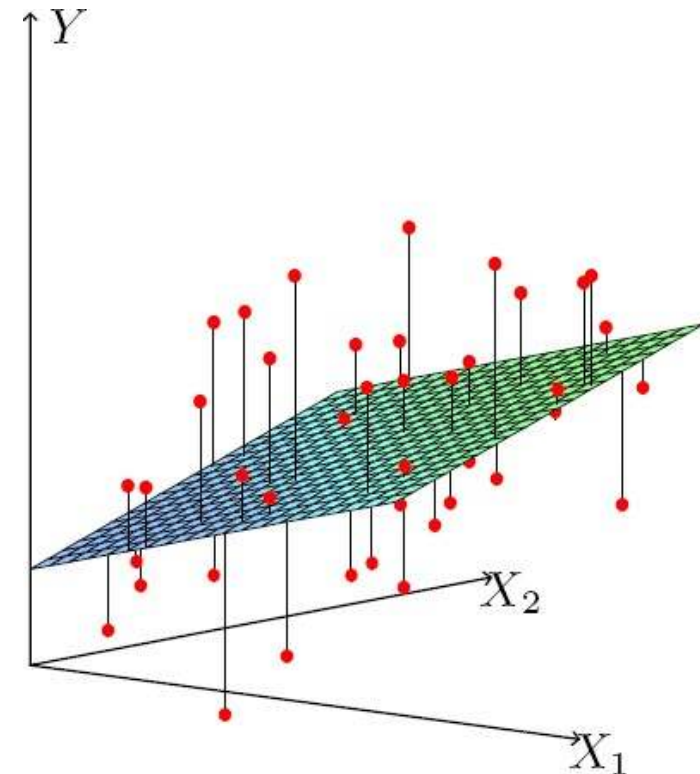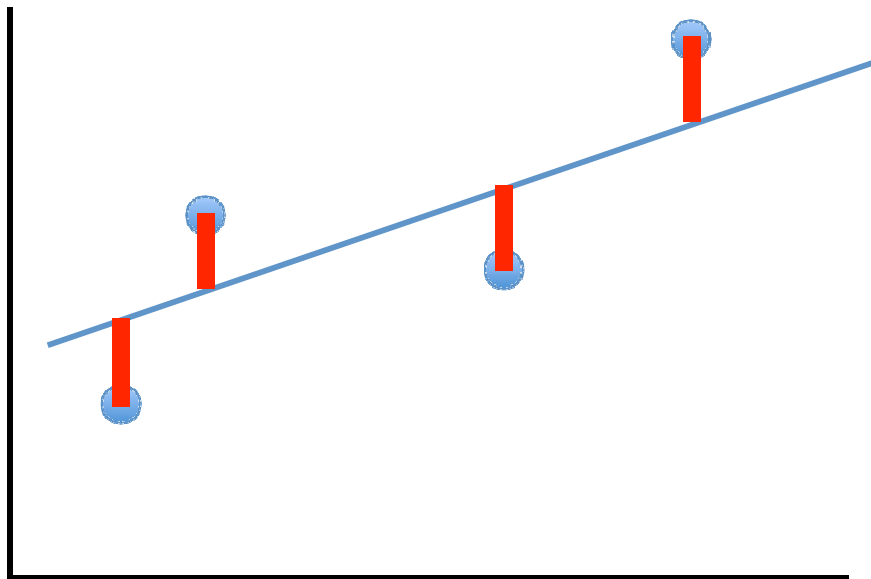
- Fit model by minimizing sum of squared errors



least squares (LSQ)
The fitted line is used as a predictor

# Least Squares Linear Regression

- Cost Function

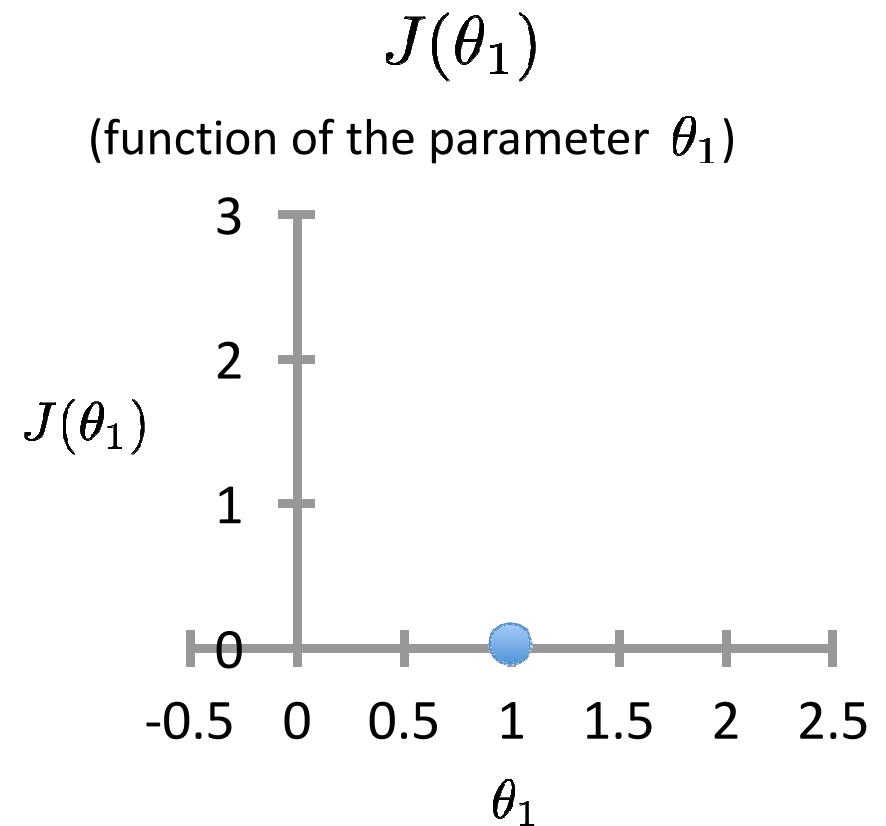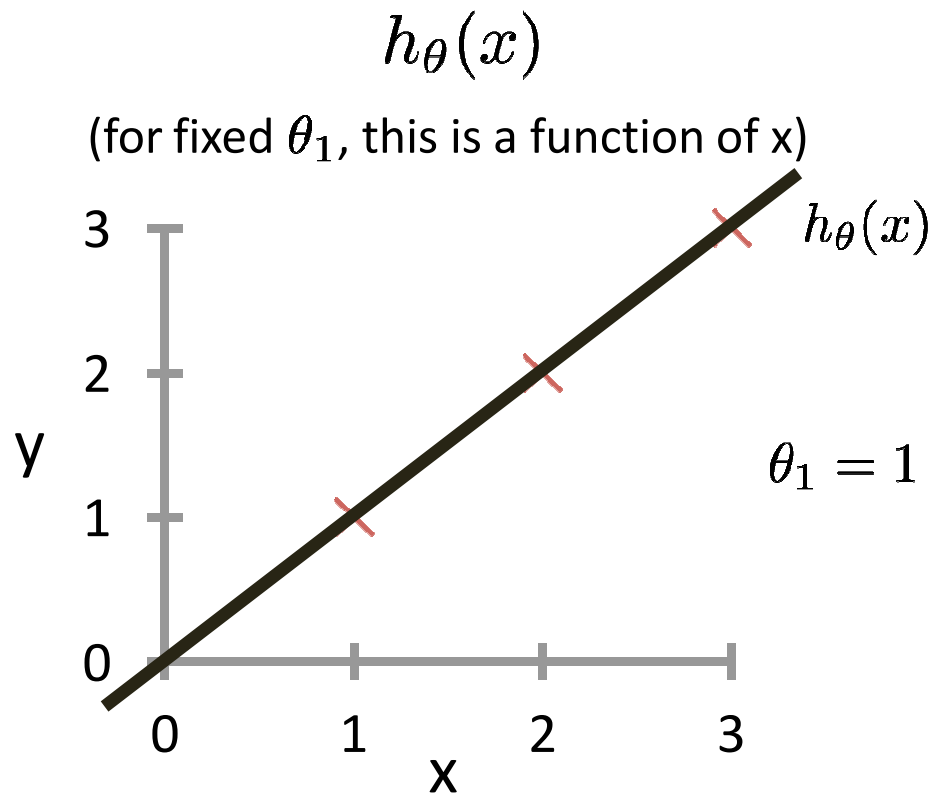$$J(\theta) = \frac{1}{2n} \sum_{i=1}^{n} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

