

Machine Learning

CS 539

Worcester Polytechnic Institute

Department of Computer Science

Instructor: Prof. Kyumin Lee

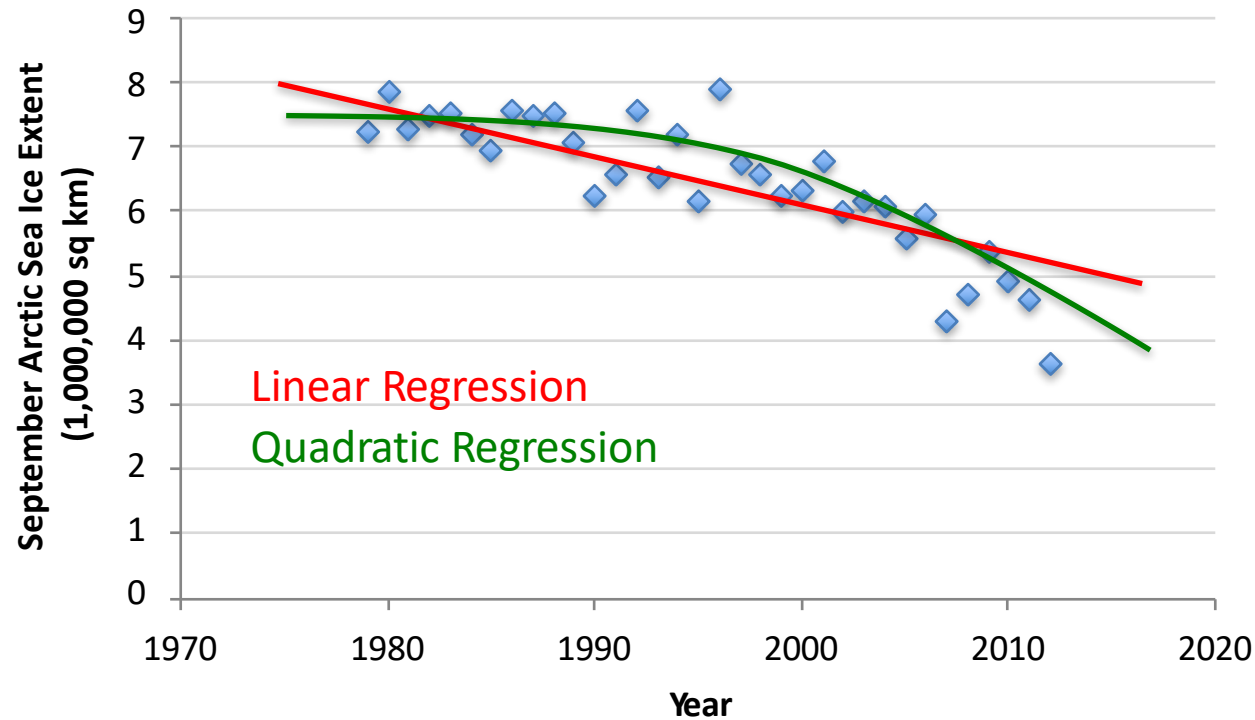
Project Teams

- Will Buchta, Sarah Semy, Deena Selitsk, and Humza Qureshi
- Kai Davidson, Sarah Kogan, Dominic Ranelli
- Bryce Lukens, Harrison Taylor, Gabriel Barbosa
- Axel Luca, Preet Patel, Xiaojun Liu
- Codey Battista, Avinash Bissoondial, Nick Chantre, Paul Crann
- John DeLeo, Jacob Donahue, Joe Molder, Iris Nycz
- Srisaranya Pujari, Kangjian Wu, Matthew Yeo, Yang Yi

