



# Machine Learning

## DSECL ZG565

Dr. Chetana Gavankar, Ph.D,  
IIT Bombay-Monash University Australia  
[Chetana.gavankar@pilani.bits-pilani.ac.in](mailto:Chetana.gavankar@pilani.bits-pilani.ac.in)



**BITS** Pilani  
Pilani Campus



**Lecture No. – 6 | Linear Regression**

**Date – 30/11/2019**

**Time – 11:00 AM – 1:00 PM**

# Session Content

---

- Regression (Andre Ng Notes)
- Bayesian linear regression (6.4 Tom Mitchell)
- Linear basis function models (3.1 Bishop)
- Bias-variance decomposition (3.2 Bishop)

# Regression

---

So far, we've been interested in learning  $P(Y|X)$  where  $Y$  has discrete values (called 'classification')

What if  $Y$  is continuous? (called 'regression')

- predict weight from gender, height, age, ...
- predict Google stock price today from Google, Yahoo, MSFT prices yesterday
- predict each pixel intensity in robot's current camera image, from previous image and previous action

# Regression

---

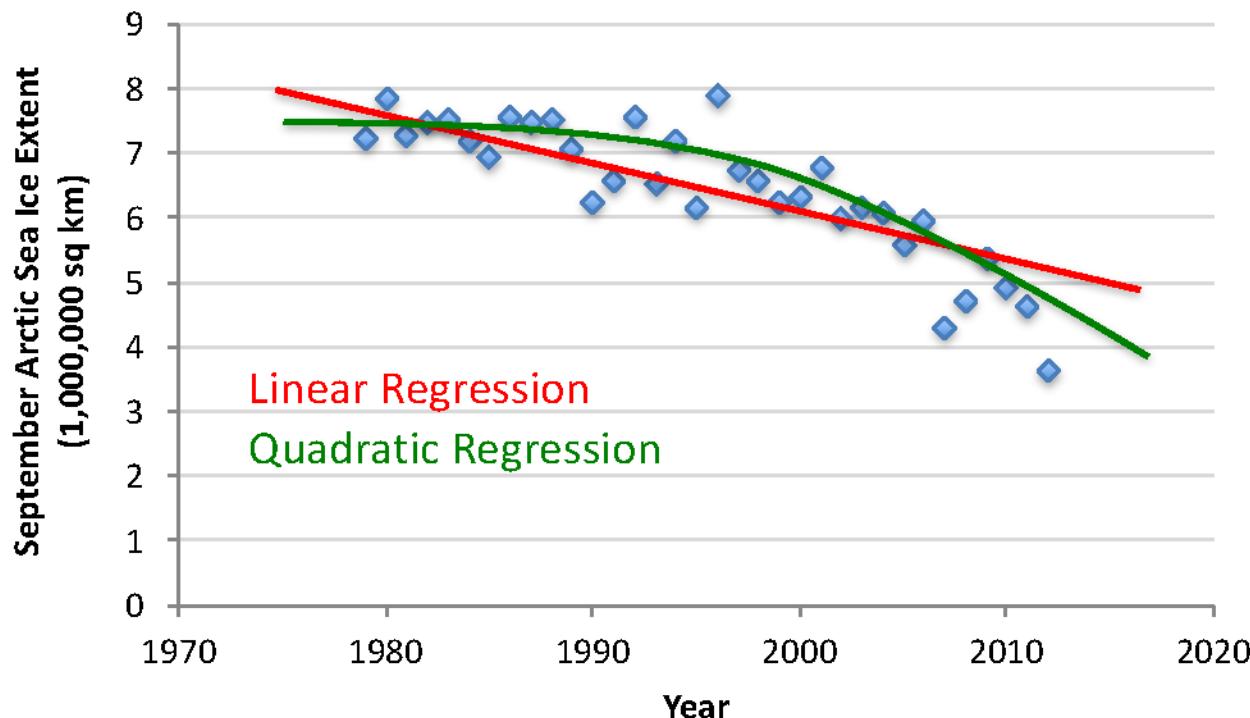
Wish to learn  $f: X \rightarrow Y$ , where  $Y$  is real, given  $\{<x^1, y^1> \dots <x^n, y^n>\}$

- Geometric Approach
  - Least squares function fitting given  $\{<x^1, y^1> \dots <x^n, y^n>\}$
- Bayesian Approach
  - Choose some parameterized form for  $P(Y|X; \theta)$ 
    - ( $\theta$  is the vector of parameters)
  - Derive learning algorithm as MCLE estimate for  $\theta$

# Geometric Approach

Given:

- Data  $\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}$  where  $\mathbf{x}^{(i)} \in \mathbb{R}^d$
- Corresponding labels  $\mathbf{y} = \{y^{(1)}, \dots, y^{(n)}\}$  where  $y^{(i)} \in \mathbb{R}$



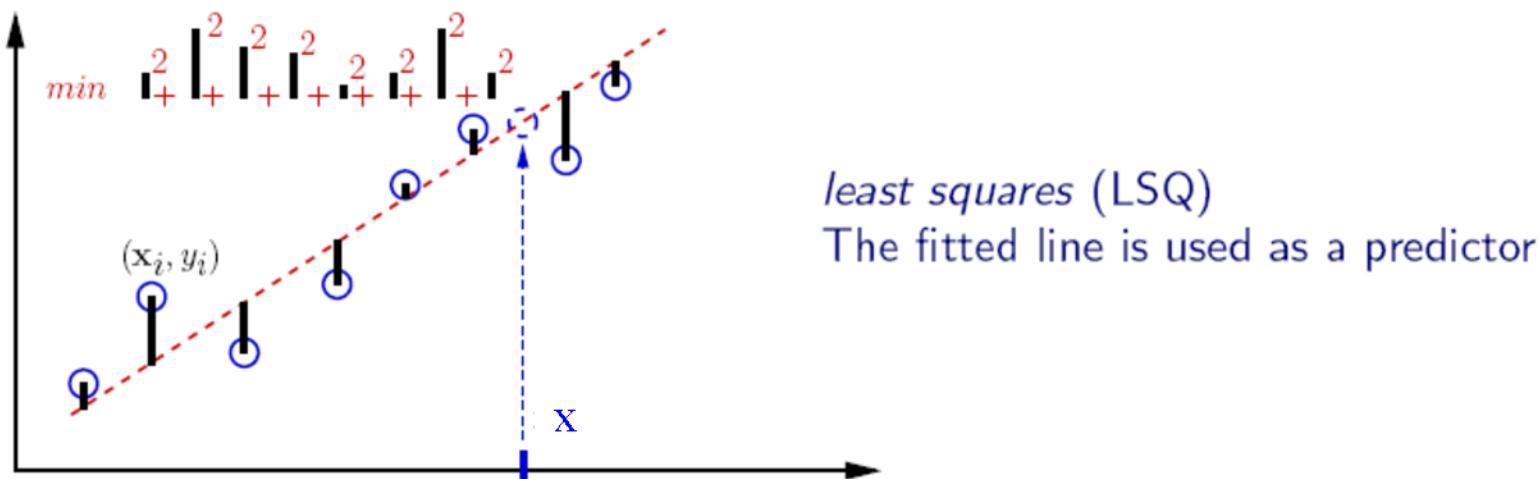
# Linear Regression

- Hypothesis:

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_d x_d = \sum_{j=0}^d \theta_j x_j$$

Assume  $x_0 = 1$

- Fit model by minimizing sum of squared errors

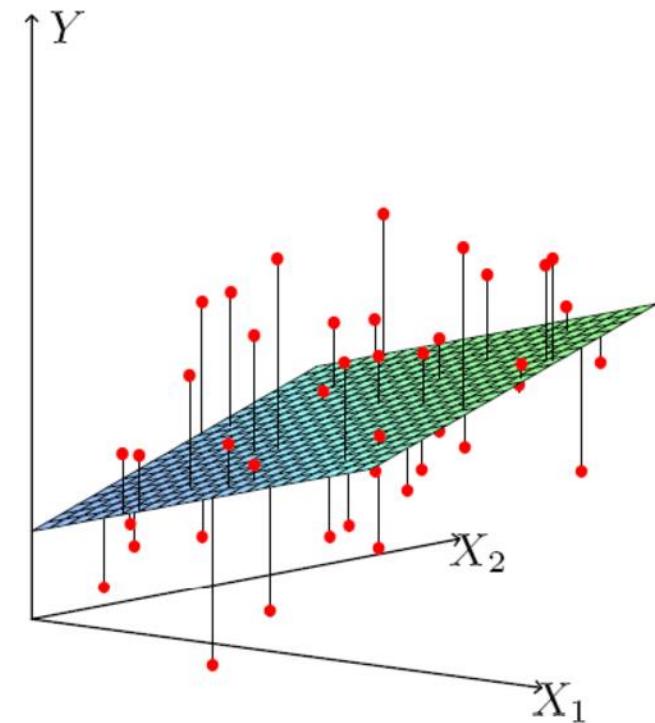
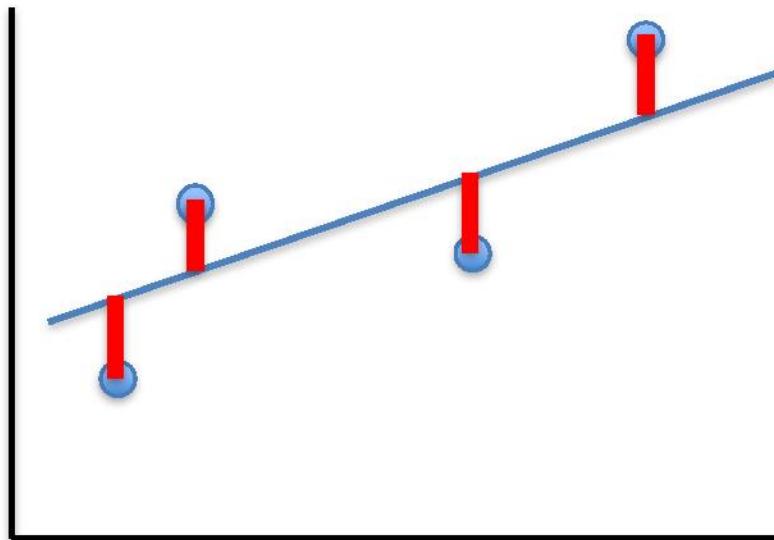


# Least Squares Linear Regression

- Cost Function

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

- Fit by solving  $\min_{\theta} J(\theta)$



# Intuition Behind Cost Function

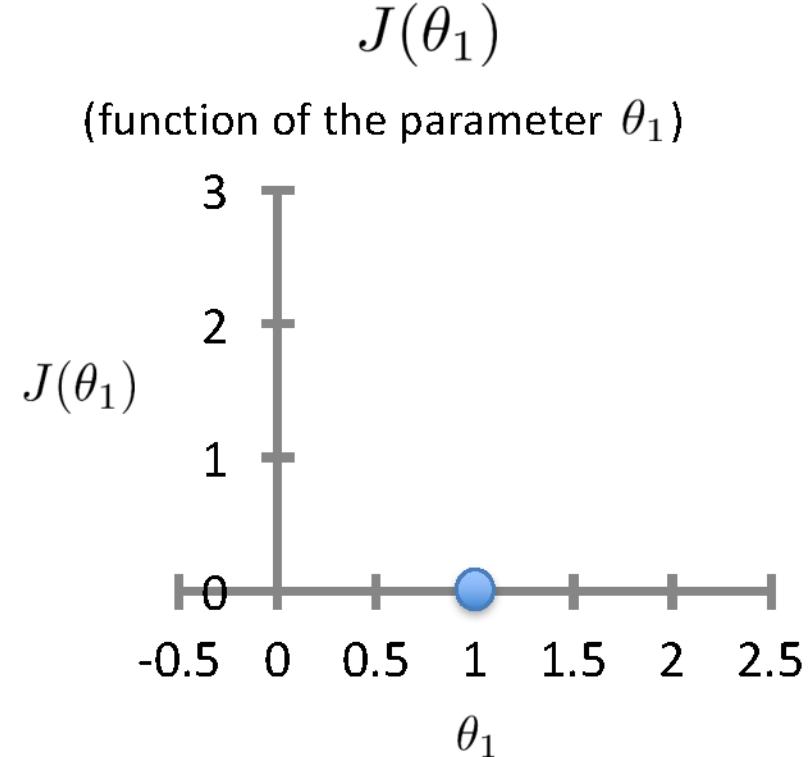
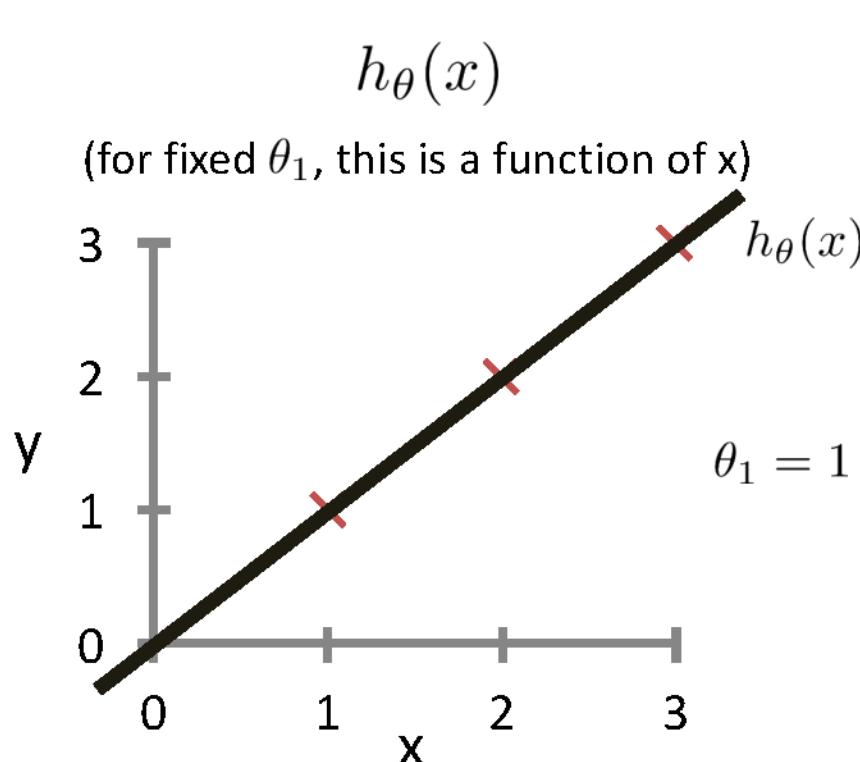
$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$$

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

# Intuition Behind Cost Function

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

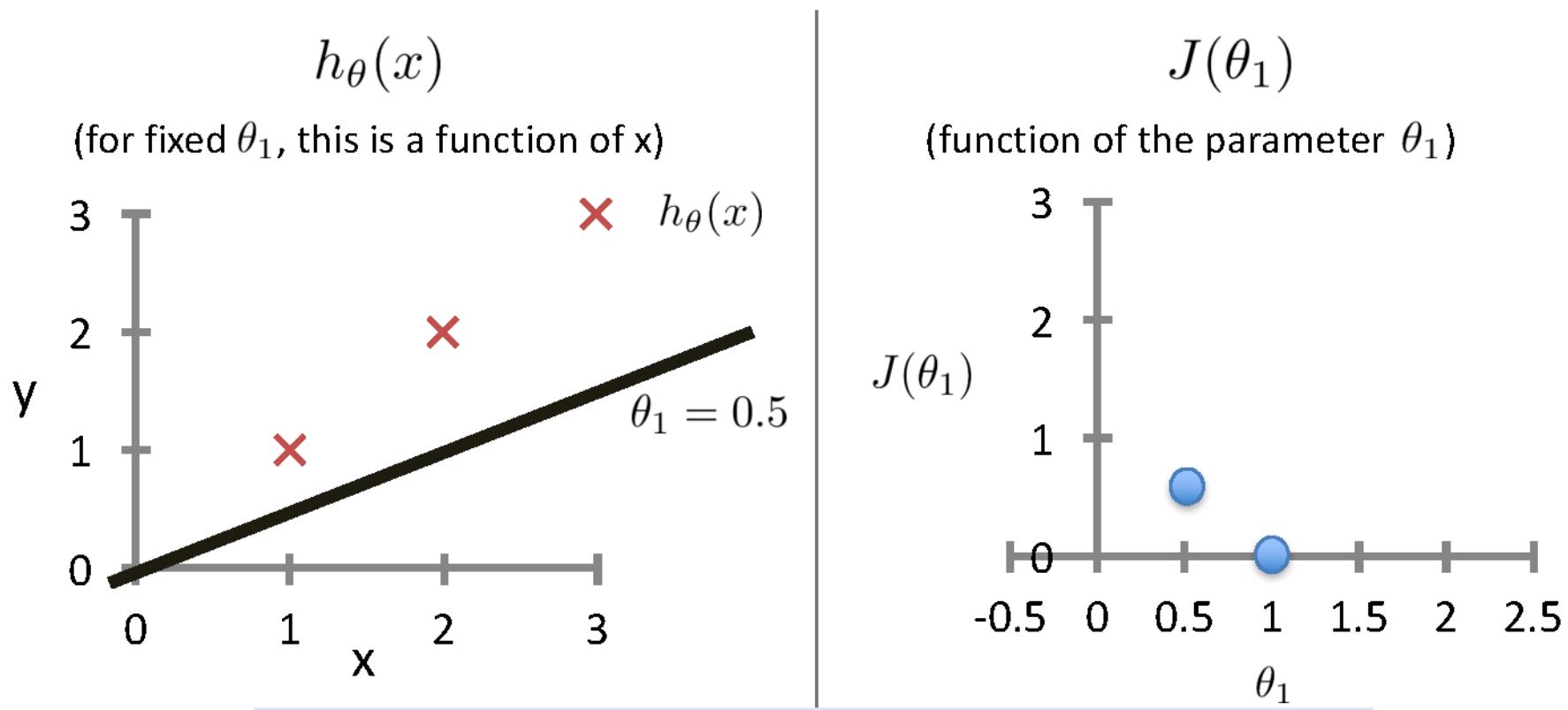
For insight on  $J()$ , let's assume  $x \in \mathbb{R}$  so  $\theta = [\theta_0, \theta_1]$