- Fit by solving

# Intuition Behind Cost Function

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

For insight on $J()$, let's assume $x^{(i)} \in R$ and $\theta = [\theta_0, \theta_1]$

$h_\theta(x)$

(for fixed $\theta_1$, this is a function of x)



$h_\theta(x)$

$\theta_1 = 1$

$J(\theta_1)$

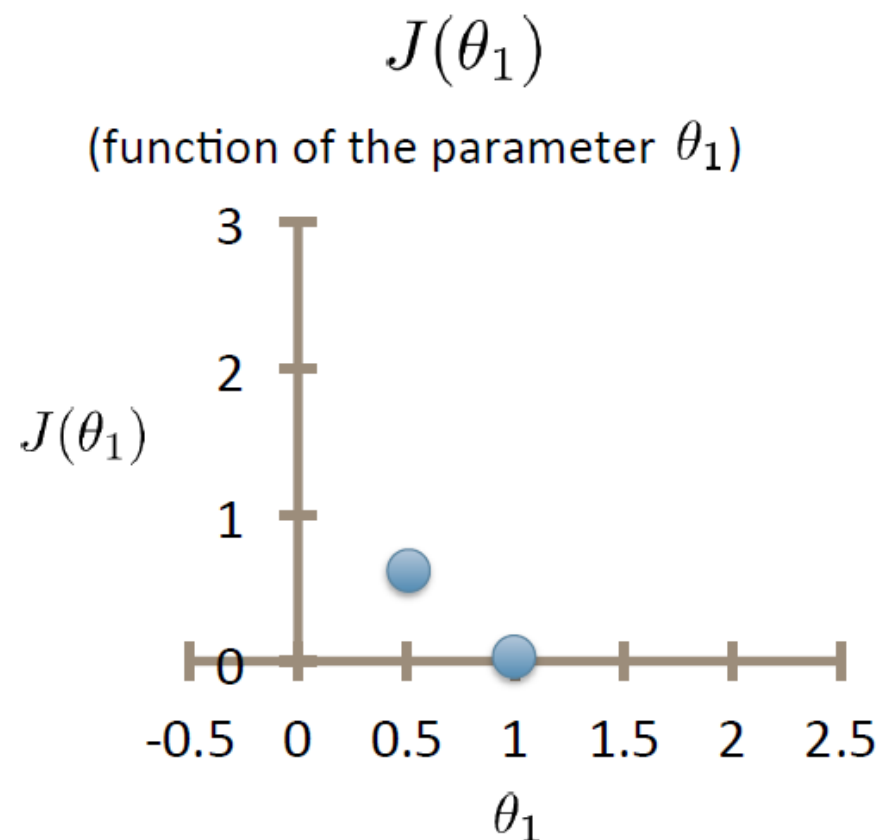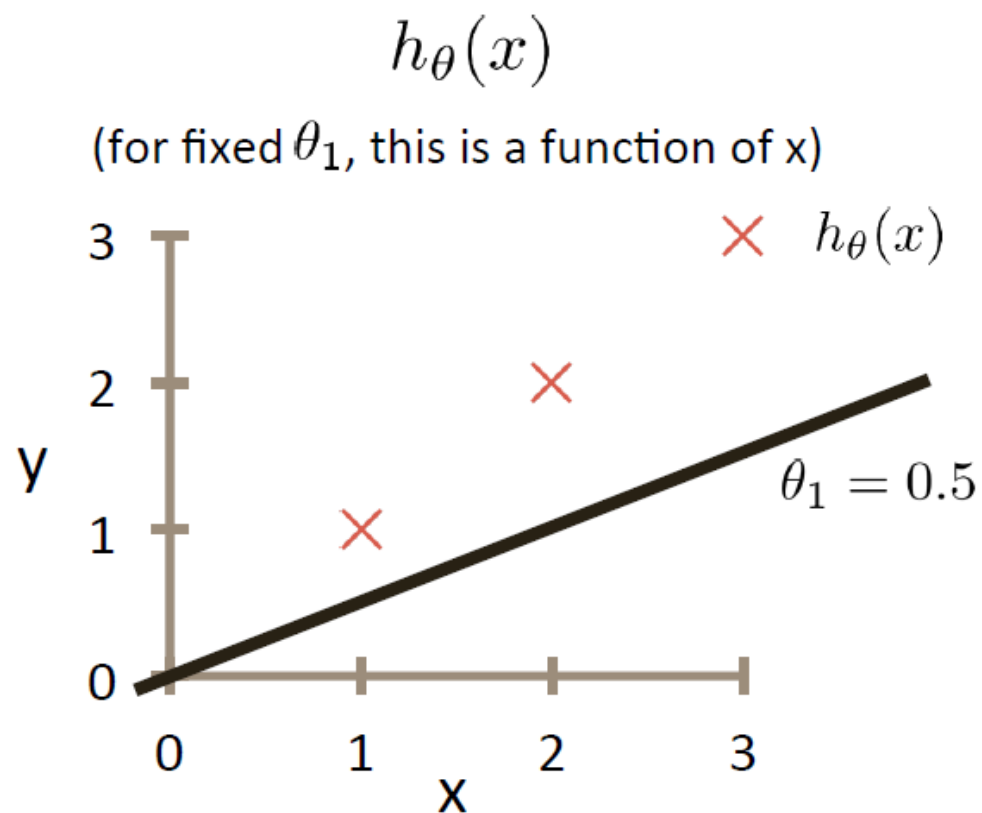(function of the parameter $\theta_1$)