Linear Regression

Regression

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



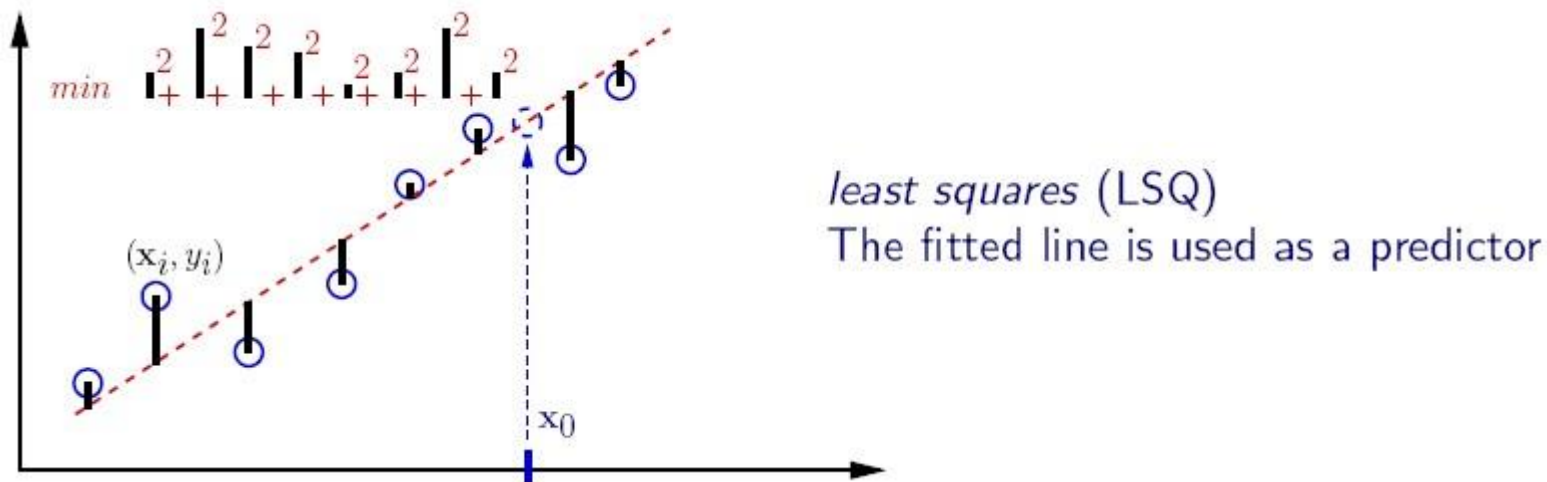
Linear Regression

- Hypothesis:

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

Assume $x_0 = 1$

- Fit model by minimizing sum of squared errors

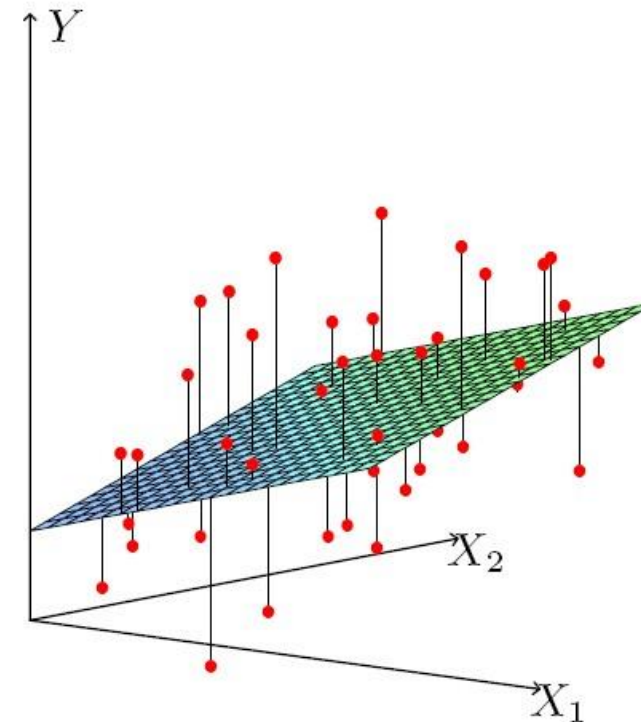
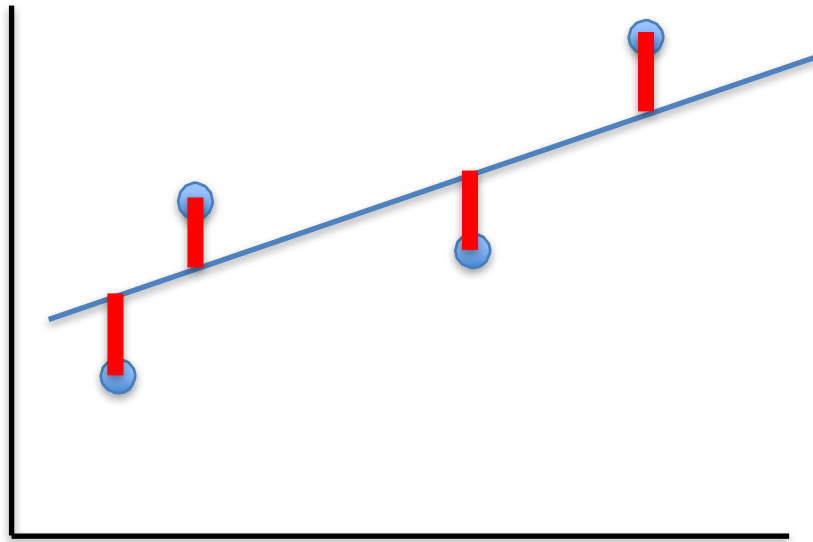