# 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  $x \in \mathbb{R}$  so  $\theta = [\theta_0, \theta_1]$



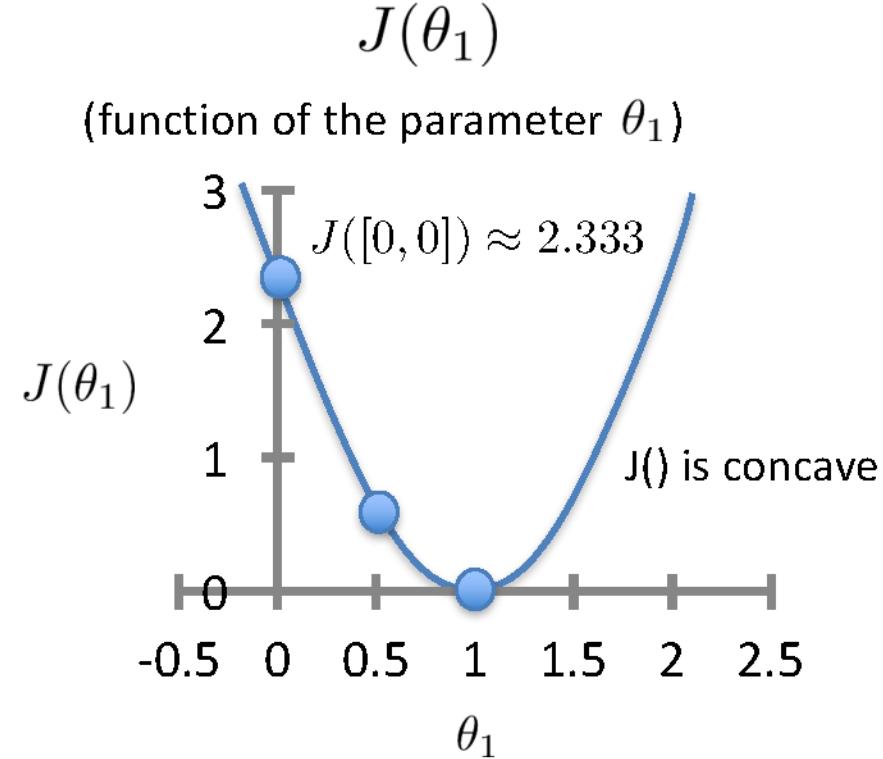
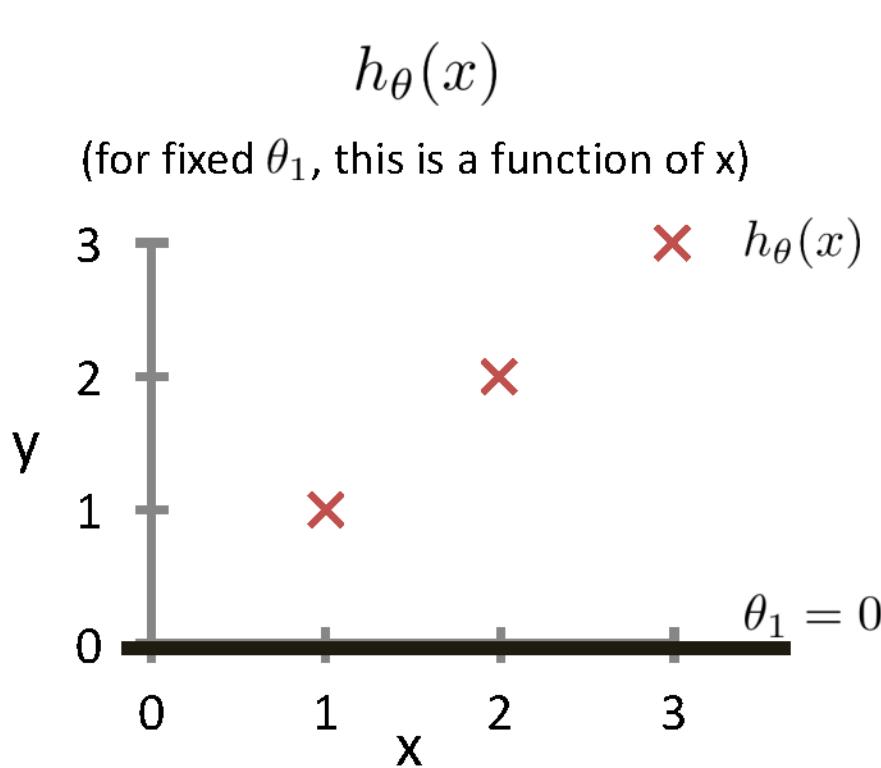
Based on example  
by Andrew Ng

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

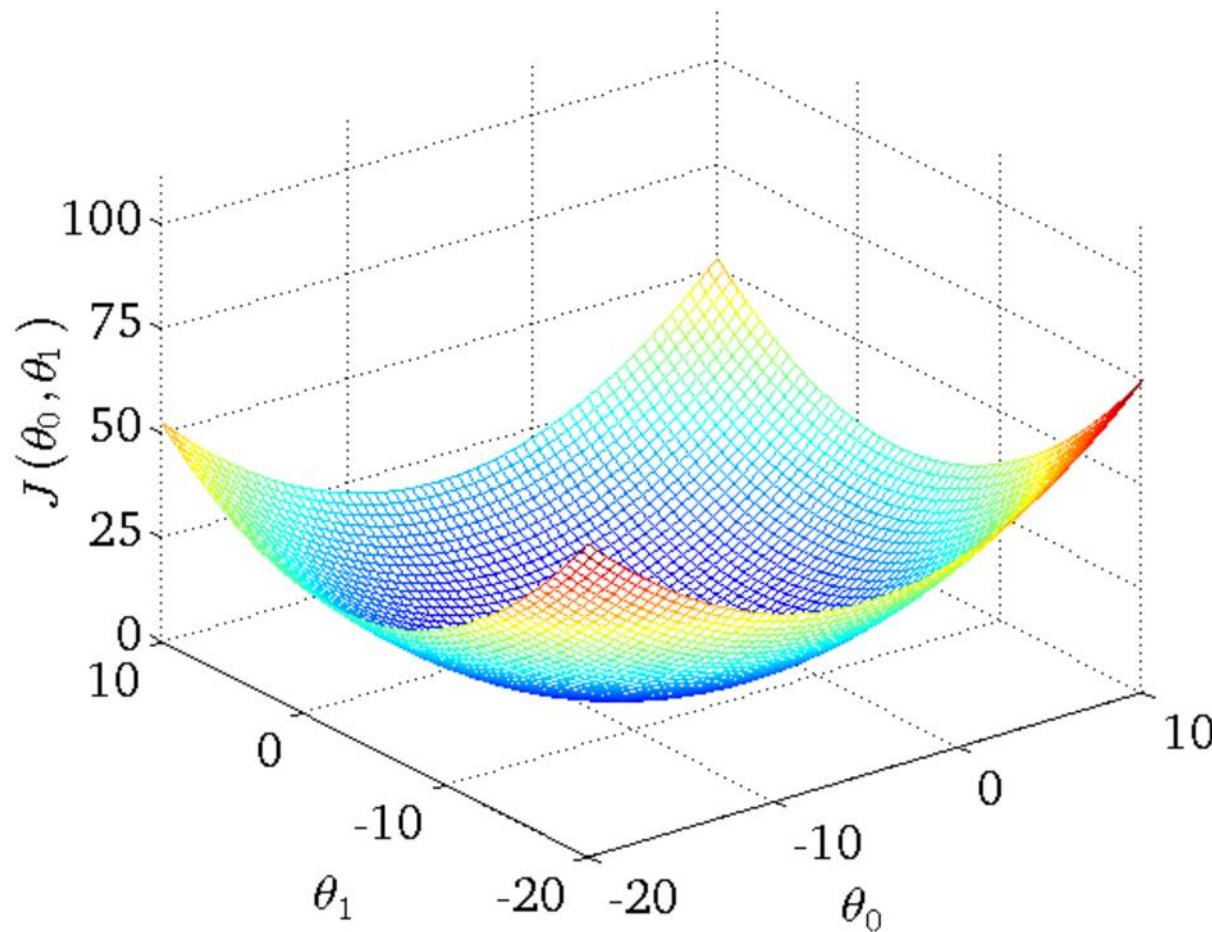
# Intuition Behind Cost Function

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

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



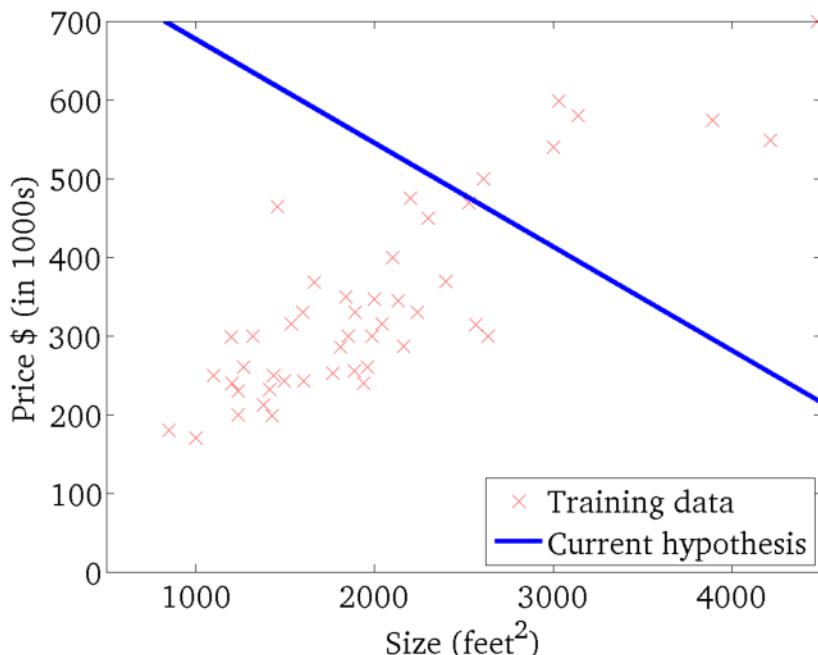
# Intuition Behind Cost Function



# Intuition Behind Cost Function

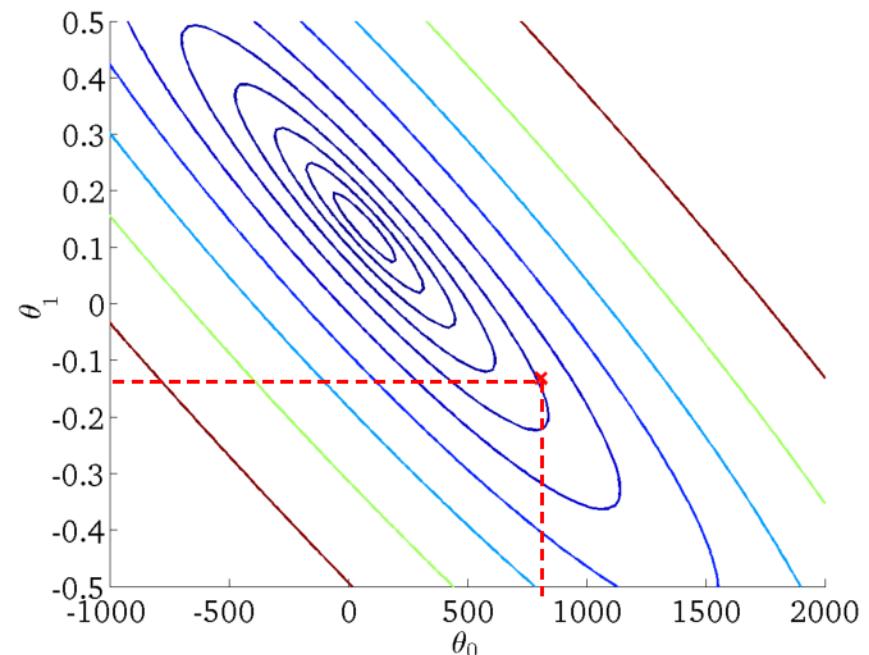
$$h_{\theta}(x)$$

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



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

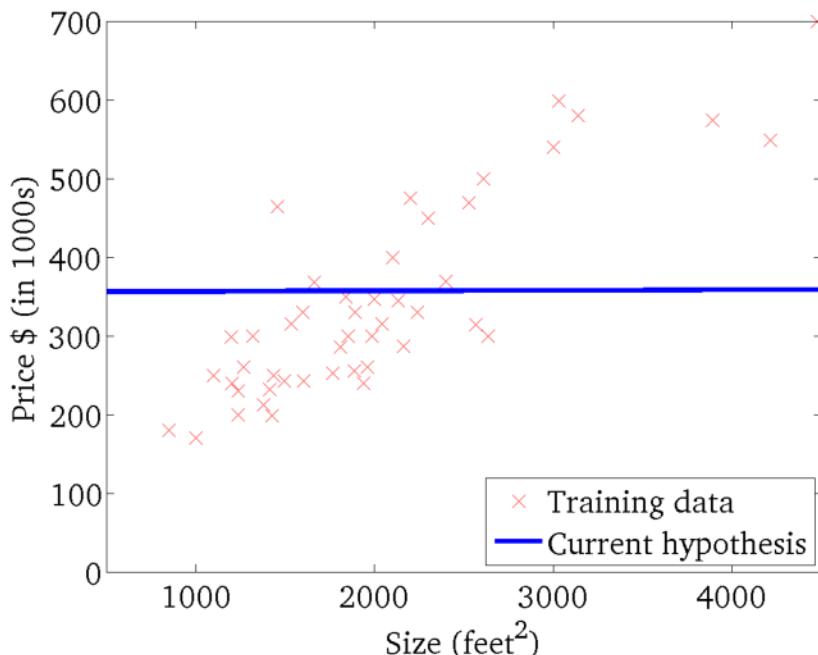
(function of the parameters  $\theta_0, \theta_1$ )