$J(\theta_1)$

$$J([0,1]) = \frac{1}{2\times 3}\left[(1-1)^2+(2-2)^2+(3-3)^2\right] = 0$$

# Intuition Behind Cost Function

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

For insight on $J()$, let's assume $x^{(i)} \in R$ and $\theta = [\theta_0, \theta_1]$

$h_\theta(x)$

(for fixed $\theta_1$, this is a function of x)

$J(\theta_1)$

(function of the parameter $\theta_1$)
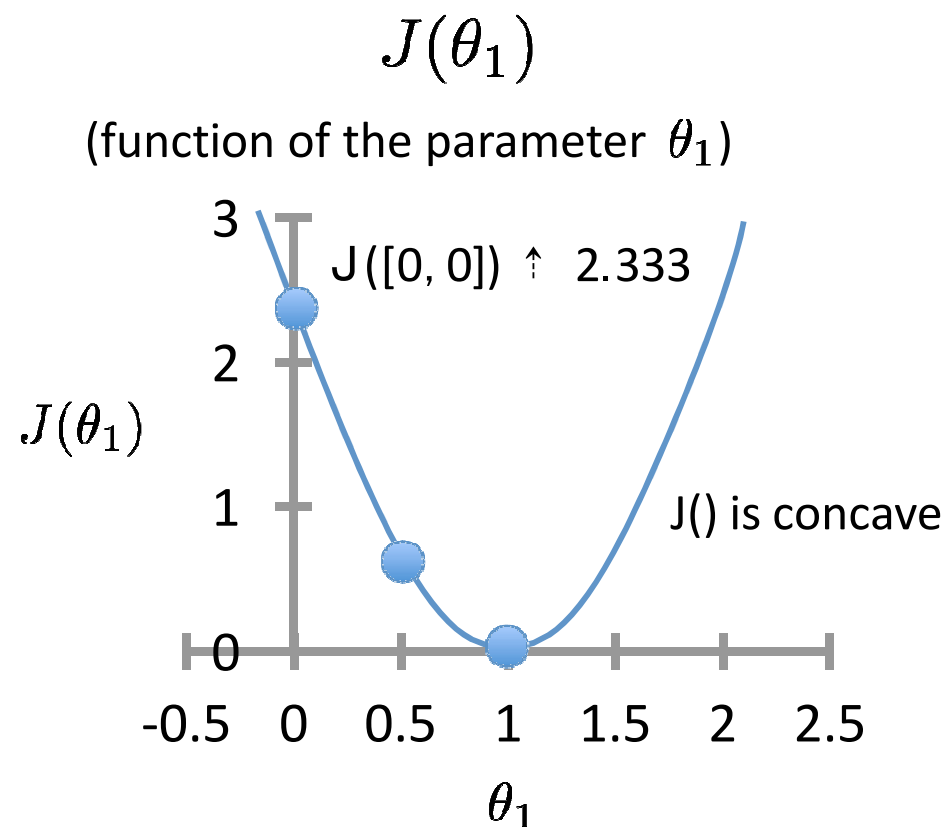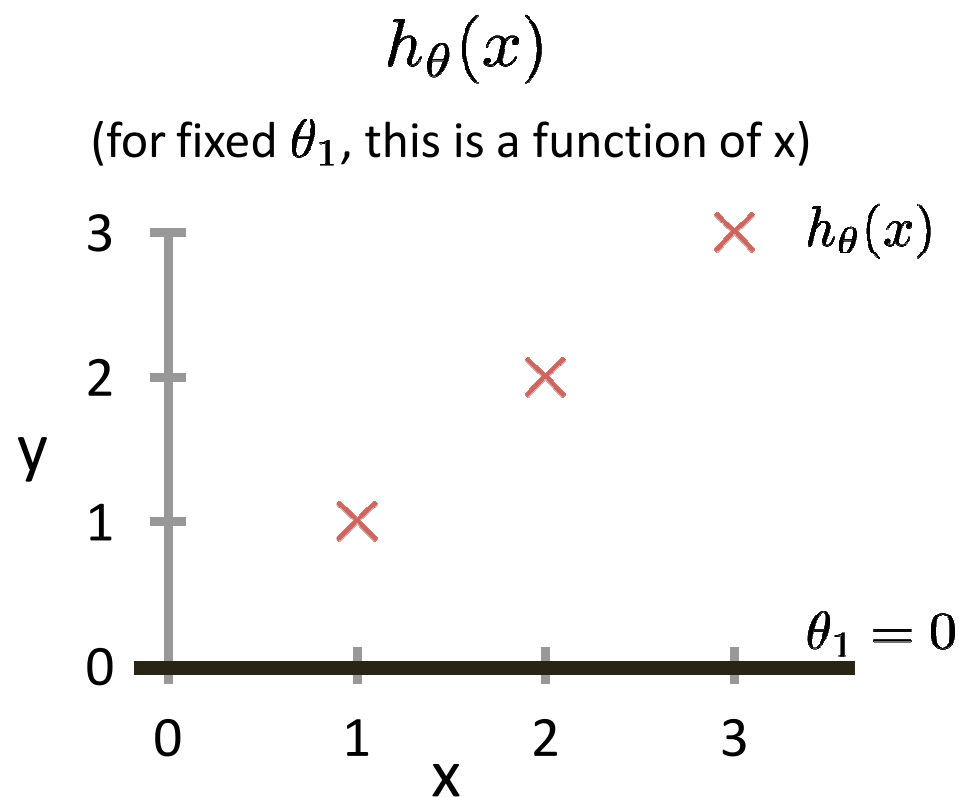


$\times \quad h_\theta(x)$

$\theta_1 = 0.5$

$$J([0,0.5]) = \frac{1}{2 \times 3} [(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2] \approx 0.58$$