Least Squares Linear Regression

- Cost Function

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

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



Intuition Behind Cost Function

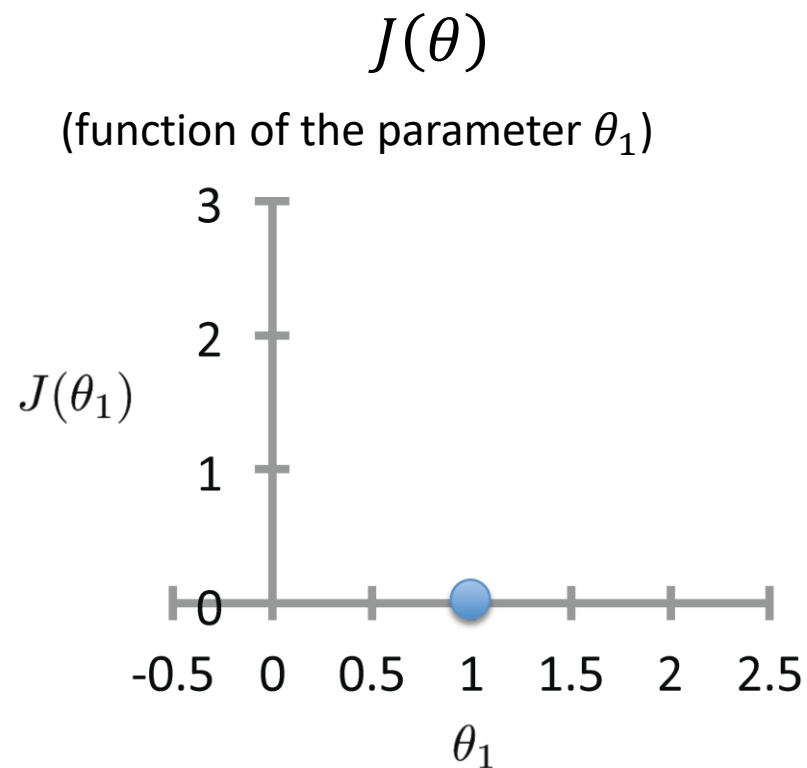
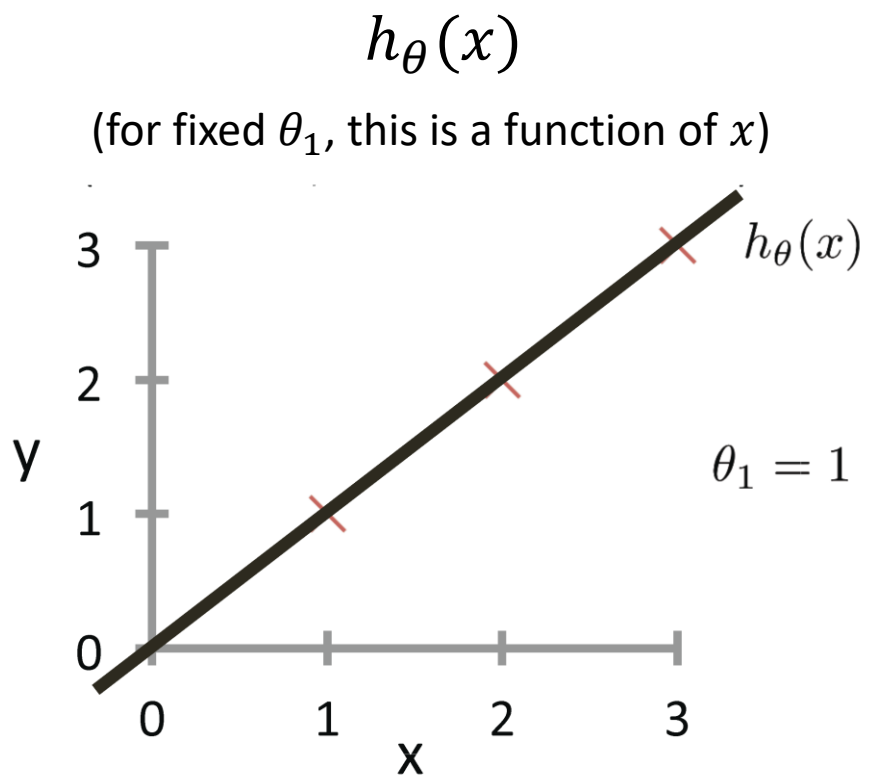
$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