# Intuition Behind Cost Function

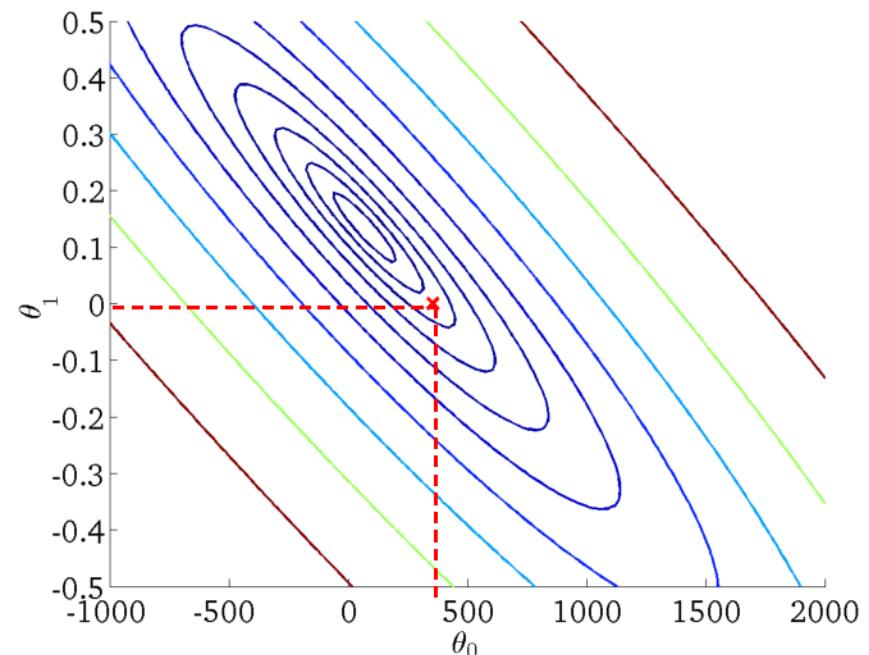
$$h_{\theta}(x)$$

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



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

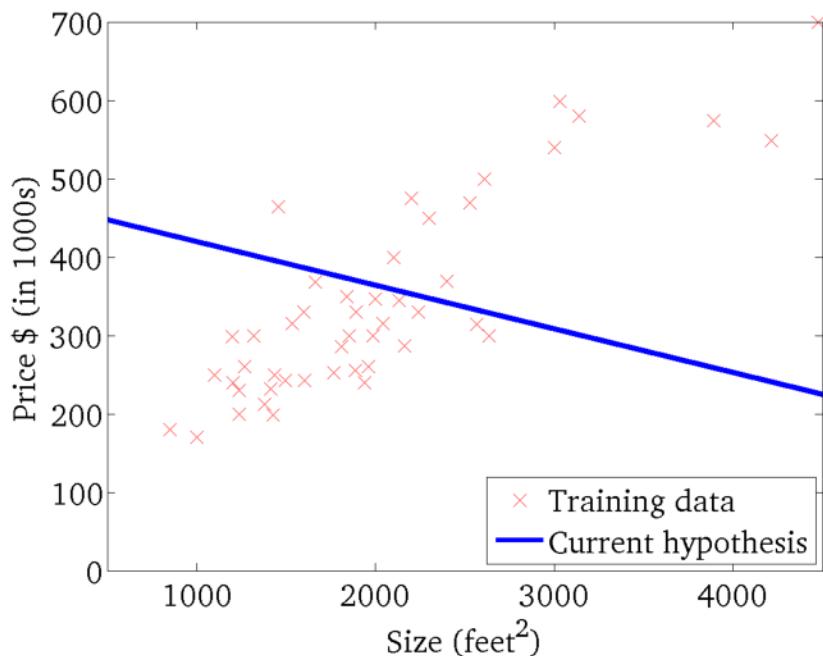
(function of the parameters  $\theta_0, \theta_1$ )



# Intuition Behind Cost Function

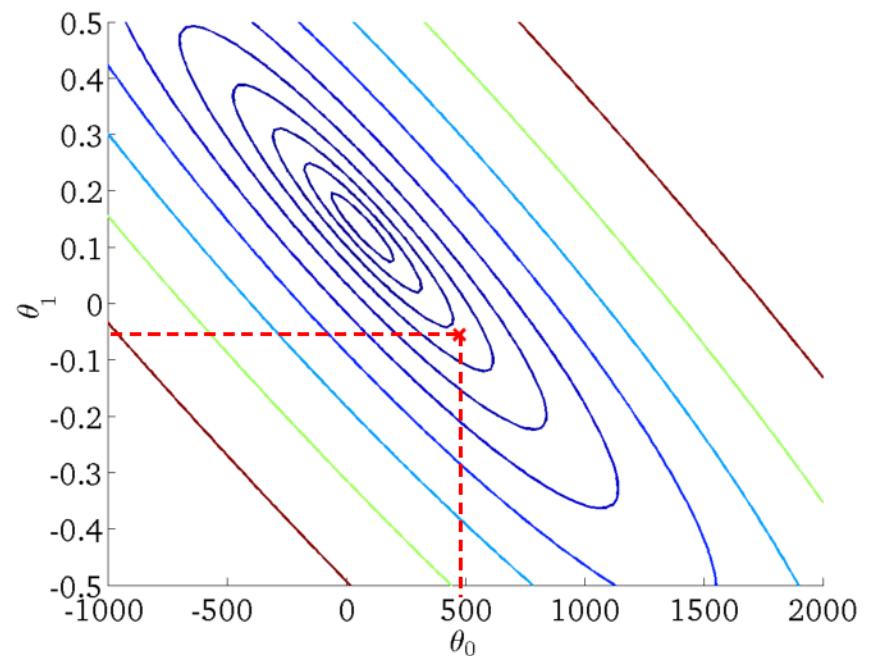
$$h_{\theta}(x)$$

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



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

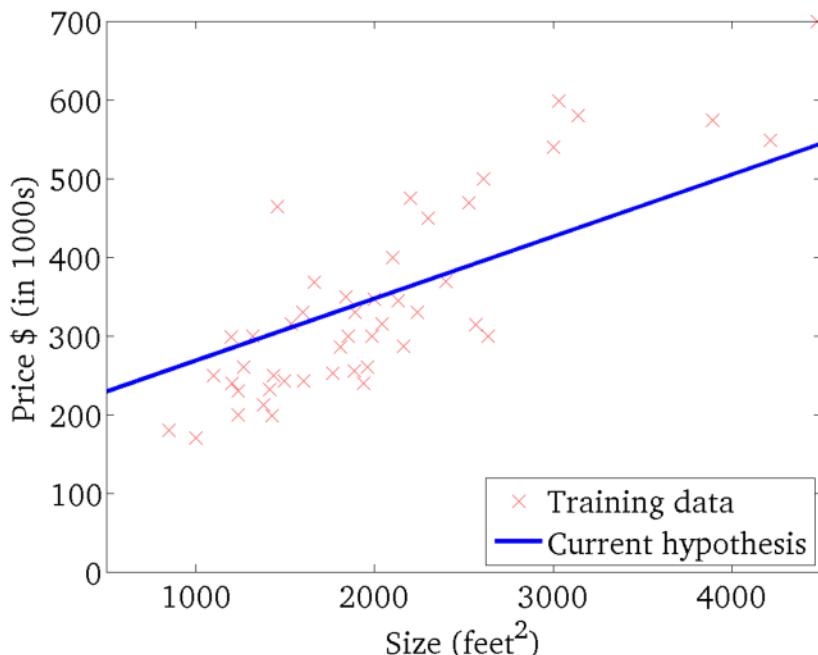
(function of the parameters  $\theta_0, \theta_1$ )