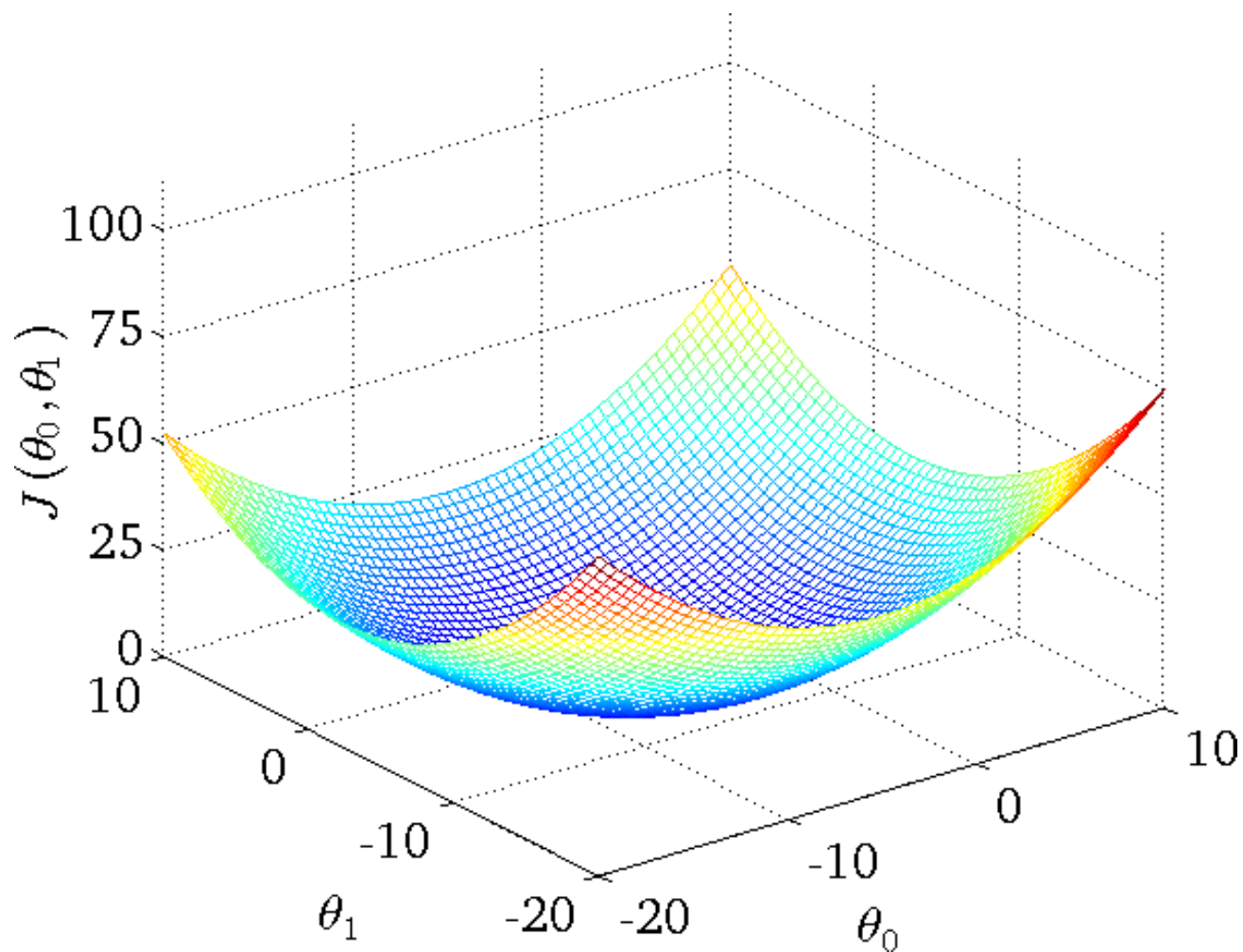# Intuition Behind Cost Function

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

For insight on $J()$, let's assume $x^{(i)} \in R$ and $\theta = [\theta_0, \theta_1]$

$h_\theta(x)$

(for fixed $\theta_1$, this is a function of x)

$J(\theta_1)$

(function of the parameter $\theta_1$)



$J([0,0]) \uparrow 2.333$

$J()$ is concave

$$J([0,0]) = \frac{1}{2 \times 3} \left[ (0-1)^2 + (0-2)^2 + (0-3)^2 \right] \approx 2.33$$
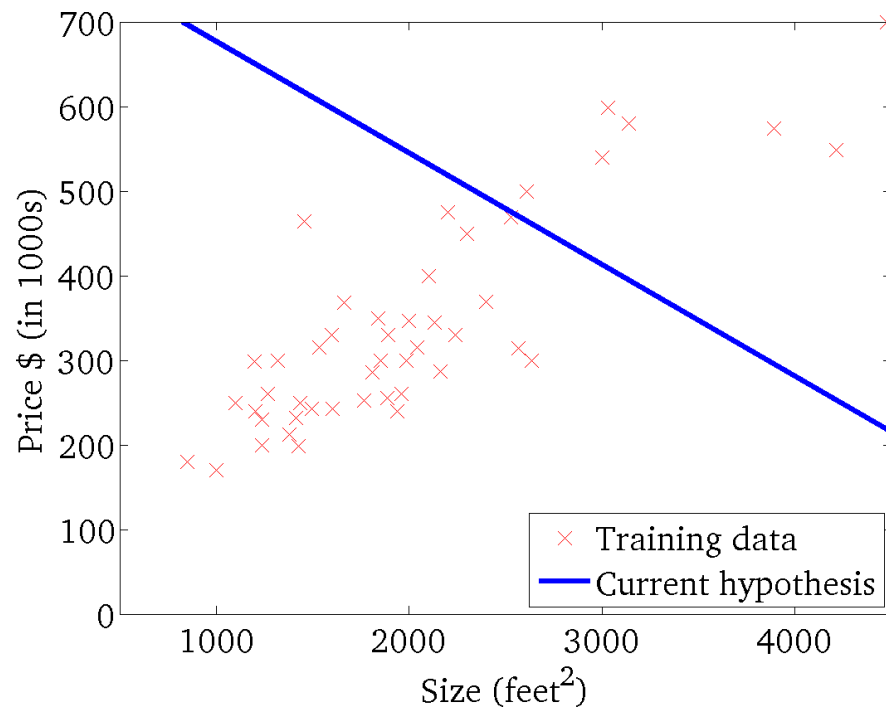
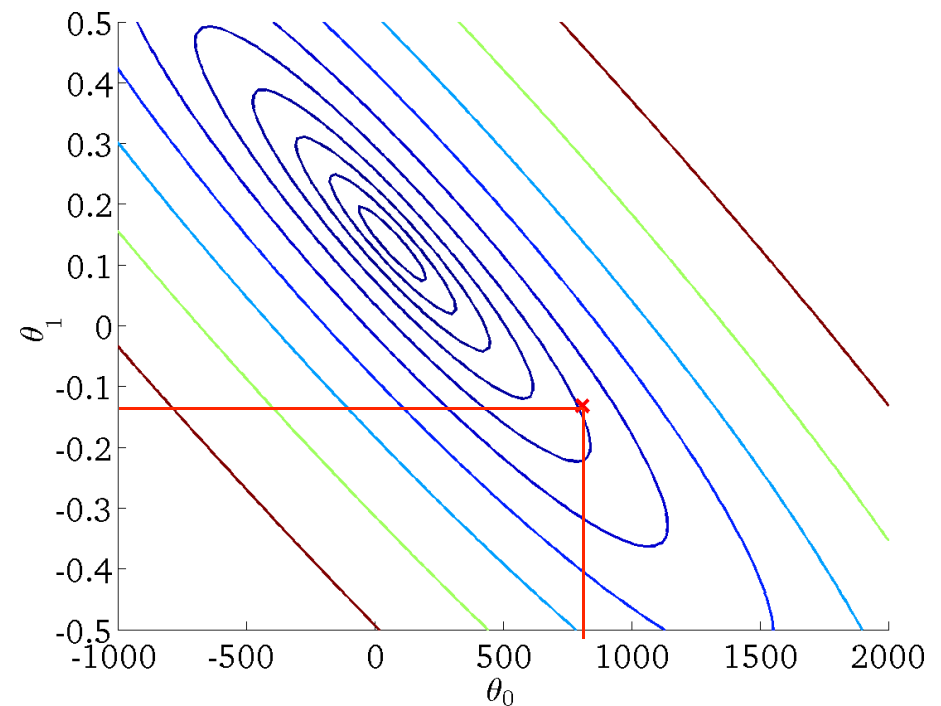# Intuition Behind Cost Function