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

Intuition Behind Cost Function

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

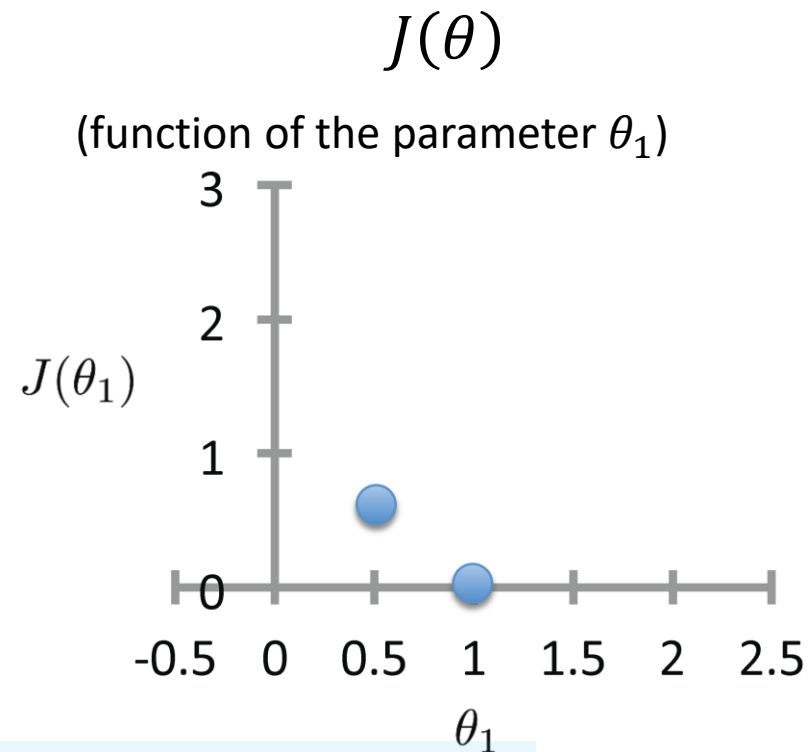
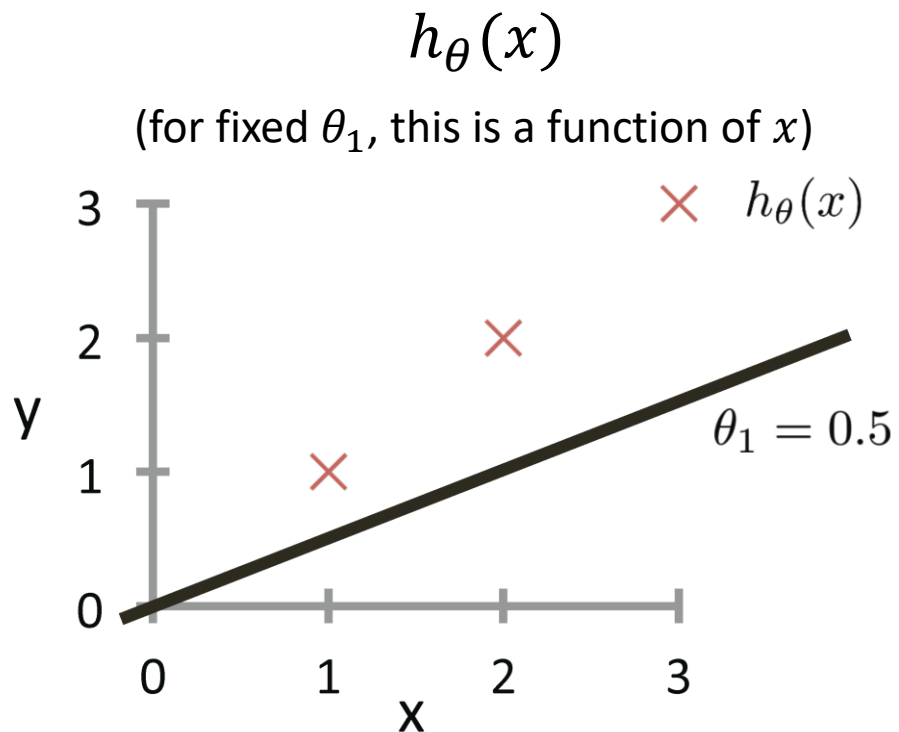
For insight on $J()$, let's assume $x \in \mathbb{R}$ so $\theta = [\theta_0, \theta_1] \rightarrow \theta_0 = 0$