# Intuition Behind Cost Function

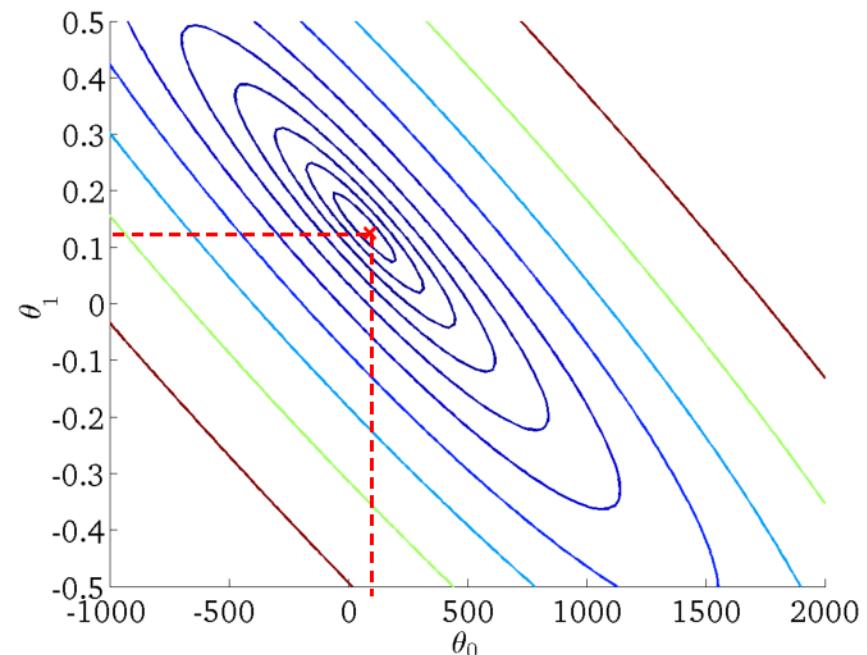
$$h_{\theta}(x)$$

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



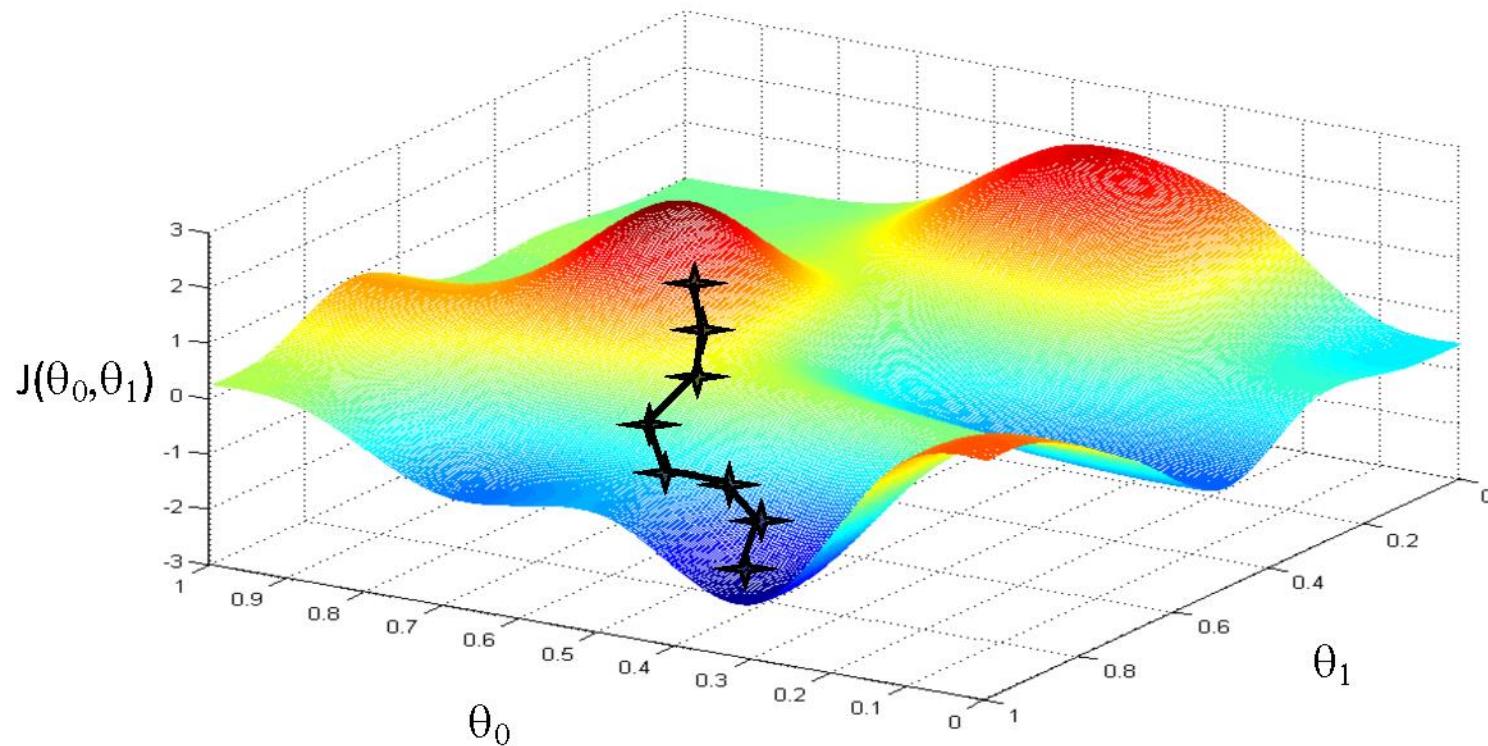
$$J(\theta_0, \theta_1)$$

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



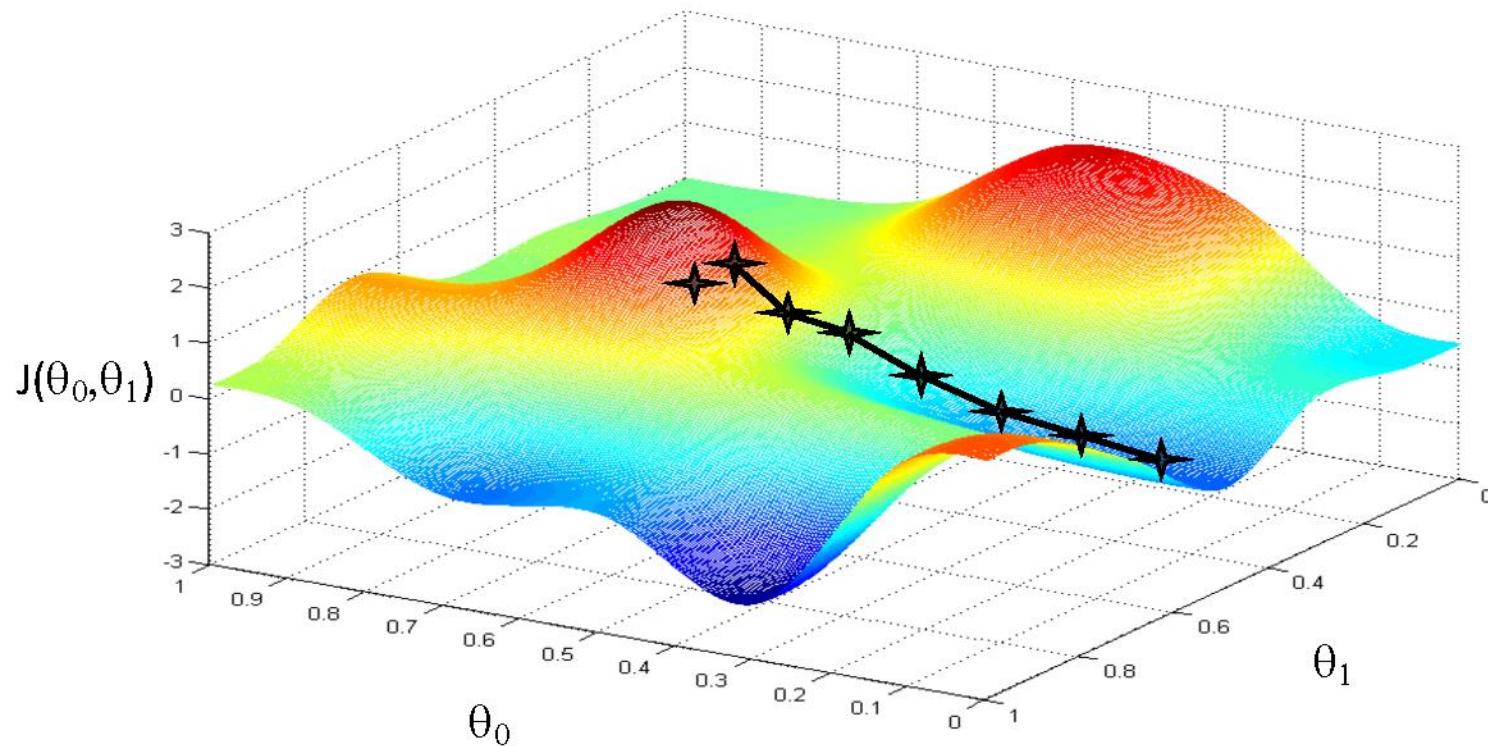
# Basic Search Procedure

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



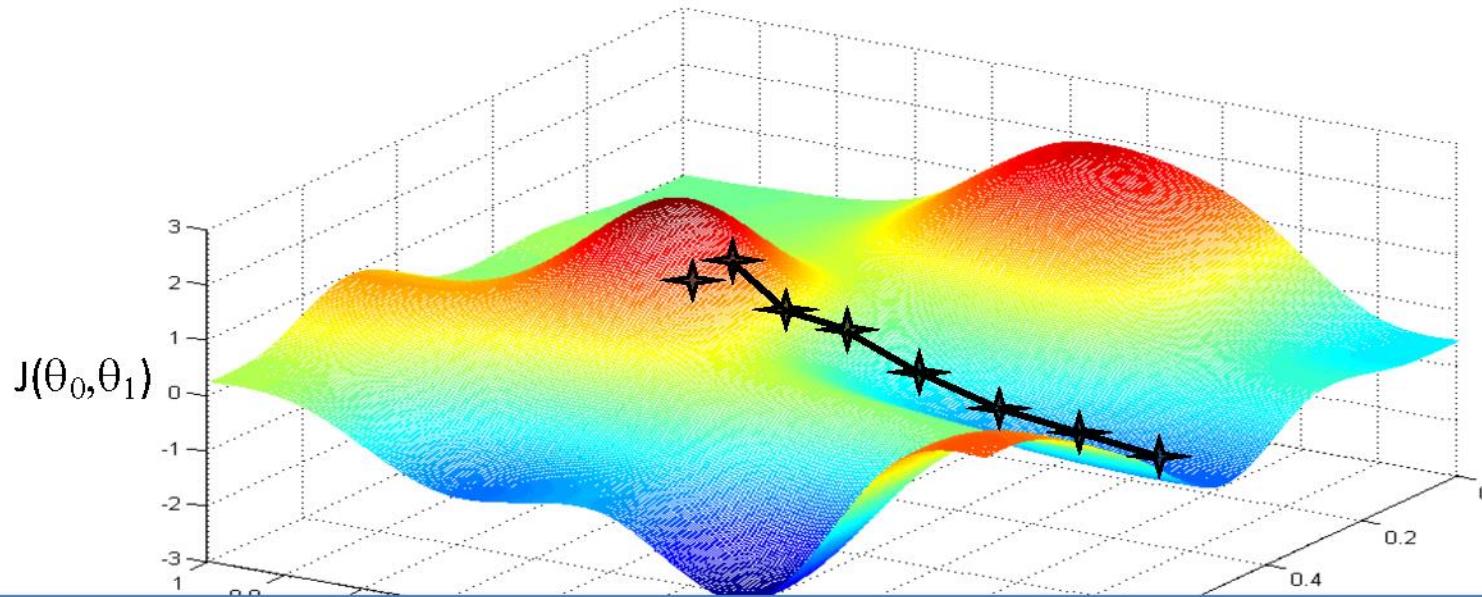
# Basic Search Procedure

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



# Basic Search Procedure

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



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

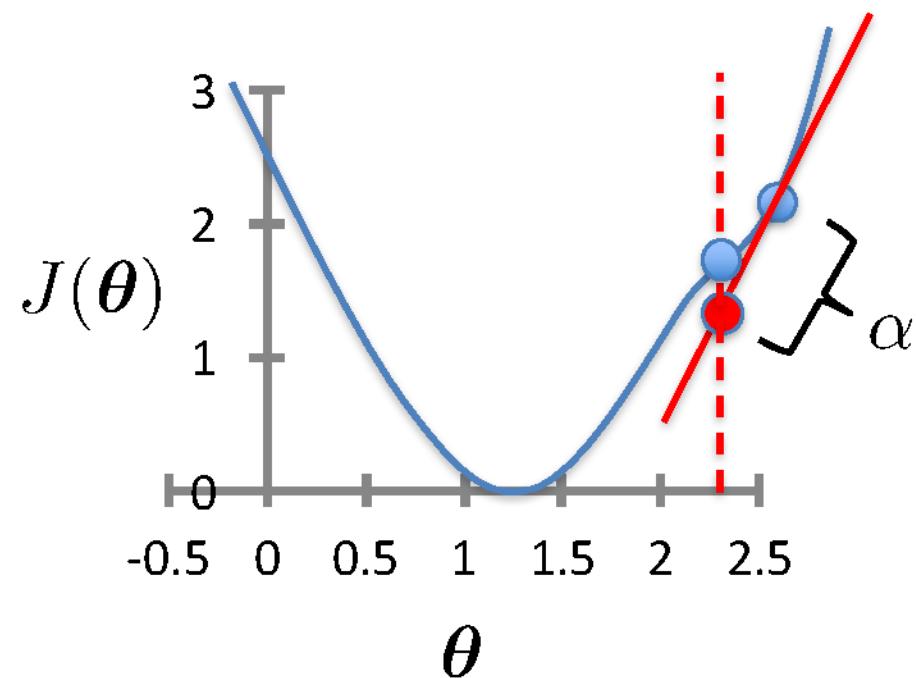
# Gradient Descent

- 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 Descent

- Initialize  $\theta$
- Repeat until convergence

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

simultaneous update  
for  $j = 0 \dots d$

For Linear Regression:

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

# Gradient Descent

- Initialize  $\theta$
- Repeat until convergence

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

simultaneous update  
for  $j = 0 \dots d$

For Linear Regression:

$$\begin{aligned}\frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^n \left( h_{\theta} (\mathbf{x}^{(i)}) - y^{(i)} \right)^2 \\ &= \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^n \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right)^2\end{aligned}$$

# Gradient Descent

- Initialize  $\theta$
- Repeat until convergence

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

simultaneous update  
for  $j = 0 \dots d$

For Linear Regression:

$$\begin{aligned}\frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^n \left( h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2 \\ &= \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^n \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right)^2 \\ &= \frac{1}{n} \sum_{i=1}^n \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right) \times \frac{\partial}{\partial \theta_j} \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right)\end{aligned}$$

# Gradient Descent

- Initialize  $\theta$
- Repeat until convergence

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

simultaneous update  
for  $j = 0 \dots d$

For Linear Regression:

$$\begin{aligned}\frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^n \left( h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2 \\ &= \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^n \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right)^2 \\ &= \frac{1}{n} \sum_{i=1}^n \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right) \times \frac{\partial}{\partial \theta_j} \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right) \\ &= \frac{1}{n} \sum_{i=1}^n \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right) x_j^{(i)}\end{aligned}$$

# Gradient Descent for Linear Regression

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n \left( h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)}$$

simultaneous  
update  
for  $j = 0 \dots d$

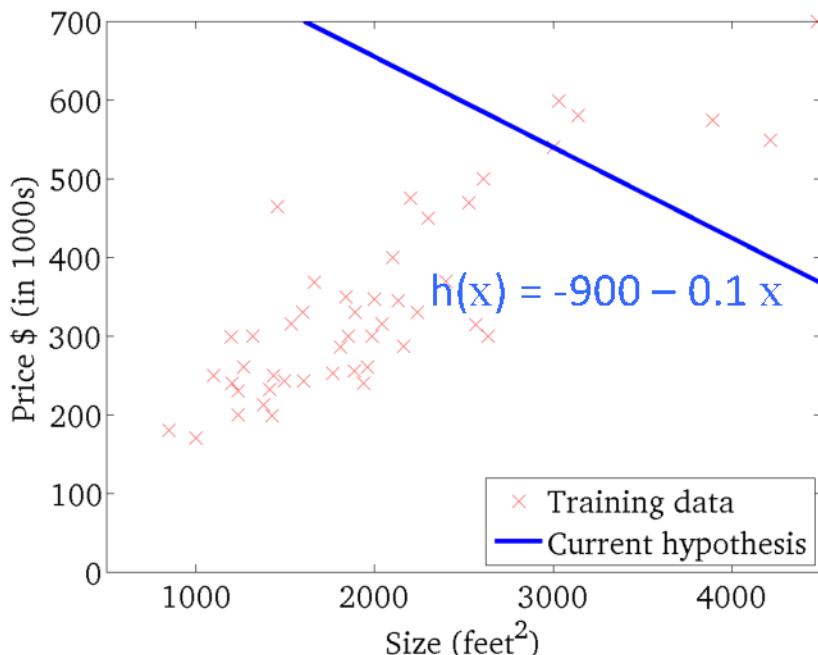
- To achieve simultaneous update
  - At the start of each GD iteration, compute  $h_{\theta}(\mathbf{x}^{(i)})$
  - Use this stored value in the update step loop
- Assume convergence when  $\|\theta_{new} - \theta_{old}\|_2 < \epsilon$

L<sub>2</sub> norm:  $\|\mathbf{v}\|_2 = \sqrt{\sum_i v_i^2} = \sqrt{v_1^2 + v_2^2 + \dots + v_{|v|}^2}$

# Gradient Descent

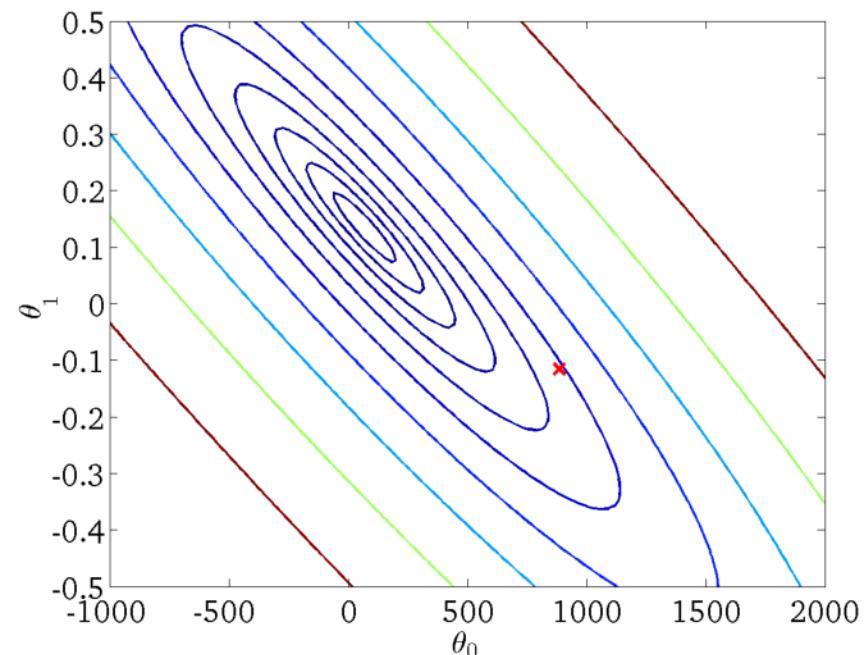
$$h_{\theta}(x)$$

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



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

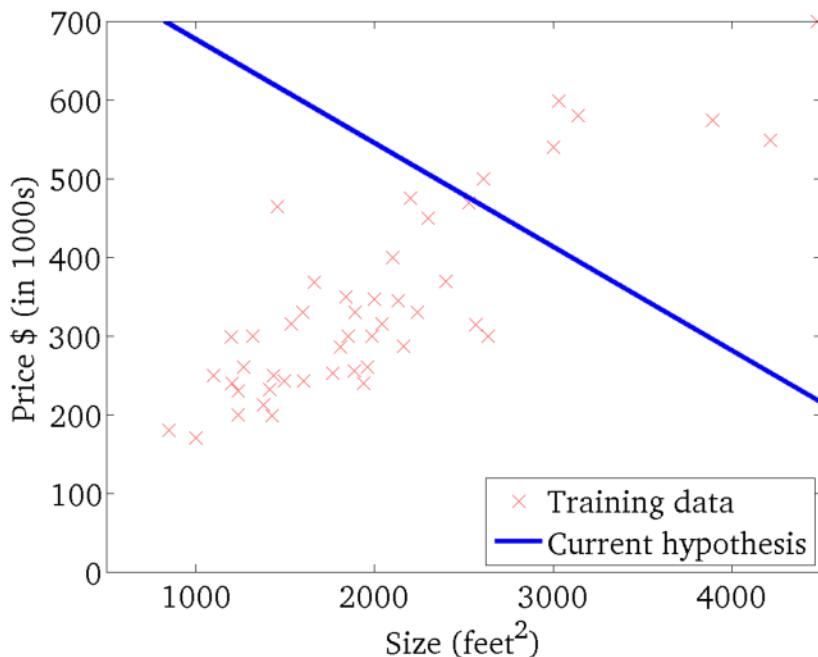
(function of the parameters  $\theta_0, \theta_1$ )



# Gradient Descent

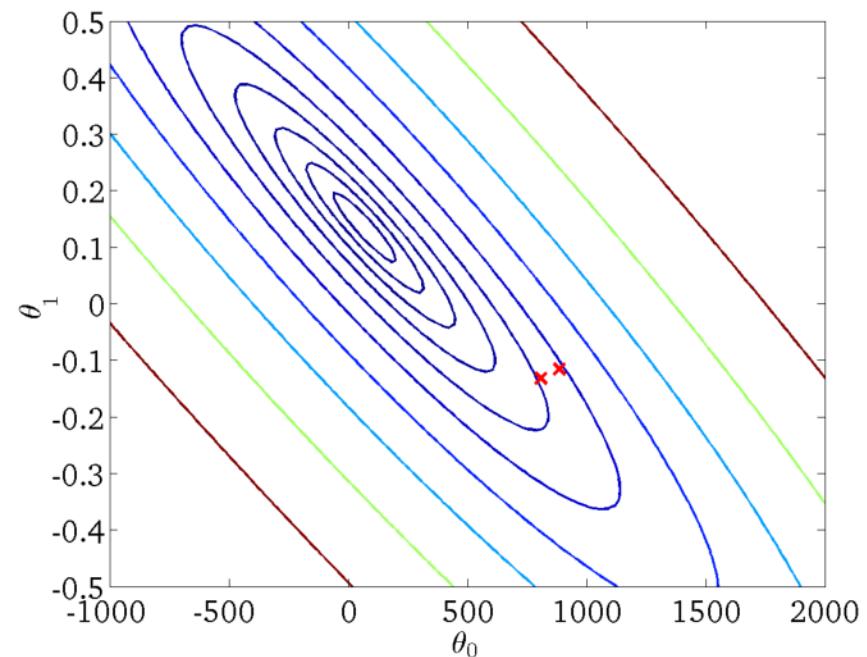
$$h_{\theta}(x)$$

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



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

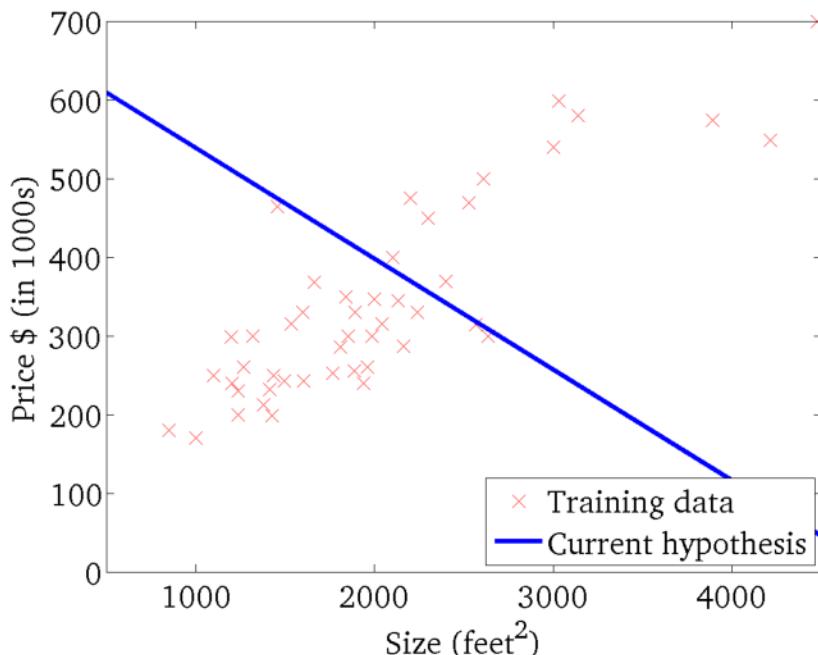
(function of the parameters  $\theta_0, \theta_1$ )



# Gradient Descent

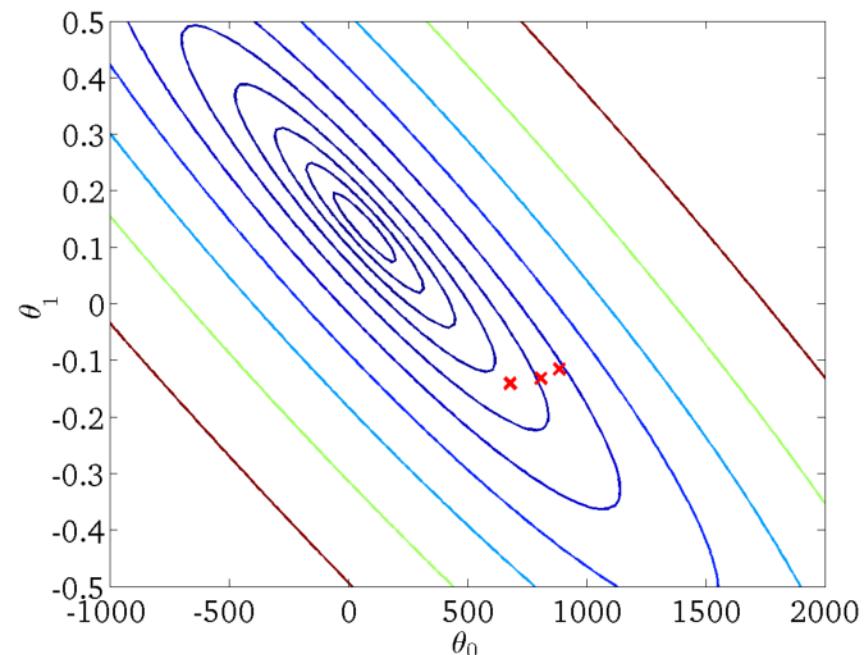
$$h_{\theta}(x)$$

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



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

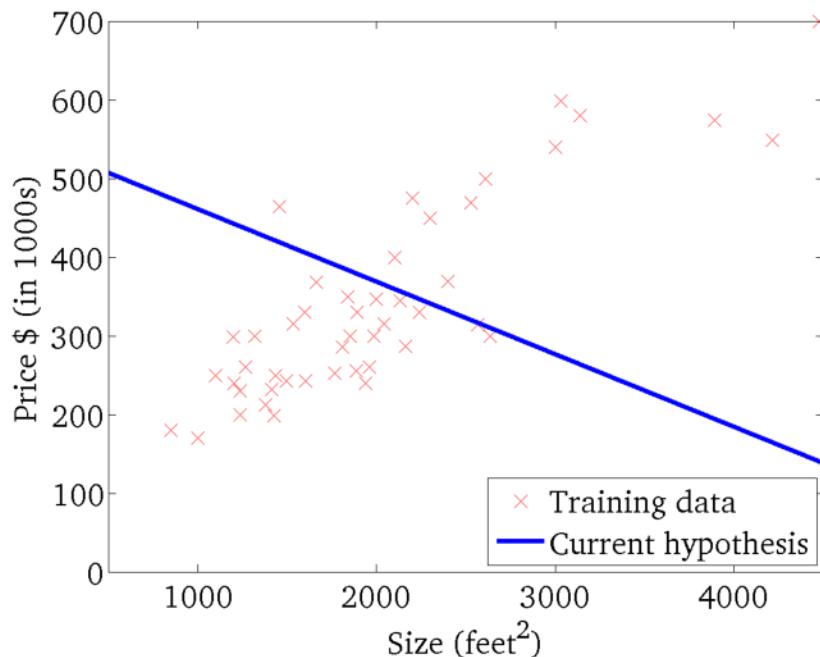
(function of the parameters  $\theta_0, \theta_1$ )