# Intuition Behind Cost Function

$$h_\theta(x)$$

(for fixed $\theta_0, \theta_1$, this is a function of x)
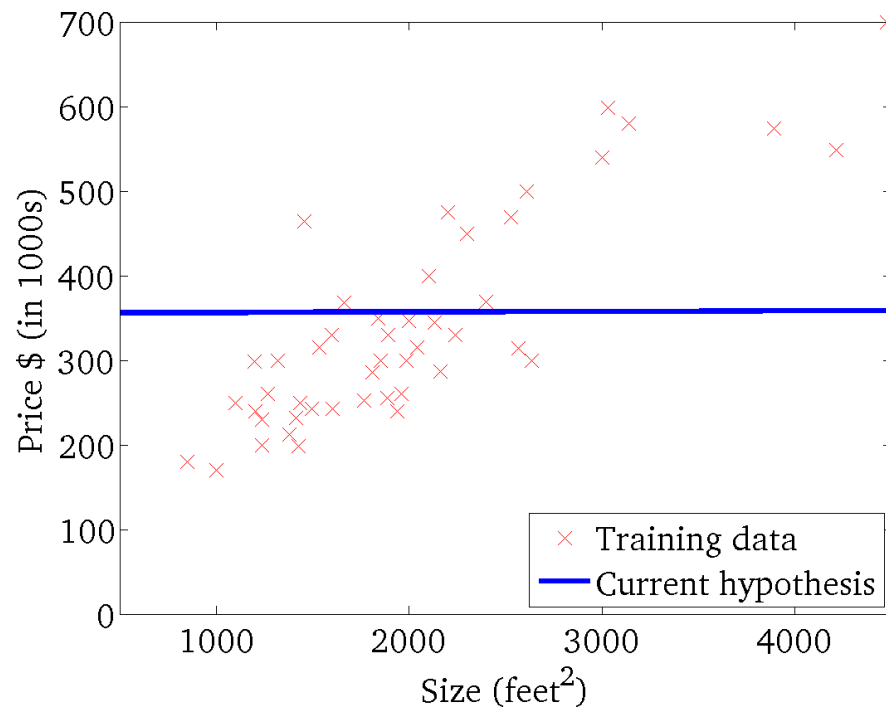
$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)

# Intuition Behind Cost Function

# Intuition Behind Cost Function

$$h_\theta(x)$$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)

# Intuition Behind Cost Function

$$h_\theta(x)$$

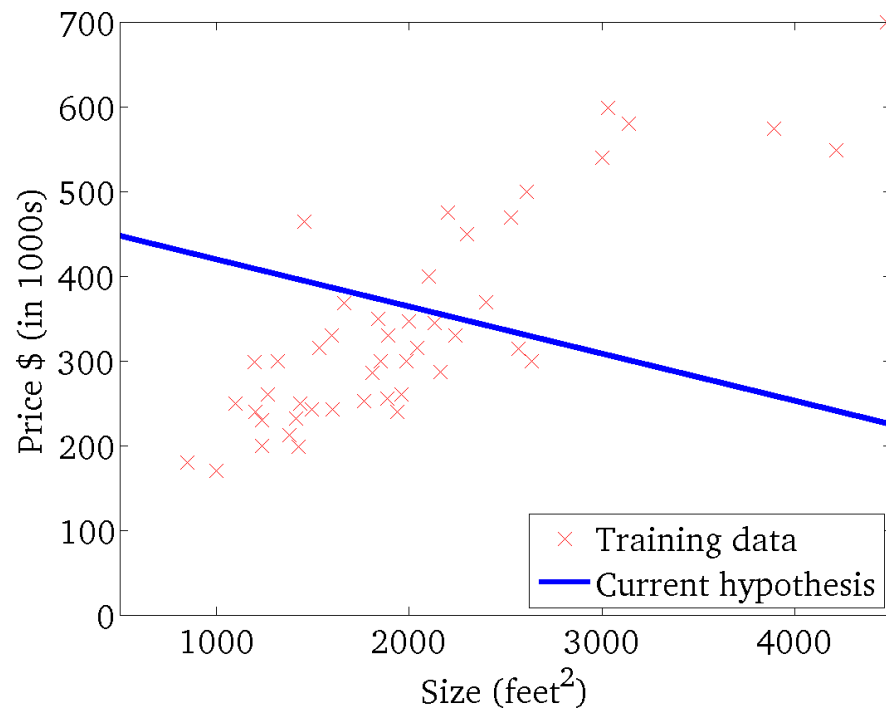(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