Intuition Behind Cost Function

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

For insight on $J()$, let's assume $x \in \mathbb{R}$ so $\theta = [\theta_0, \theta_1] \rightarrow \theta_0 = 0$

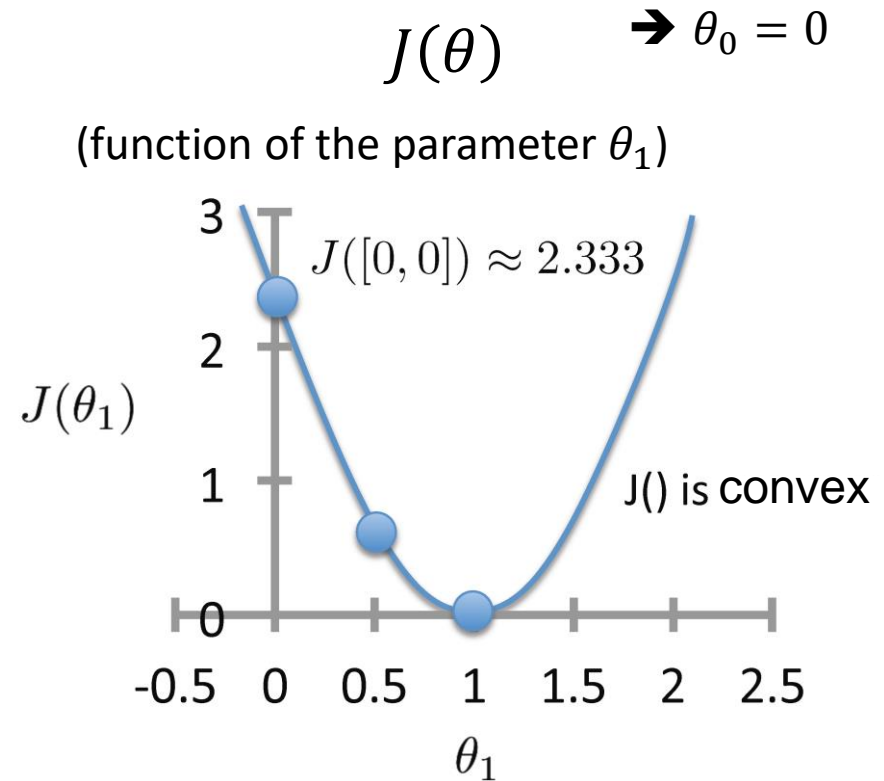
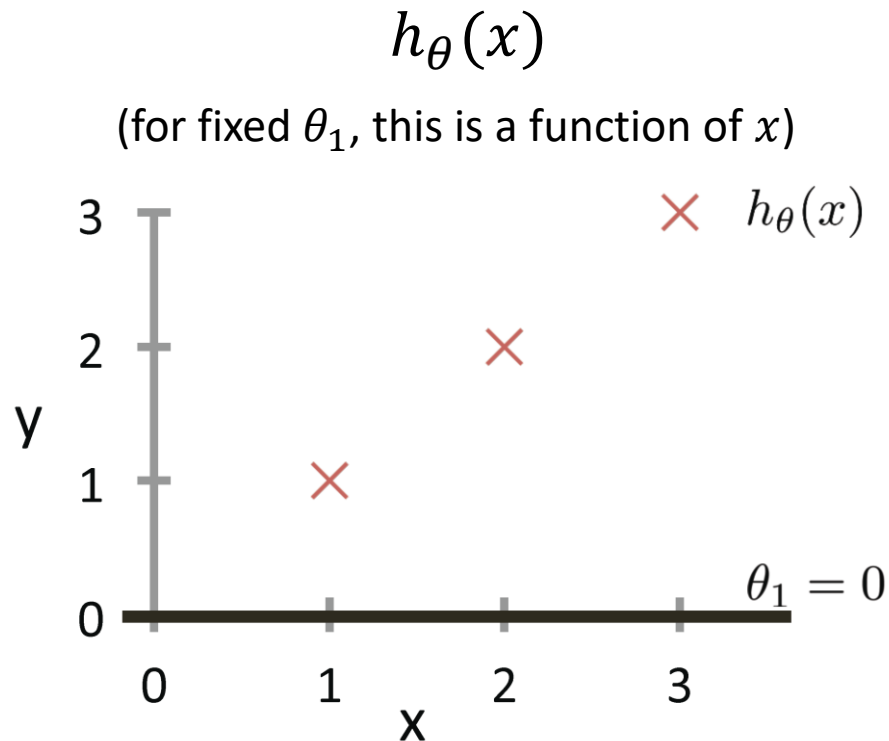