# Gradient Descent

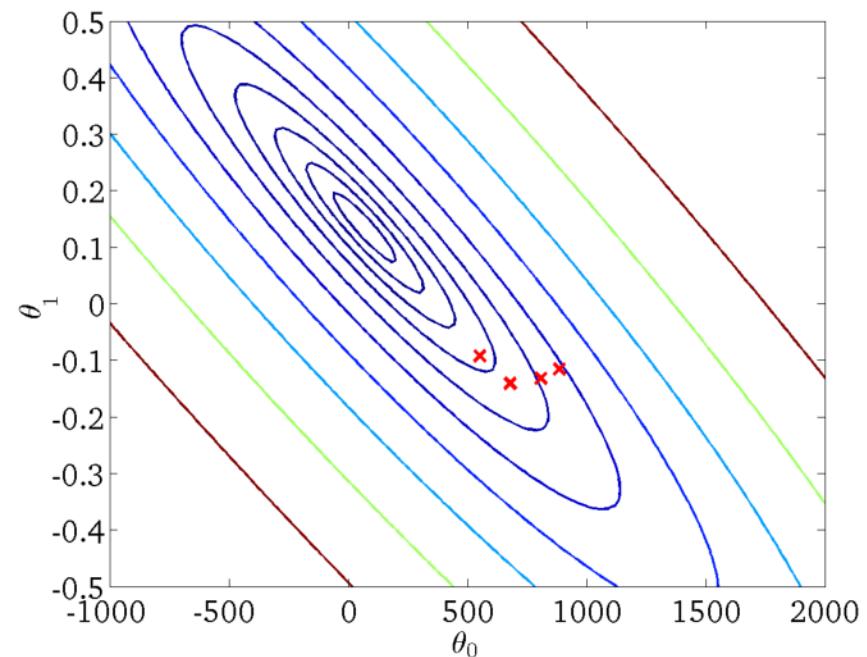
$$h_{\theta}(x)$$

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



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

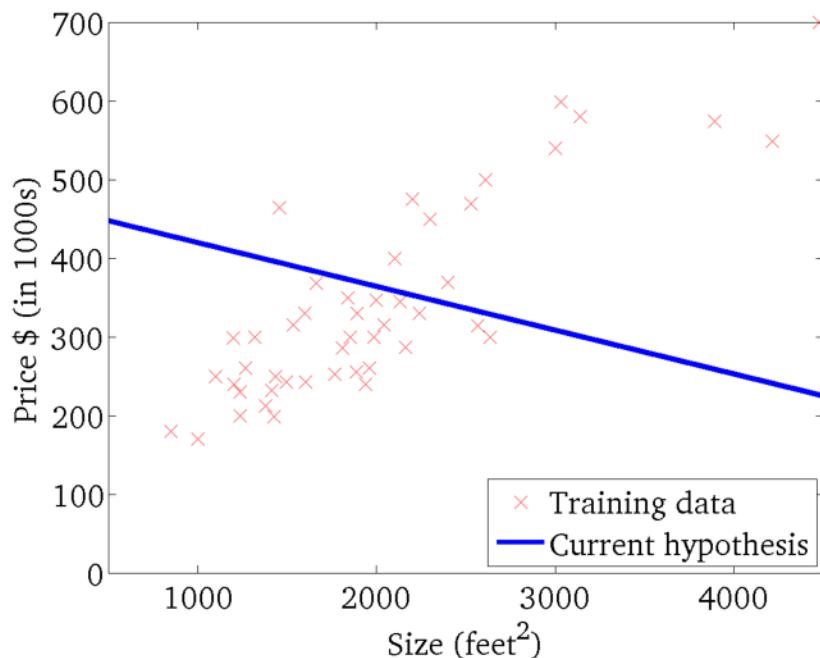
(function of the parameters  $\theta_0, \theta_1$ )



# Gradient Descent

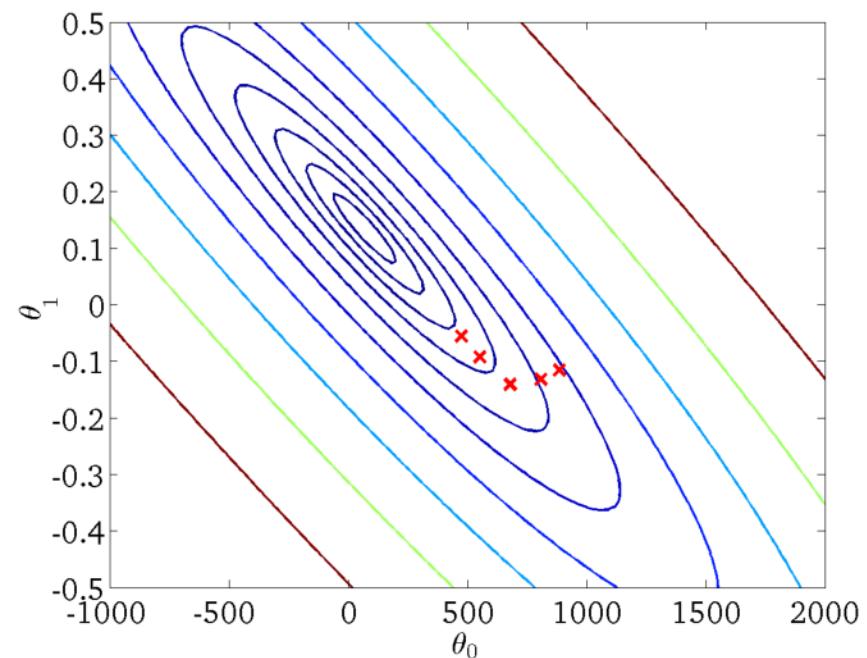
$$h_{\theta}(x)$$

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



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

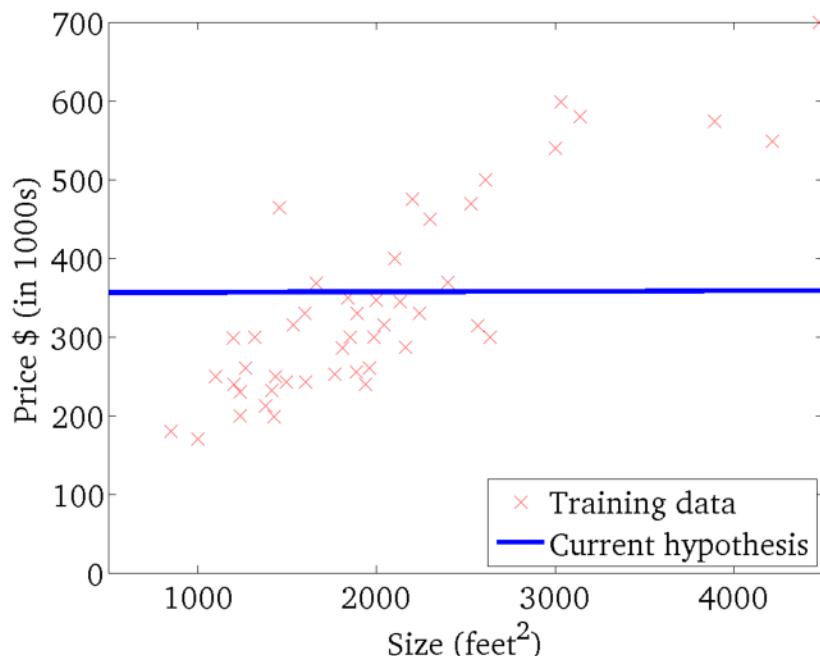
(function of the parameters  $\theta_0, \theta_1$ )



# Gradient Descent

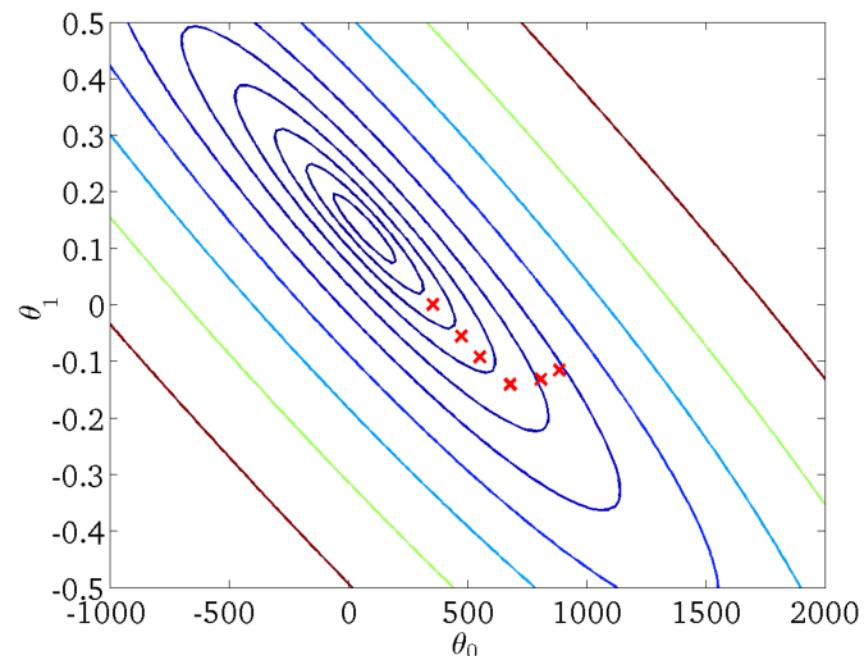
$$h_{\theta}(x)$$

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



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

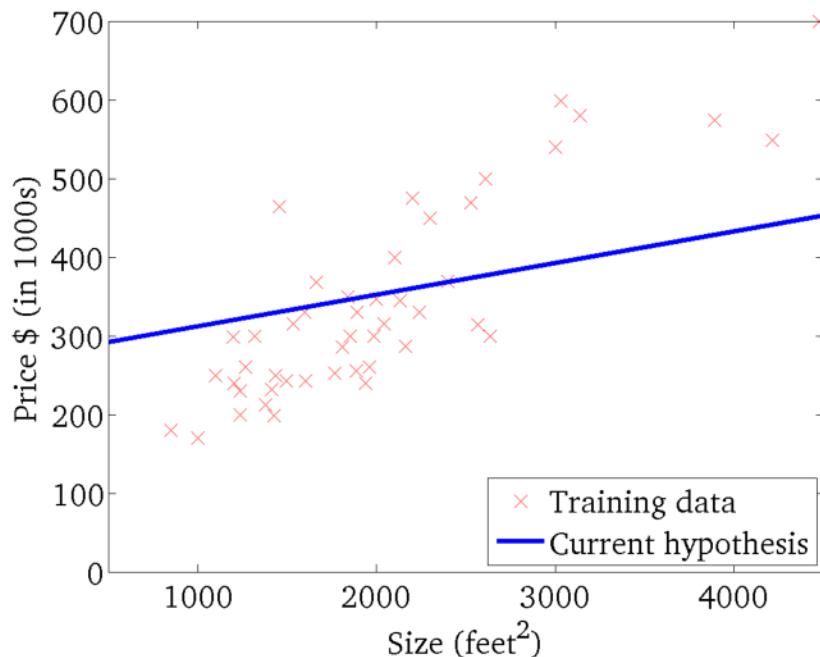
(function of the parameters  $\theta_0, \theta_1$ )



# Gradient Descent

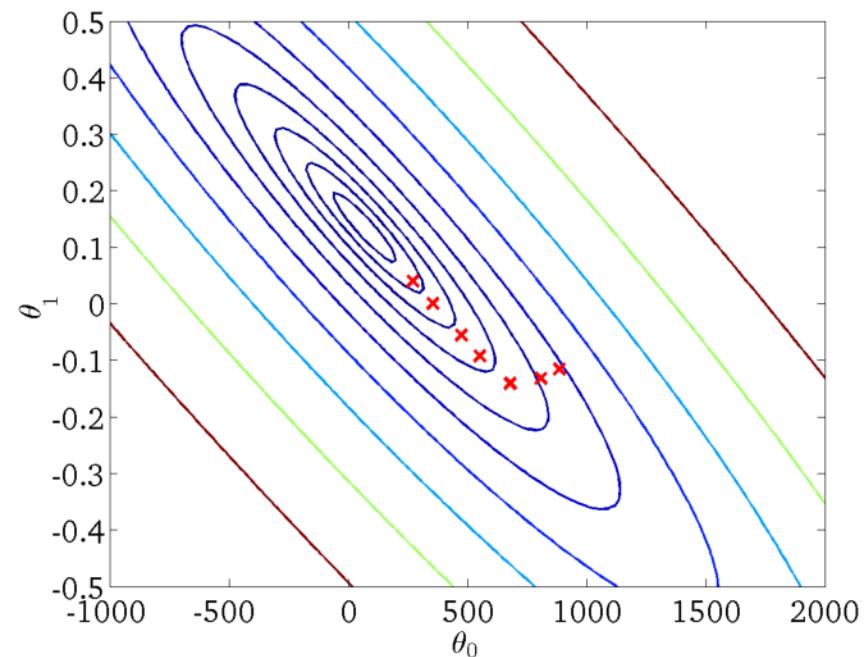
$$h_{\theta}(x)$$

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



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

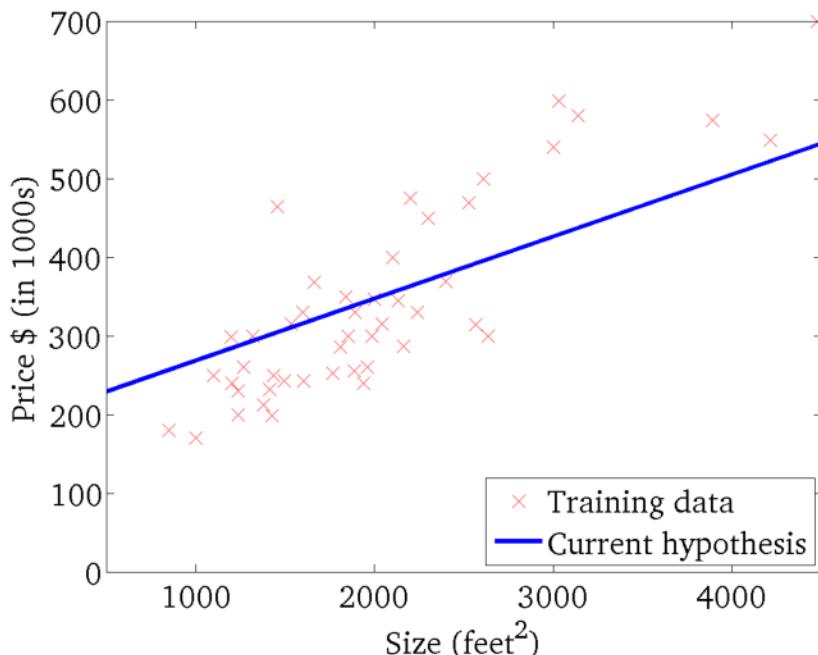
(function of the parameters  $\theta_0, \theta_1$ )