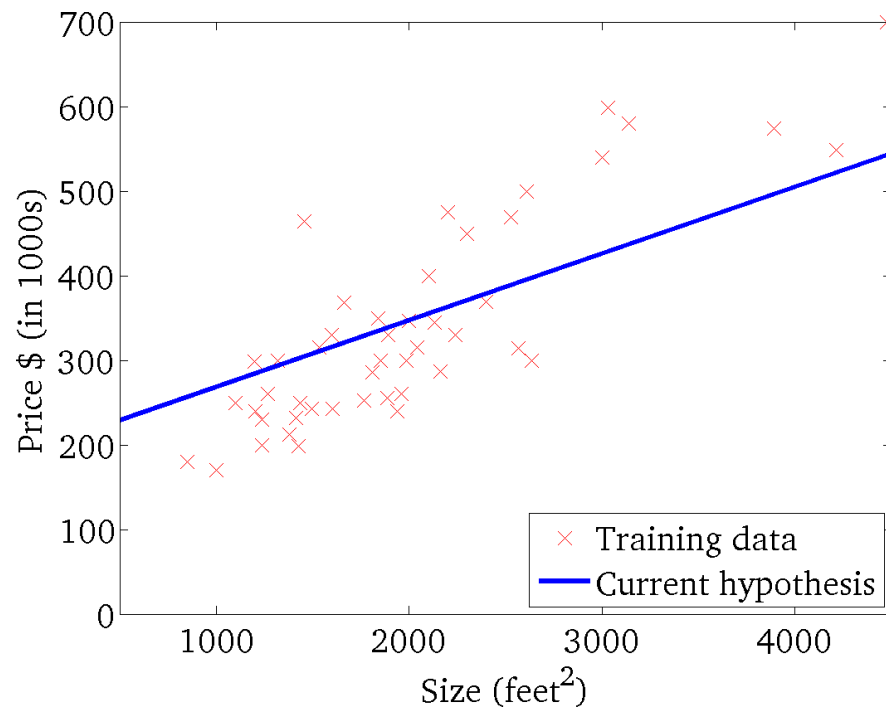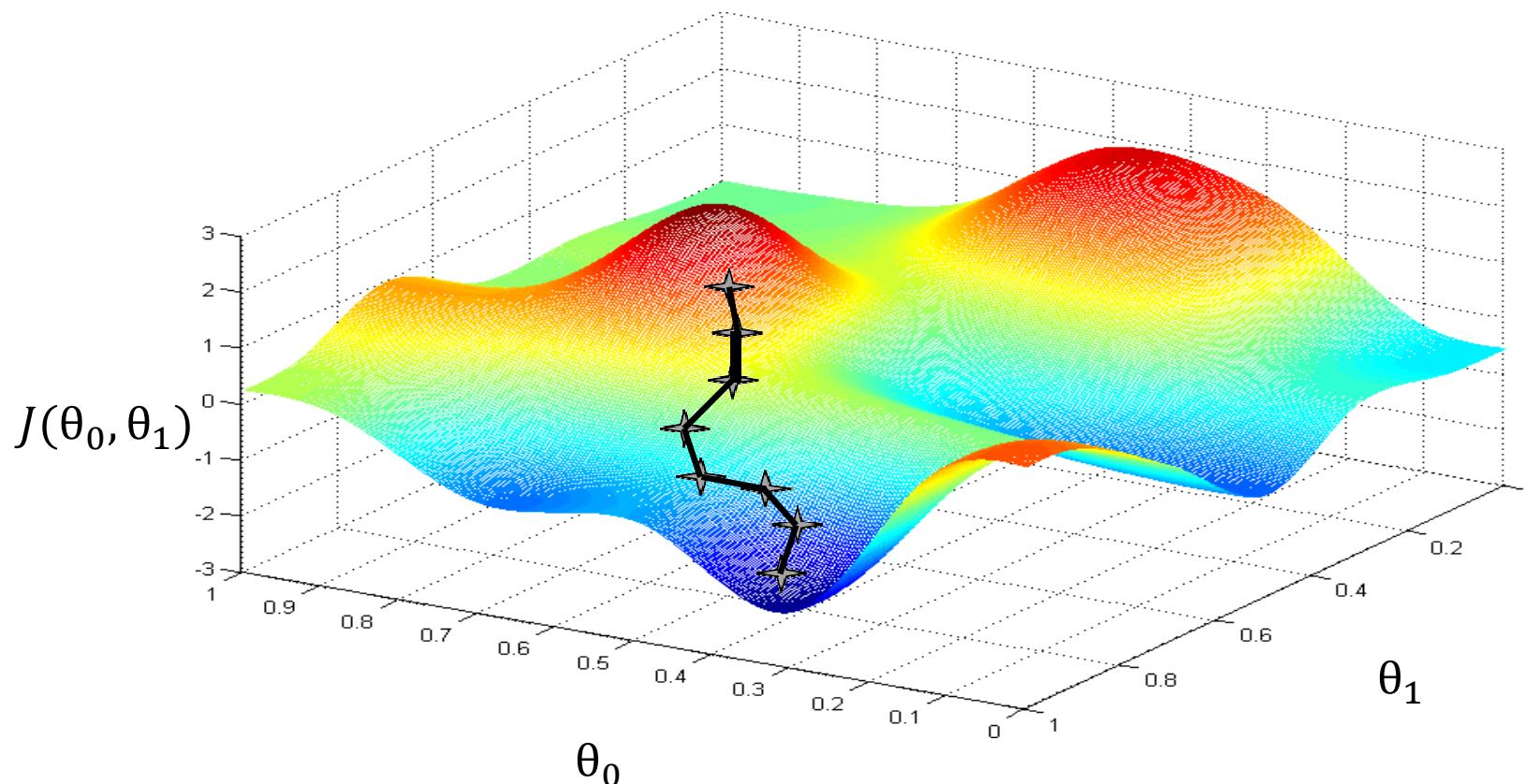(function of the parameters $\theta_0, \theta_1$)

# Basic Search Procedure

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



$J(\theta_0, \theta_1)$

# Basic Search Procedure

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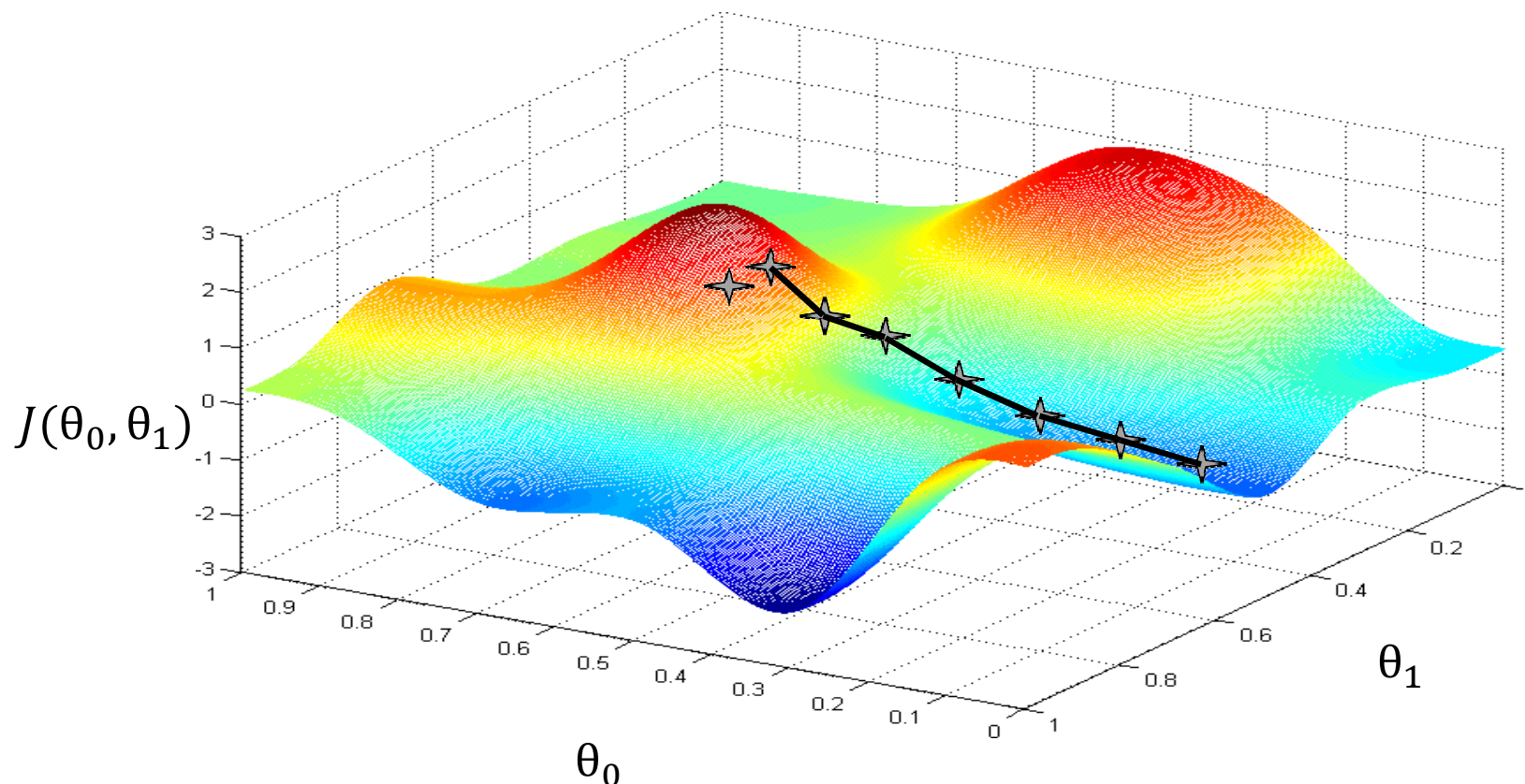