$$J([0, 0.5]) = \frac{1}{2 \times 3} [(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2] \approx 0.58$$

Intuition Behind Cost Function

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

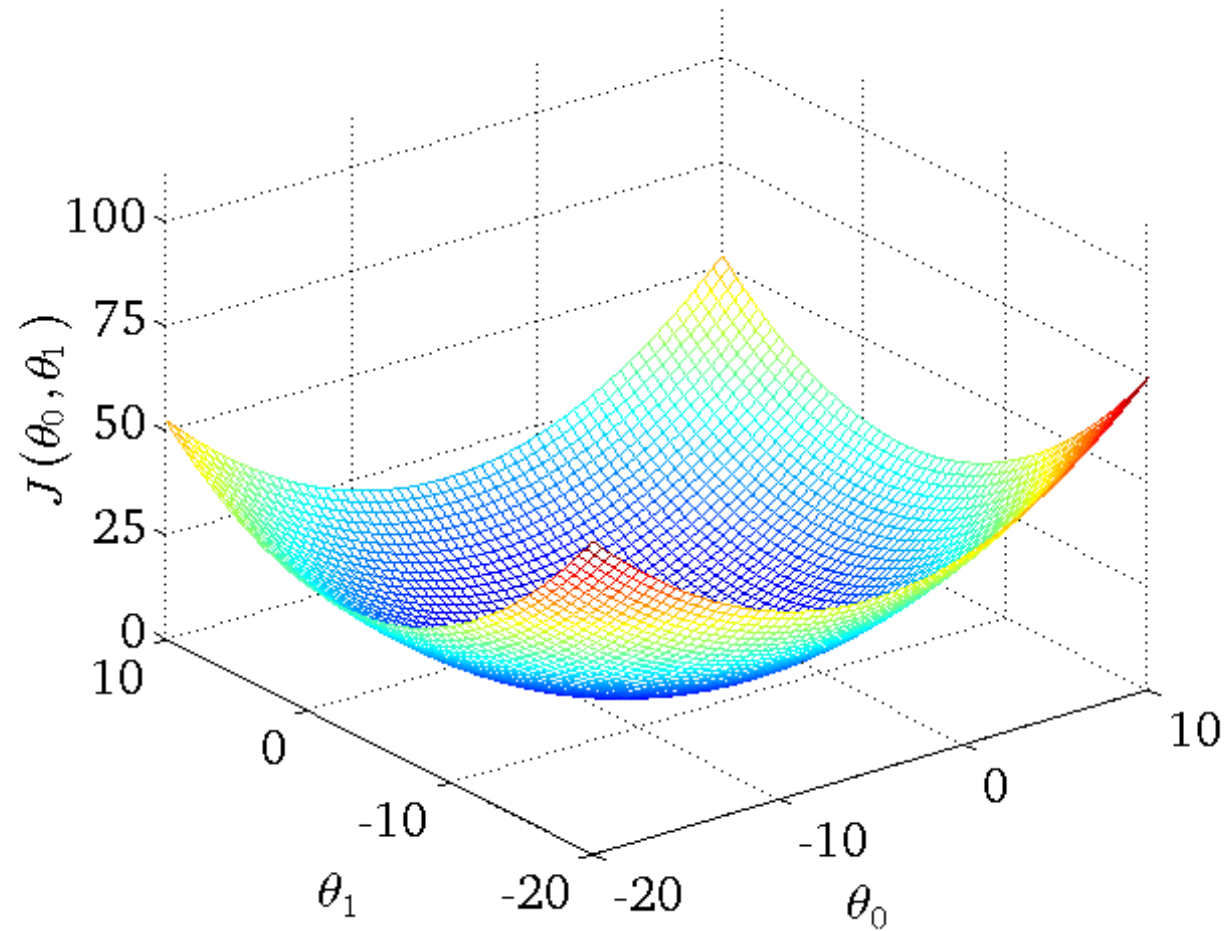
For insight on $J()$, let's assume $x \in \mathbb{R}$ so $\theta = [\theta_0, \theta_1]$