# Gradient Descent

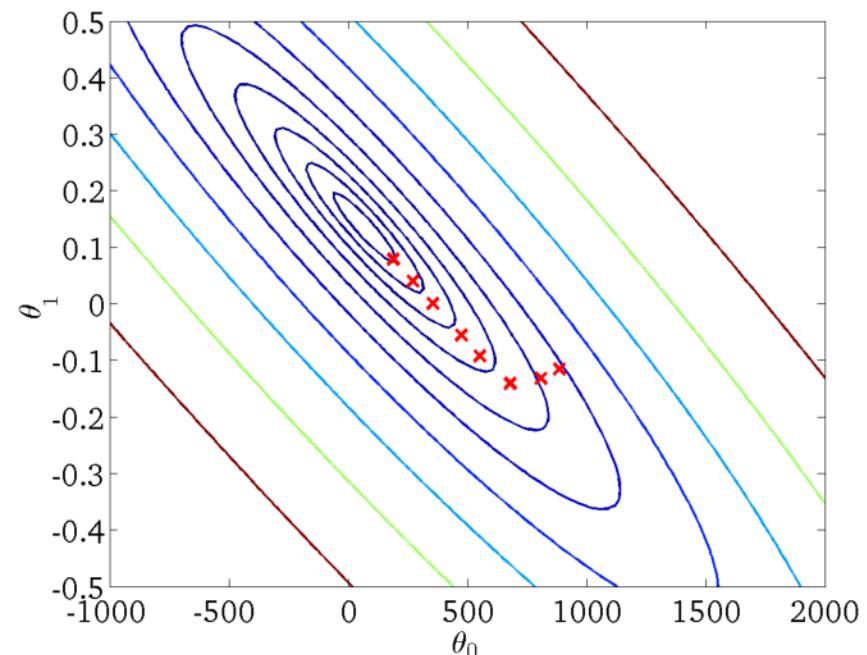
$$h_{\theta}(x)$$

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



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

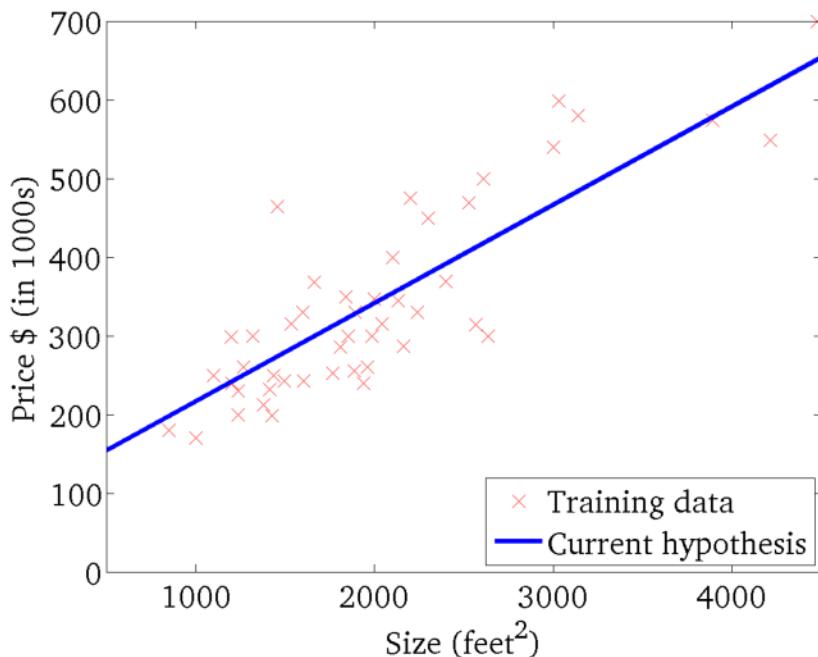
(function of the parameters  $\theta_0, \theta_1$ )



# Gradient Descent

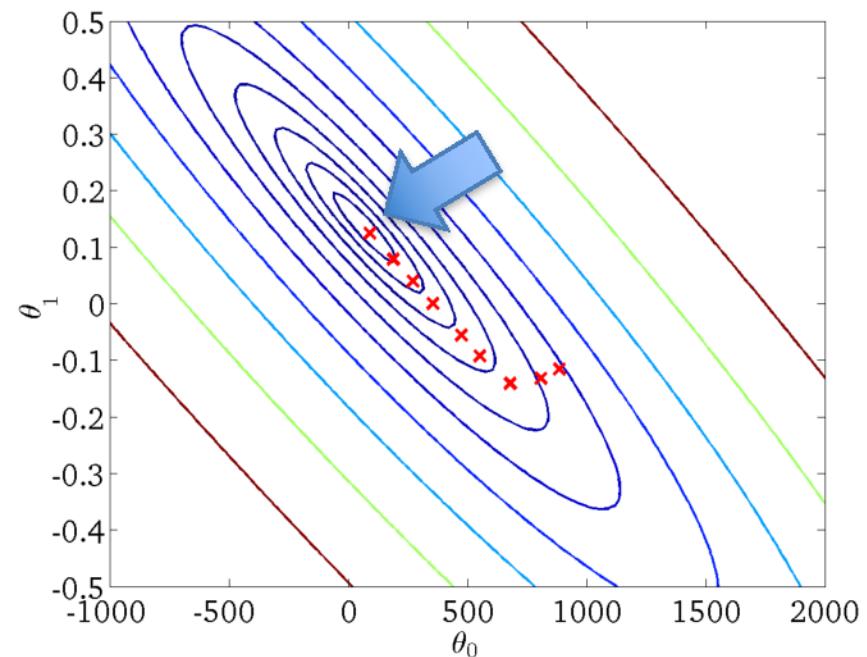
$$h_{\theta}(x)$$

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



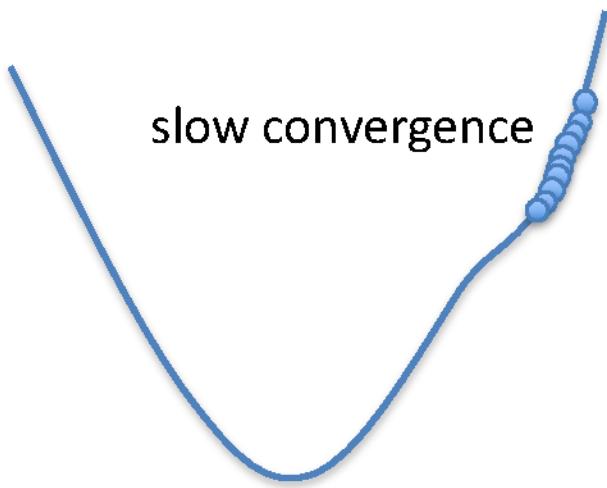
$$J(\theta_0, \theta_1)$$

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

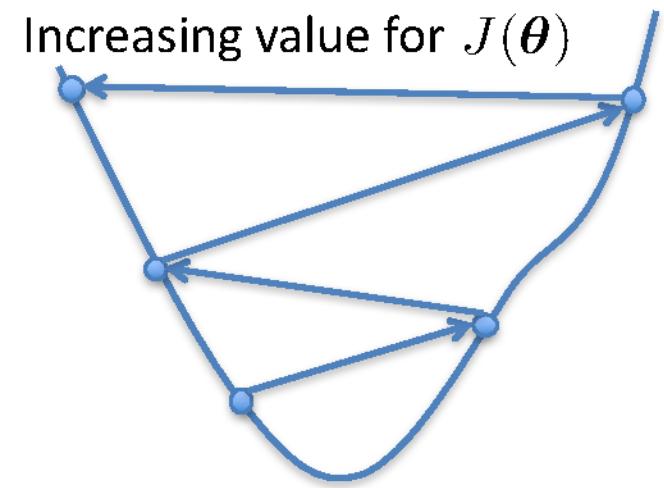


# Choosing $\alpha$

$\alpha$  too small



$\alpha$  too large



- 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  $\alpha$

---

# Learning Model Parameters – Closed Form Solution (using vectorization)

# Vectorization

- Benefits of vectorization
  - More compact equations
  - Faster code (using optimized matrix libraries)

- Consider our model:

$$h(\mathbf{x}) = \sum_{j=0}^d \theta_j x_j$$

- Let

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad \mathbf{x}^\top = [ 1 \quad x_1 \quad \dots \quad x_d ]$$

- Can write the model in vectorized form as  $h(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{x}$

# Vectorization

- Consider our model for n instances:

$$h(\mathbf{x}^{(i)}) = \sum_{j=0}^d \theta_j x_j^{(i)}$$

- Let

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & \dots & x_d^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix}$$
$$\mathbb{R}^{(d+1) \times 1} \qquad \qquad \qquad \mathbb{R}^{n \times (d+1)}$$

- Can write the model in vectorized form as  $h_{\theta}(\mathbf{x}) = \mathbf{X}\theta$

# Vectorization

- For the linear regression cost function:

$$\begin{aligned} J(\boldsymbol{\theta}) &= \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2 \\ &= \frac{1}{2n} \sum_{i=1}^n \left( \boldsymbol{\theta}^\top \mathbf{x}^{(i)} - y^{(i)} \right)^2 \\ &= \frac{1}{2n} \underbrace{(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^\top}_{\mathbb{R}^{1 \times n}} \underbrace{(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})}_{\mathbb{R}^{n \times 1}} \end{aligned}$$

$\mathbb{R}^{n \times (d+1)}$   
 $\mathbb{R}^{(d+1) \times 1}$

Let:

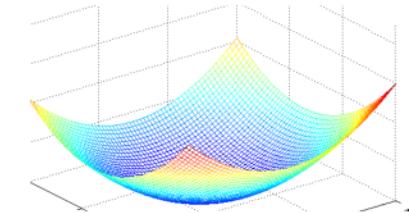
$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

# Closed Form Solution

- Instead of using GD, solve for optimal  $\theta$  analytically
  - Notice that the solution is when  $\frac{\partial}{\partial \theta} J(\theta) = 0$
- Derivation:

$$\begin{aligned}\mathcal{J}(\theta) &= \frac{1}{2n} (\mathbf{X}\theta - \mathbf{y})^\top (\mathbf{X}\theta - \mathbf{y}) \\ &\propto \theta^\top \mathbf{X}^\top \mathbf{X} \theta - \boxed{\mathbf{y}^\top \mathbf{X} \theta} - \boxed{\theta^\top \mathbf{X}^\top \mathbf{y}} + \mathbf{y}^\top \mathbf{y} \\ &\propto \theta^\top \mathbf{X}^\top \mathbf{X} \theta - 2\theta^\top \mathbf{X}^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y}\end{aligned}$$

1 x 1



Take derivative and set equal to 0, then solve for  $\theta$ :

$$\frac{\partial}{\partial \theta} (\theta^\top \mathbf{X}^\top \mathbf{X} \theta - 2\theta^\top \mathbf{X}^\top \mathbf{y} + \cancel{\mathbf{y}^\top \mathbf{y}}) = 0$$

$$(\mathbf{X}^\top \mathbf{X})\theta - \mathbf{X}^\top \mathbf{y} = 0$$

$$(\mathbf{X}^\top \mathbf{X})\theta = \mathbf{X}^\top \mathbf{y}$$

Closed Form Solution:

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

# Closed Form Solution

- Can obtain  $\theta$  by simply plugging  $X$  and  $y$  into

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

$$X = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & \dots & x_d^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

- If  $X^T X$  is not invertible (i.e., singular), may need to:
  - Use pseudo-inverse instead of the inverse
    - In python, `numpy.linalg.pinv(a)`
  - Remove redundant (not linearly independent) features
  - Remove extra features to ensure that  $d \leq n$

# Gradient Descent vs Closed Form

## Gradient Descent

- Requires multiple iterations
- Need to choose  $\alpha$
- Works well when  $n$  is large
- Can support incremental learning

## Closed Form Solution

- Non-iterative
- No need for  $\alpha$
- Slow if  $n$  is large
  - Computing  $(X^T X)^{-1}$  is roughly  $O(n^3)$

# Bayesian linear regression

# Bayesian analysis

---

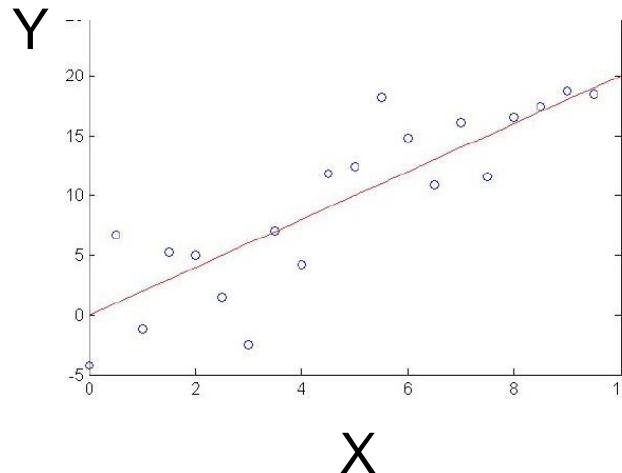
- Bayesian analysis will show that
  - under certain assumptions any learning algorithm that minimizes the squared error between the output hypothesis predictions and the training data will output a maximum likelihood hypothesis

# Maximum likelihood and least-squared error hypotheses

---

- A set of  $m$  training examples is provided, where the target value of each example is corrupted by random noise drawn according to a Normal probability distribution.
- Each training example is a pair of the form  $(x_i, d_i)$  where  $d_i = f(x_i) + e_i$ . Here  $f(x_i)$  is the noise-free value of the target function and  $e_i$  is a random variable representing the noise.
  - values of the  $e_i$  are drawn independently and that they are distributed according to a Normal distribution with zero mean

# Choose parameterized form for $P(Y|X; \theta)$



Assume Y is some deterministic  $f(X)$ , plus random noise

$$y = f(x) + \epsilon \quad \text{where } \epsilon \sim N(0, \sigma)$$

Therefore Y is a random variable that follows the distribution

$$p(y|x) = N(f(x), \sigma)$$

and the expected value of y for any given x is  $f(x)$

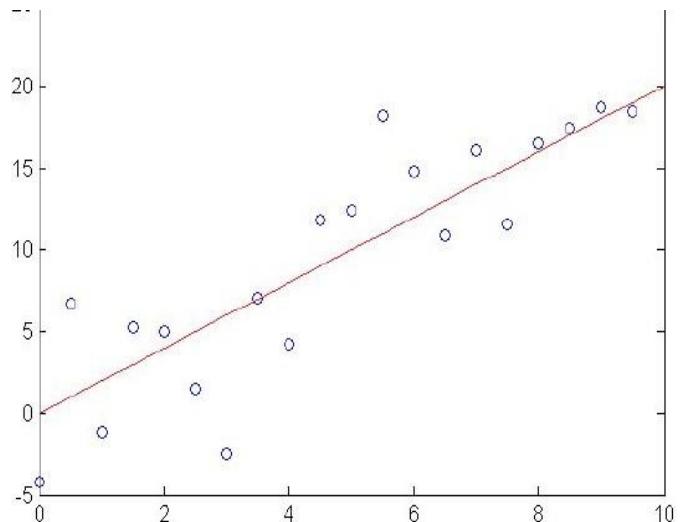
# Consider Linear Regression

$$p(y|x) = N(f(x), \sigma)$$

E.g., assume  $f(x)$  is linear function of  $x$

$$p(y|x) = N(w_0 + w_1x, \sigma)$$

$$E[y|x] = w_0 + w_1x$$



Notation: to make our parameters explicit, let's write

$$W = \langle w_0, w_1 \rangle$$

$$p(y|x; W) = N(w_0 + w_1x, \sigma)$$

# Training Linear Regression : Maximum Conditional Likelihood Estimate (MCLE)

$$p(y|x; W) = N(w_0 + w_1x, \sigma)$$

How can we learn  $W$  from the training data?