$J(\theta_0, \theta_1)$

$\theta_0$

$\theta_1$

# Basic Search Procedure

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



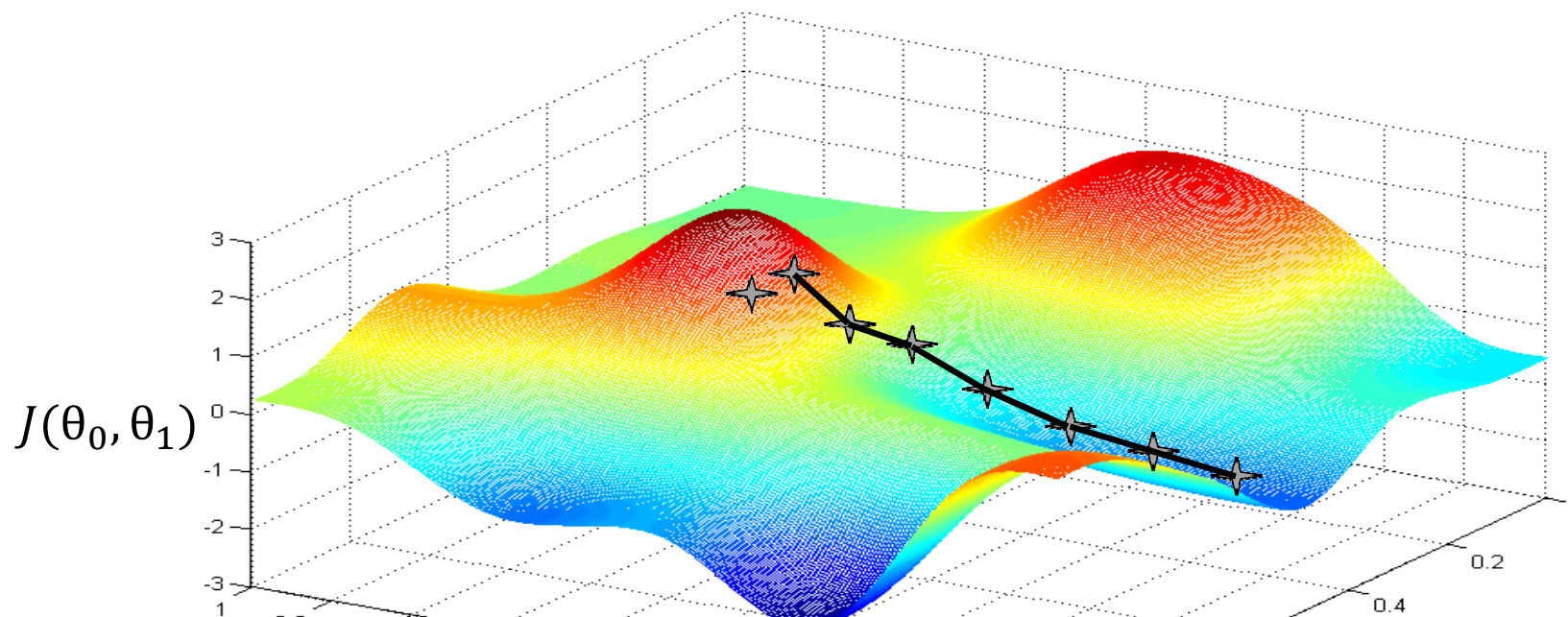$J(\theta_0, \theta_1)$

Since the least squares objective function is convex (concave), we don't need to worry about local minima

# Gradient Descent

- Initialize  θ

- Repeat until convergence

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

simultaneous update
for $j$ = 0 ... d
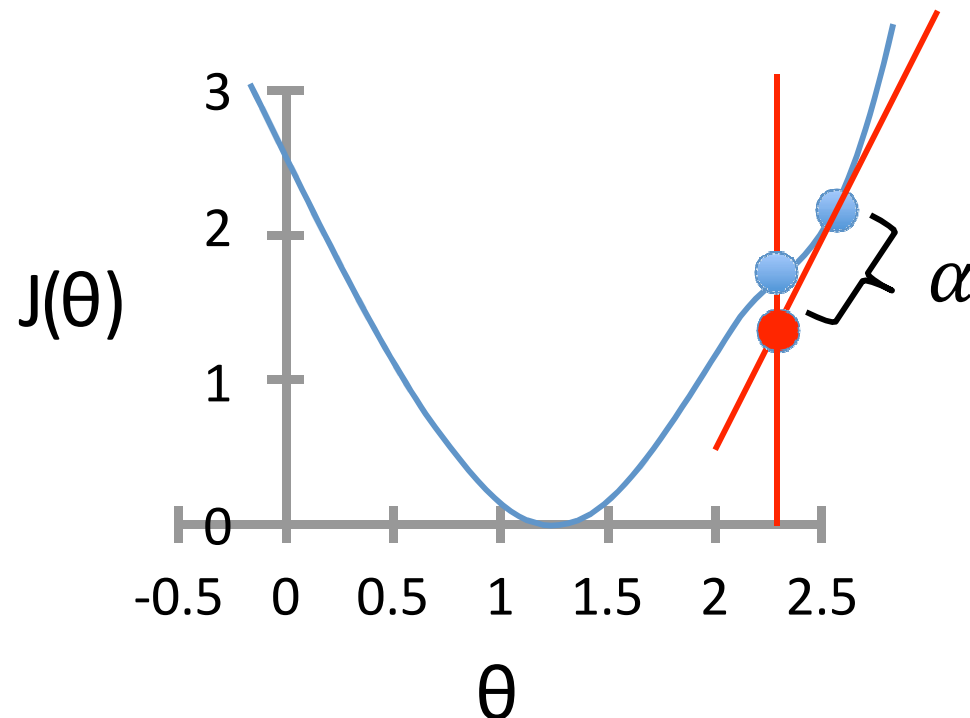
learning rate (small)
e.g., α = 0.05