<http://mathworld.wolfram.com/ConvexFunction.html>

<https://www.desmos.com/calculator/kreo2ssqj8>

Intuition Behind Cost Function (3-D surface plot)



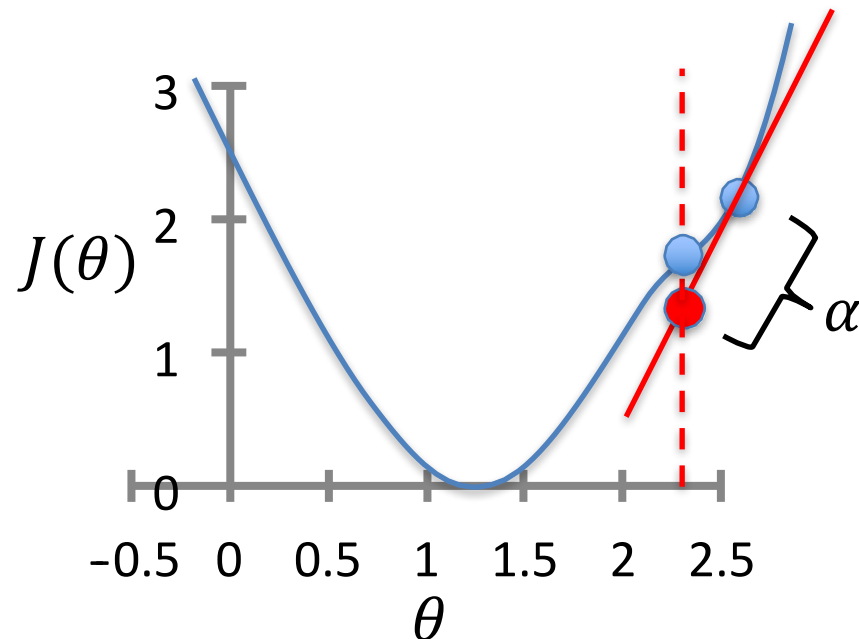
Gradient Descent

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

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

To move from equation [1.1] to [1.2], we need to apply two basic derivative rules:

Scalar multiple rule: $\frac{d}{dx}(\alpha u) = \alpha \frac{du}{dx}$

Sum rule: $\frac{d}{dx} \sum u = \sum \frac{du}{dx}$

Moving from [1.2] to [1.3], we apply both the power rule and the chain rule:

Power rule: $\frac{d}{dx} u^n = n u^{n-1} \frac{du}{dx}$

Chain rule: $\frac{d}{dx} f(g(x)) = f'(g(x))g'(x)$

Gradient Descent for Linear Regression

- Initialize θ
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y_j^{(i)}) x_j^{(i)} \quad \begin{array}{l} \text{simultaneous update} \\ \text{for } j = 0 \dots d \end{array}$$

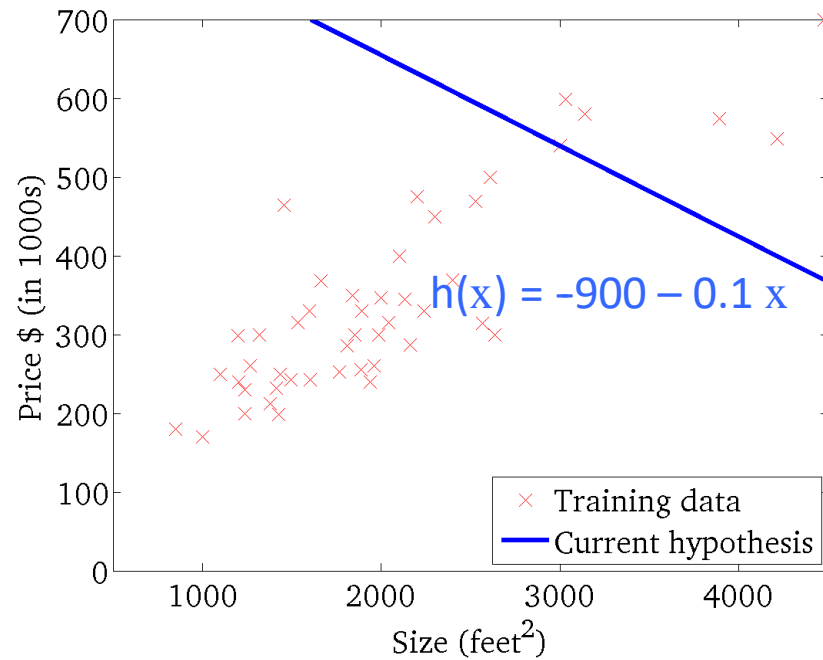
- To achieve simultaneous update
 - At the start of each GD iteration, compute $h_{\theta}(x^{(i)})$
 - Use this stored value in the update step loop
- Assume convergence when $\|\theta_{new} - \theta_{old}\|_2 < \epsilon$