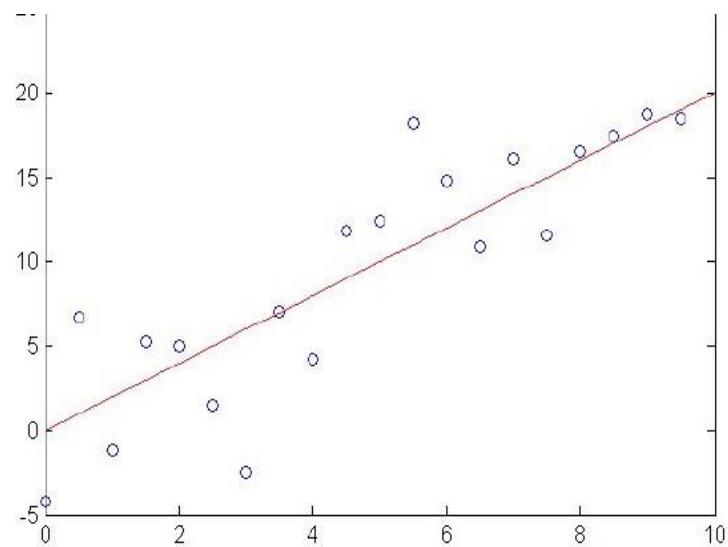
Learn Maximum Conditional Likelihood Estimate!

$$W_{MCLE} = \arg \max_W \prod_l p(y^l|x^l, W)$$

$$W_{MCLE} = \arg \max_W \sum_l \ln p(y^l|x^l, W)$$

where

$$p(y|x; W) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}(\frac{y-f(x;W)}{\sigma})^2}$$



# Training Linear Regression: MCLE

Learn Maximum Conditional Likelihood Estimate

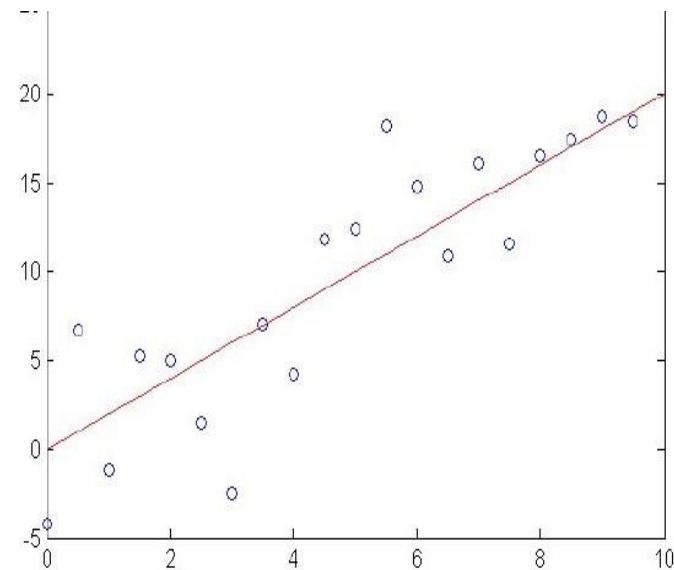
$$W_{MCLE} = \arg \max_W \sum_l \ln p(y^l | x^l, W)$$

where

$$p(y|x; W) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}(\frac{y-f(x;W)}{\sigma})^2}$$

so:

$$W_{MCLE} = \arg \min_W \sum_l (y - f(x; W))^2$$



# Training Linear Regression: MCLE, MLE



$$h_{ML} = \operatorname{argmax}_{h \in H} \sum_{i=1}^m \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} (d_i - h(x_i))^2$$

The first term in this expression is a constant independent of  $h$ , and can therefore be discarded, yielding

$$h_{ML} = \operatorname{argmax}_{h \in H} \sum_{i=1}^m -\frac{1}{2\sigma^2} (d_i - h(x_i))^2$$

Maximizing this negative quantity is equivalent to minimizing the corresponding positive quantity.

$$h_{ML} = \operatorname{argmin}_{h \in H} \sum_{i=1}^m \frac{1}{2\sigma^2} (d_i - h(x_i))^2$$

Finally, we can again discard constants that are independent of  $h$ .

$$h_{ML} = \operatorname{argmin}_{h \in H} \sum_{i=1}^m (d_i - h(x_i))^2 \tag{6.6}$$

# Training Linear Regression

Maximum Conditional Likelihood Estimate is equivalent to minimizing the squared error loss

$$W_{MCLE} = \arg \min_W \sum_l (y - f(x; W))^2$$

Can we derive gradient descent rule for training?

$$\begin{aligned}\frac{\partial \sum_l (y - f(x; W))^2}{\partial w_i} &= \sum_l 2(y - f(x; W)) \frac{\partial (y - f(x; W))}{\partial w_i} \\ &= \sum_l -2(y - f(x; W)) \frac{\partial f(x; W)}{\partial w_i}\end{aligned}$$

# Regression – What you should know

---

Under general assumption

$$p(y|x; W) = N(f(x; W), \sigma)$$

- We can use gradient descent as a general learning algorithm
  - as long as our objective function is differentiable wrt W
  - though we might learn local optima
- Almost nothing we said here required that  $f(x)$  be linear in  $x$

# Extending Linear Regression to More Complex Models

- The inputs  $\mathbf{X}$  for linear regression can be:
  - Original quantitative inputs
  - Transformation of quantitative inputs
    - e.g. log, exp, square root, square, etc.
  - Polynomial transformation
    - example:  $y = \beta_0 + \beta_1 \cdot x + \beta_2 \cdot x^2 + \beta_3 \cdot x^3$
  - Basis expansions
  - Dummy coding of categorical inputs
  - Interactions between variables
    - example:  $x_3 = x_1 \cdot x_2$

This allows use of linear regression techniques to fit non-linear datasets.

# Linear Basis Function Models

# Linear Basis Function

---

- Simplest linear model for regression is one that involves a linear combination of the input variables

$$y(x, w) = w_0 + w_1 x_1 + \dots + w_D x_D$$

- Extend the class of models by considering linear combinations of fixed nonlinear functions of the input variables, of the form

$$y(x, w) = w_0 + \sum_{j=1}^{M-1} w_j \varphi_j(x)$$

- where  $\varphi_j(x)$  are known as basis functions.
- By denoting the maximum value of the index j by M - 1, the total number of parameters in this model will be M.

# Linear Basis Function

---

- Convenient to define an additional dummy ‘basis function’  $\phi_0(x)=1$ . So,

$$y(x, w) = \sum_{j=1}^{M-1} w_j \varphi_j(x) = \mathbf{w}^\top \boldsymbol{\Phi}_j(x)$$

where  $w = (w_0, \dots, w_{M-1})^T$  and  $\boldsymbol{\Phi} = (\phi_0, \phi_1, \dots, \phi_n)$

- If the original variables comprise the vector  $x$ , then the features can be expressed in terms of the basis functions  $\{\phi_j(x)\}$

# Linear Basis Function Models

- Generally,

$$h_{\theta}(\mathbf{x}) = \sum_{j=0}^d \theta_j \phi_j(\mathbf{x})$$

basis function

- Typically,  $\phi_0(\mathbf{x}) = 1$  so that  $\theta_0$  acts as a bias
- In the simplest case, we use linear basis functions :

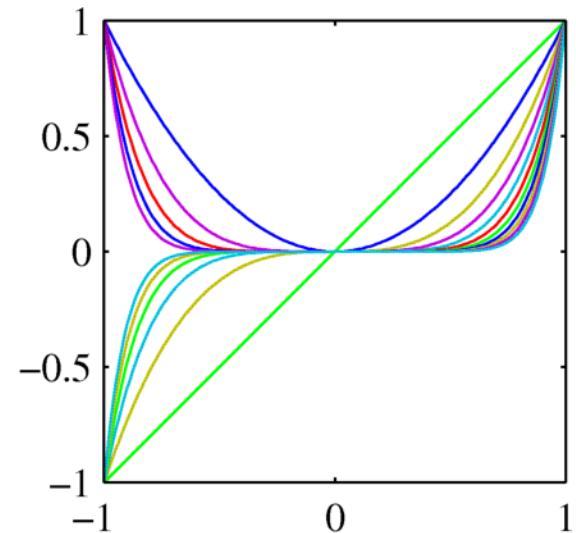
$$\phi_j(\mathbf{x}) = x_j$$

# Linear Basis Function Models

- Polynomial basis functions:

$$\phi_j(x) = x^j$$

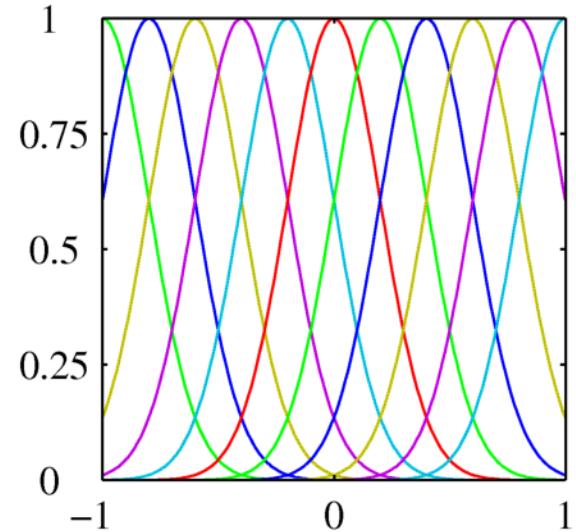
- These are global; a small change in  $x$  affects all basis functions



- Gaussian basis functions:

$$\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\}$$

- These are local; a small change in  $x$  only affect nearby basis functions.  $\mu_j$  and  $s$  control location and scale (width).



# Linear Basis Function Models

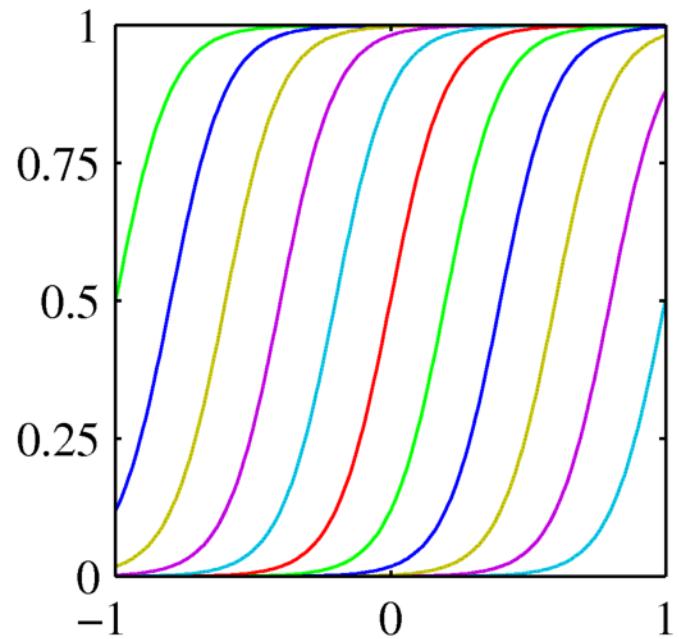
- Sigmoidal basis functions:

$$\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right)$$

where

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

- These are also local; a small change in  $x$  only affects nearby basis functions.  $\mu_j$  and  $s$  control location and scale (slope).



# Linear Basis Function Models

- Basic Linear Model:

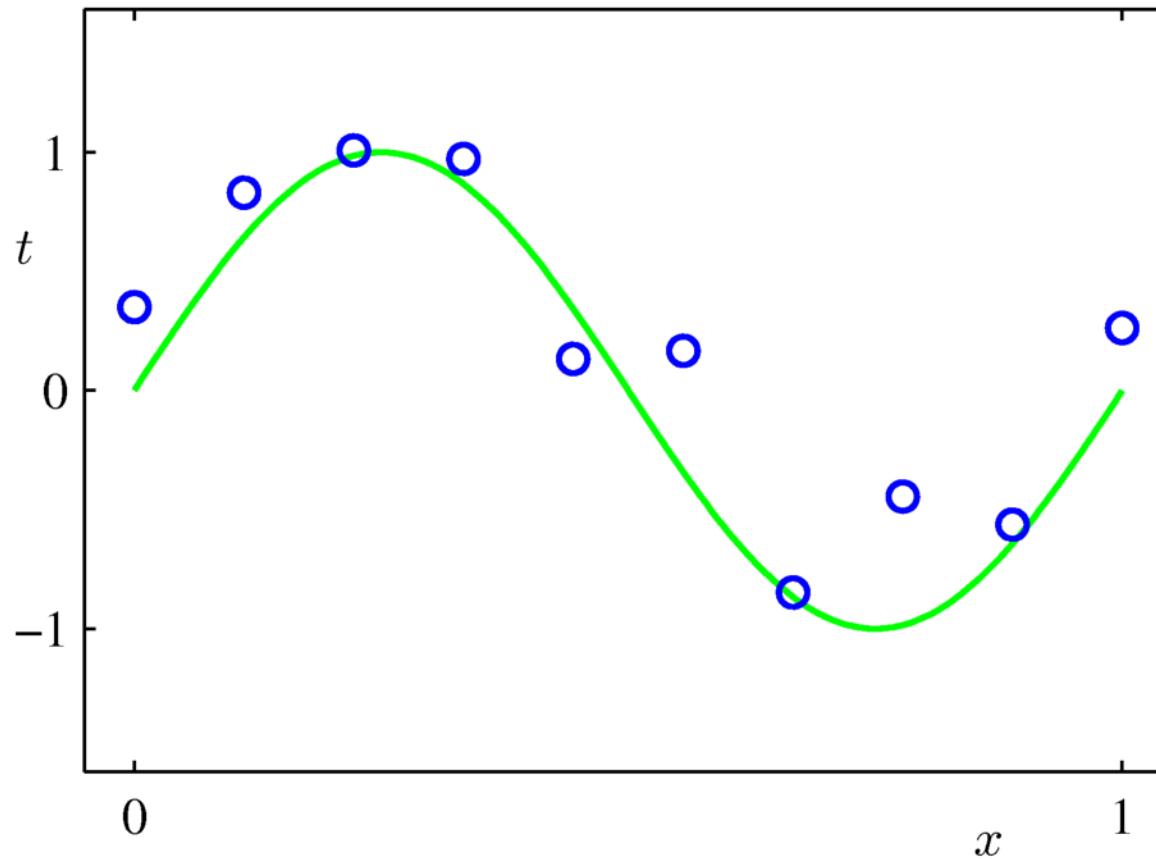
$$h_{\theta}(\mathbf{x}) = \sum_{j=0}^d \theta_j x_j$$

- Generalized Linear Model:

$$h_{\theta}(\mathbf{x}) = \sum_{j=0}^d \theta_j \phi_j(\mathbf{x})$$

- Once we have replaced the data by the outputs of the basis functions, fitting the generalized model is exactly the same problem as fitting the basic model
  - Unless we use the kernel trick – more on that when we cover support vector machines
  - Therefore, there is no point in cluttering the math with basis functions

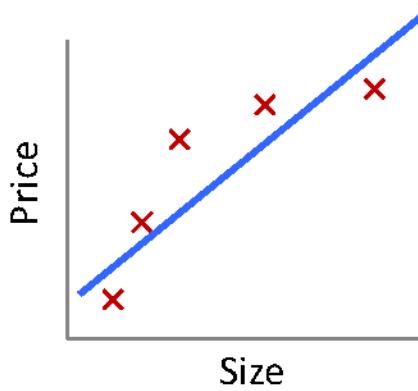
# Example of Fitting a Polynomial Curve with a Linear Model



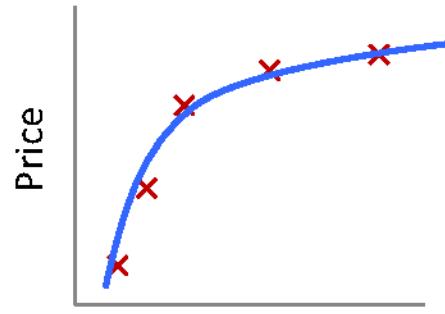
$$y = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_p x^p = \sum_{j=0}^p \theta_j x^j$$

# Bias-Variance Decomposition

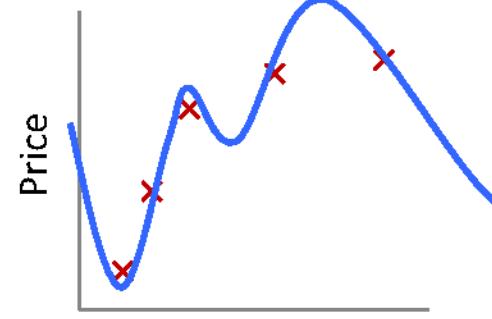
# Quality of Fit



Underfitting  
(high bias)



Correct fit



Overfitting  
(high variance)