# Gradient Descent

- Initialize  θ

- Repeat until convergence

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

simultaneous update
for $j = 0 \ldots d$

# Gradient Descent

- Initialize $\theta$
- Repeat until convergence

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

simultaneous update
for $j = 0 \ldots d$

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

For linear regression: $\dfrac{\partial}{\partial \theta_j} J(\theta) = \dfrac{\partial}{\partial \theta_j} \dfrac{1}{2n} \sum_{i=1}^{n} (h_\theta(x^{(i)}) - y^{(i)})^2$

# Gradient Descent

- Initialize  θ

- Repeat until convergence

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

simultaneous update
for $j = 0 \ldots d$

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_d x_d = \sum_{j=0}^{d} \theta_j x_j$$

For linear regression:

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^{n} \left(h_\theta(x^{(i)}) - y^{(i)}\right)^2$$

$$= \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^{n} \left(\sum_{j=0}^{d} \theta_j x_j^{(i)} - y^{(i)}\right)^2$$

# Gradient Descent

- Initialize  θ
- Repeat until convergence

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

simultaneous update
for $j = 0 \dots d$

For linear regression:

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^{n} \left( \sum_{j=0}^{d} \theta_j x_j^{(i)} - y^{(i)} \right)^2$$

$$= \frac{1}{n} \sum_{i=1}^{n} \left( \sum_{j=0}^{d} \theta_j x_j^{(i)} - y^{(i)} \right) \frac{\partial}{\partial \theta_j} \left( \sum_{j=0}^{d} \theta_j x_j^{(i)} - y^{(i)} \right)$$

$$= \frac{1}{n} \sum_{i=1}^{n} \left( \sum_{j=0}^{d} \theta_j x_j^{(i)} - y^{(i)} \right) x_j^{(i)}$$

# Gradient Descent

- Initialize  θ
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \propto \frac{1}{n}\sum_{i=1}^{n}\left(h_\theta(x^{(i)}) - y^{(i)}\right)x_j^{(i)}$$
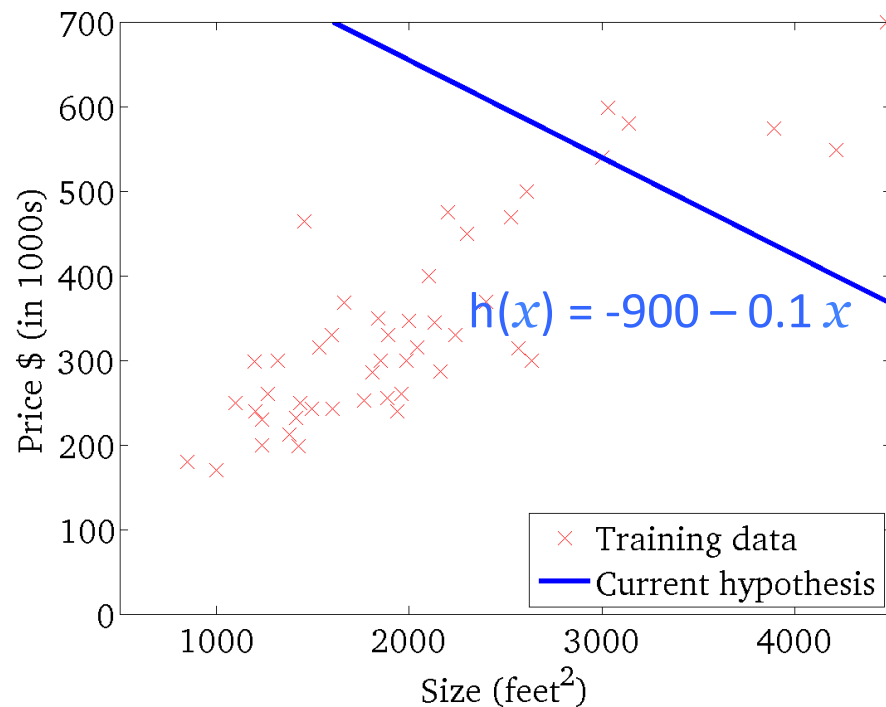
simultaneous update
for $j$ = 0 … d

- To achieve simultaneous update
  - At the start of each GD iteration, compute $h_\theta(x^{(i)})$
  - Use this stored value in the update step loop

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

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

# Gradient Descent

$h_\theta(x)$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)



h$(x) = -900 - 0.1\,x$

Training data: $\times$
Current hypothesis: ——

# Gradient Descent



$h_\theta(x)$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

# Gradient Descent

$$h_\theta(x)$$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)

# Gradient Descent

$$h_\theta(x)$$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)

# Gradient Descent

$h_\theta(x)$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

# Gradient Descent

$$h_\theta(x)$$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)

# Gradient Descent

$$h_\theta(x)$$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)

# Gradient Descent

$$h_\theta(x)$$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)

# Gradient Descent

$$h_\theta(x)$$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)
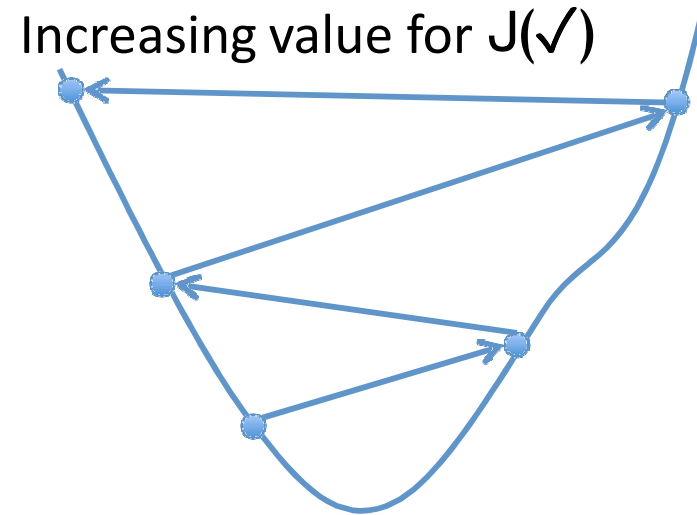
# Choosing α

## α too small

slow convergence



## α too large

Increasing value for J(✓)



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

To see if gradient descent is working, print out J(✓) each iteration
- The value should decrease at each iteration
- If it doesn't, adjust α