## Overfitting:

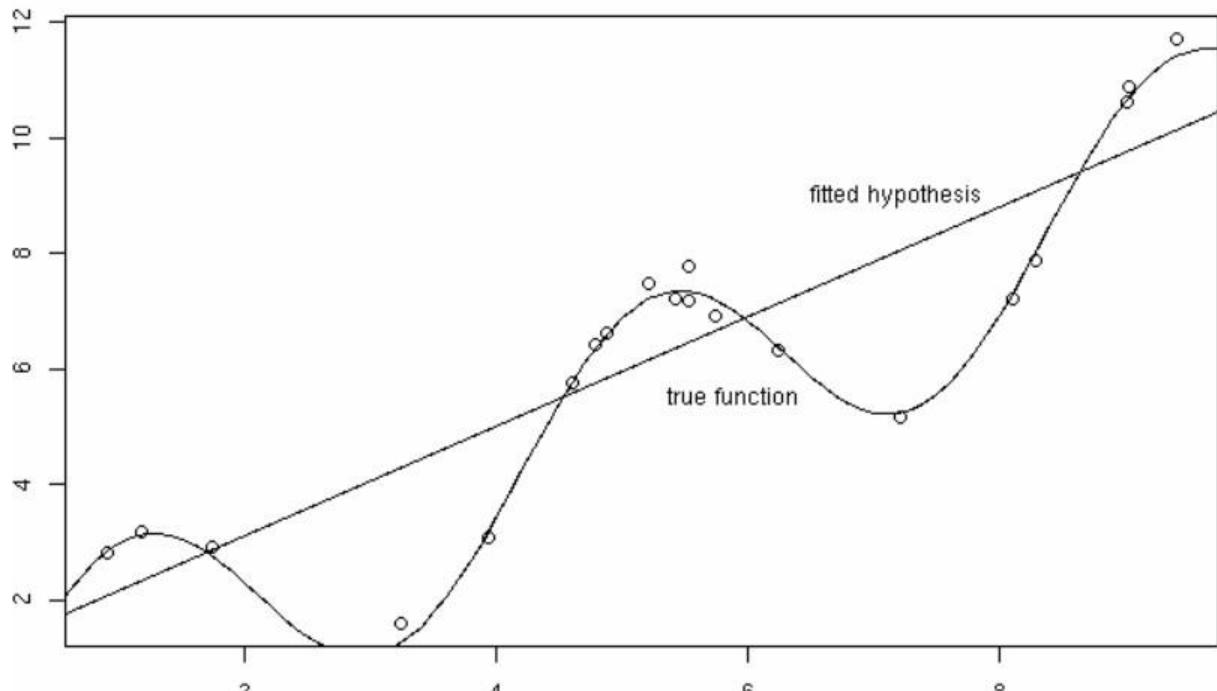
- The learned hypothesis may fit the training set very well ( $J(\theta) \approx 0$ )
- ...but fails to generalize to new examples

# Bias-Variance Tradeoff

- **Bias:** difference between  
**what you expect to learn** and **truth**
  - Measures how well you expect to represent true solution
  - Decreases with more complex model
  
- **Variance:** difference between  
**what you expect to learn** and  
**what you learn from a particular dataset**
  - Measures how sensitive learner is to specific dataset
  - Increases with more complex model

# Example

Tom Dietterich, Oregon St



True Function :  
 $y=f(x)$

$$y = x + 2 \sin(1.5x) + N(0,0.2)$$

# Bias – Variance decomposition of error

$$E_{D,\varepsilon} \left\{ (f(x) + \varepsilon - h_D(x))^2 \right\}$$

dataset and noise      true function      noise      learned from D

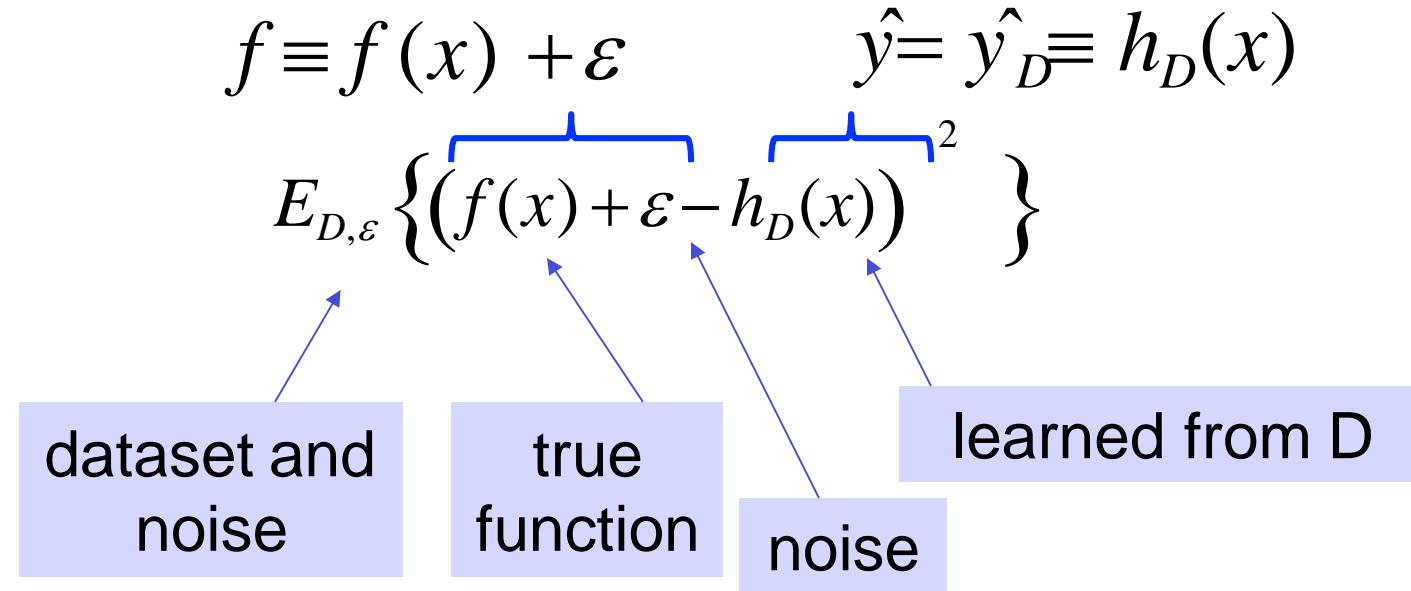
Fix test case  $x$ , then do this experiment:

1. Draw size  $n$  sample  $D = (x_1, y_1), \dots, (x_n, y_n)$
2. Train linear regressor  $h_D$  using  $D$
3. Draw one test example  $(x, f(x) + \varepsilon)$
4. Measure squared error of  $h_D$  on that example  $x$
5. What's the expected error?

# Bias – Variance decomposition of error

Notation - to simplify this

$$f \equiv f(x) + \varepsilon \quad \hat{y} = \hat{y}_D \equiv h_D(x)$$

$$E_{D,\varepsilon} \left\{ (f(x) + \varepsilon - h_D(x))^2 \right\}$$


dataset and noise      true function      noise      learned from D

$$h \equiv E_D\{h_D(x)\}$$

long-term expectation of learner's prediction  
on this  $x$  averaged over many data sets  $D$

# Bias – Variance decomposition of error

$$\begin{aligned}
 & E_{D,\varepsilon} \left\{ (\hat{y}_D - y)^2 \right\} & h = E_D \left\{ \hat{y}_D \right\} \\
 & = E \left\{ ([f-h] + [h-\hat{y}_D])^2 \right\} & \hat{y}_D = \hat{y}_D(x) \\
 & = E \left\{ [f-h]^2 + [h-\hat{y}_D]^2 + 2[f-h][h-\hat{y}_D] \right\} & f \equiv f(x) + \varepsilon \\
 & = E \left\{ [f-h]^2 + [h-\hat{y}_D]^2 + 2[fh - f\hat{y}_D - h^2 + h\hat{y}_D] \right\} & \\
 & = E[(f-h)^2] + E[(h-\hat{y}_D)^2] + 2 \left( \cancel{E[fh]} - \cancel{E[f\hat{y}_D]} - \cancel{E[h^2]} + \cancel{E[h\hat{y}_D]} \right)
 \end{aligned}$$

# Bias – Variance decomposition of error

$$\begin{aligned}
 & E_{D,\varepsilon} \left\{ (\hat{y}_D - y)^2 \right\} & h = E_D \left\{ \hat{y}_D \right\} \\
 & = E \left\{ ([f-h] + [h-\hat{y}_D])^2 \right\} & \hat{y}_D = \hat{y}_D(x) \\
 & = E \left\{ [f-h]^2 + [h-\hat{y}_D]^2 + 2[f-h][h-\hat{y}_D] \right\} & f \equiv f(x) + \varepsilon \\
 & = E \left\{ [f-h]^2 + [h-\hat{y}_D]^2 + 2[fh - f\hat{y}_D - h^2 + h\hat{y}_D] \right\} & \\
 & = E[(f-h)^2] + E[(h-\hat{y}_D)^2] + 2 \left( \cancel{E[fh]} - \cancel{E[f\hat{y}_D]} - \cancel{E[h^2]} + \cancel{E[h\hat{y}_D]} \right)
 \end{aligned}$$

# Bias – Variance decomposition of error

$$h = E_D \{ \hat{y}_D \} \quad \hat{y} = \hat{y}_D$$

$$f \equiv f(x) + \varepsilon$$

$$\begin{aligned}
 & 2(E[fh] - E[fy] - E[h^2] + E[hy]) \\
 &= E[f]E[h] - E[f]E[y] - E[h]E[h] + E[h]E[y] \\
 &= E[f]E[E[y]] - E[f]E[y] - E[h]E[E[y]] + E[h]E[y] \\
 &= E[f]E[y] - E[f]E[y] - E[h]E[y] + E[h]E[y] \\
 &= 0
 \end{aligned}$$

# Bias – Variance decomposition of error

$$\begin{aligned}
 & E_{D,\varepsilon} \{(f - \hat{y})^2\} \\
 &= E \left\{ ([f - h] + [h - \hat{y}])^2 \right\} \\
 &= E \left\{ [f - h]^2 + [h - \hat{y}]^2 + 2[f - h][h - \hat{y}] \right\} \\
 &= E[(f - h)^2] + E[(h - \hat{y})^2]
 \end{aligned}$$

$$\begin{aligned}
 h &\equiv E_D\{h_D(x)\} \\
 \hat{y} &\equiv \hat{y}_D \equiv h_D(x) \\
 f &\equiv f(x) + \varepsilon
 \end{aligned}$$

**BIAS<sup>2</sup>**

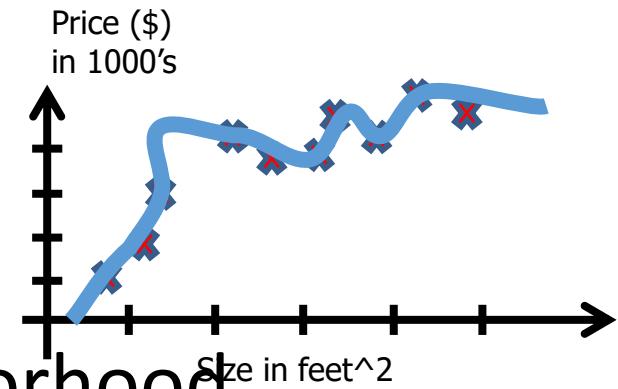
Squared difference between best possible prediction for  $x$ ,  $f(x)$ , and our “long-term” expectation for what the learner will do if we averaged over many datasets  $D$ ,  $E_D[h_D(x)]$

**VARIANCE**

Squared difference btwn our long-term expectation for the learners performance,  $E_D[h_D(x)]$ , and what we expect in a representative run on a dataset D ( $\hat{y}$ )

# Addressing overfitting

- $x_1$  = size of house
- $x_2$  = no. of bedrooms
- $x_3$  = no. of floors
- $x_4$  = age of house
- $x_5$  = average income in neighborhood
- $x_6$  = kitchen size
- :
- $x_{100}$



Slide credit: Andrew Ng

# Addressing overfitting

---

- **1. Reduce number of features.**
  - Manually select which features to keep.
  - Model selection algorithm
  
- **2. Regularization.**
  - Keep all the features, but reduce magnitude/values of parameters  $\theta_j$ .
  - Works well when we have a lot of features, each of which contributes a bit to predicting  $y$ .

Slide credit: Andrew Ng

# 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)

# Regularization

- Linear regression objective function

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$

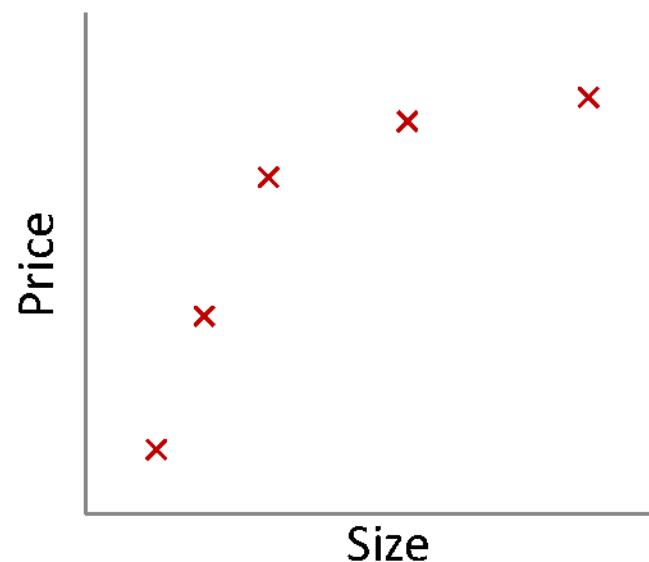

The equation represents the cost function for linear regression. It consists of two parts: a sum of squared differences between the predicted values  $h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})$  and the actual values  $y^{(i)}$  for all data points  $i$ , scaled by  $\frac{1}{2n}$ ; and a regularization term, which is a sum of the squares of all model parameters  $\theta_j$  for  $j$  from 1 to  $d$ , scaled by  $\frac{\lambda}{2}$ . The regularization parameter  $\lambda$  controls the trade-off between fitting the data well and keeping the model parameters small.

- $\lambda$  is the regularization parameter ( $\lambda \geq 0$ )
- No regularization on  $\theta_0$ !

# Understanding Regularization

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$

- What happens if we set  $\lambda$  to be huge (e.g.,  $10^{10}$ )?

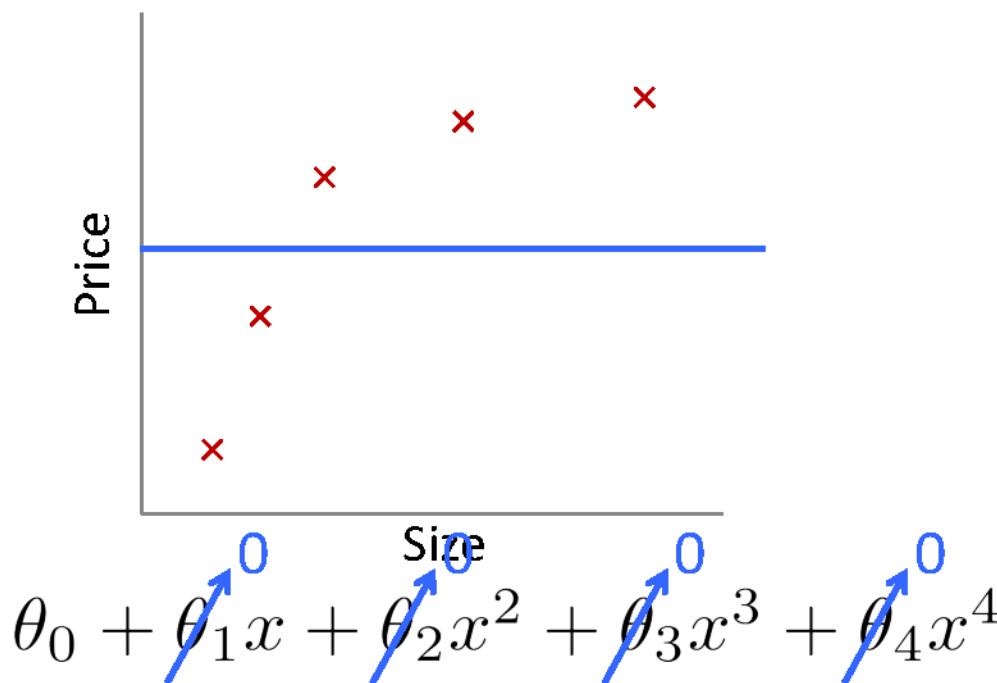


$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

# Understanding Regularization

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$

- What happens if we set  $\lambda$  to be huge (e.g.,  $10^{10}$ )?



# Regularized Linear Regression

- Cost Function

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$

- Fit by solving  $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

- Gradient update:

$$\frac{\partial}{\partial \theta_0} J(\boldsymbol{\theta})$$

$$\theta_0 \leftarrow \theta_0 - \alpha \frac{1}{n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)$$

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$$

$$\theta_j \leftarrow \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)} - \lambda \theta_j$$

# Regularized Linear Regression

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$

$$\theta_0 \leftarrow \theta_0 - \alpha \frac{1}{n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)$$

$$\theta_j \leftarrow \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)} - \lambda \theta_j$$

- We can rewrite the gradient step as:

$$\theta_j \leftarrow \theta_j (1 - \alpha \lambda) - \alpha \frac{1}{n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)}$$

---

## Funny but simple explanation about bias and variance

- <https://www.youtube.com/watch?v=EuBBz3bl-aA>

# In our next session

We will cover:

Tom Mitchell - Chapter 3

- Decision Tree
- Handling overfitting
- Continuous values
- Missing Values
- Random Forest