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

Gradient Descent

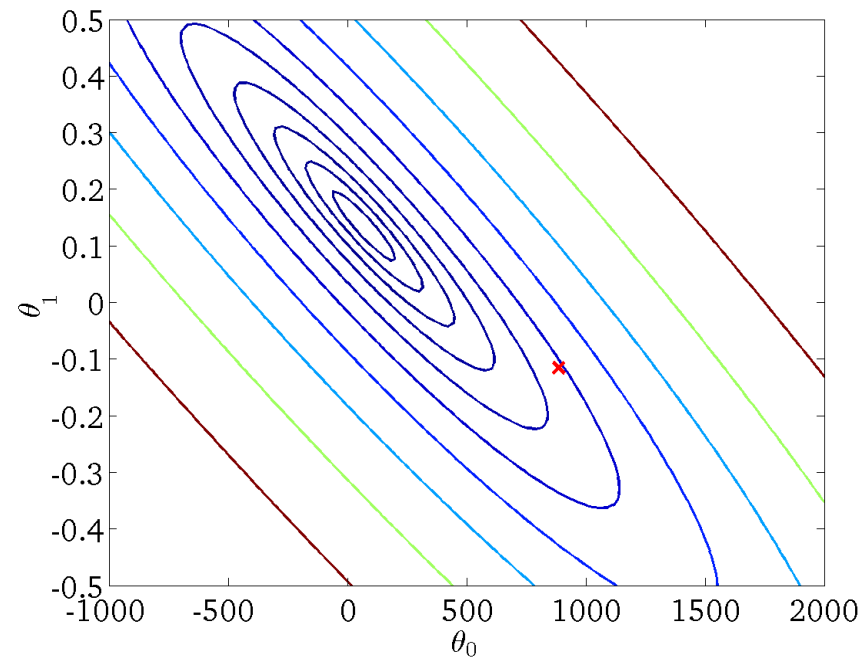
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



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

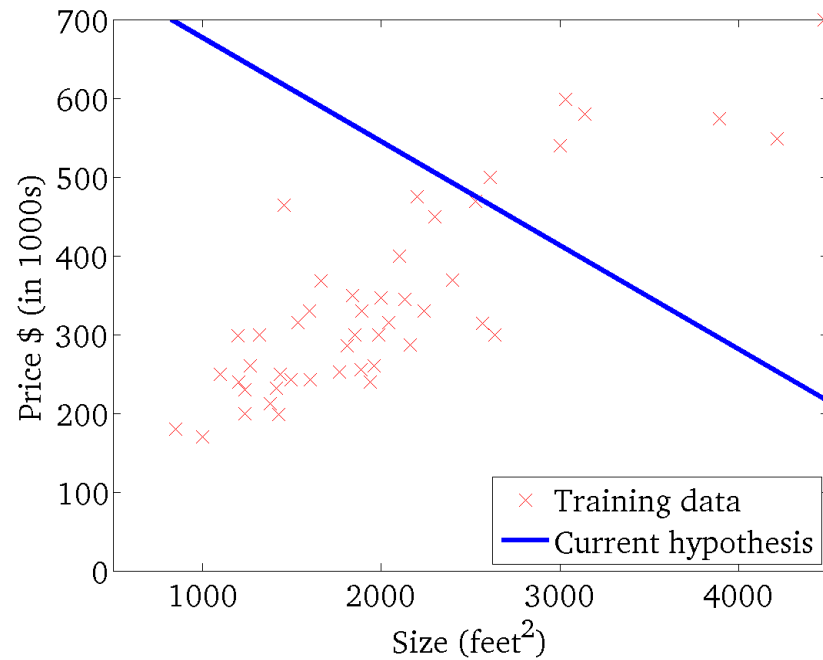
(function of the parameters θ_0, θ_1)



Gradient Descent

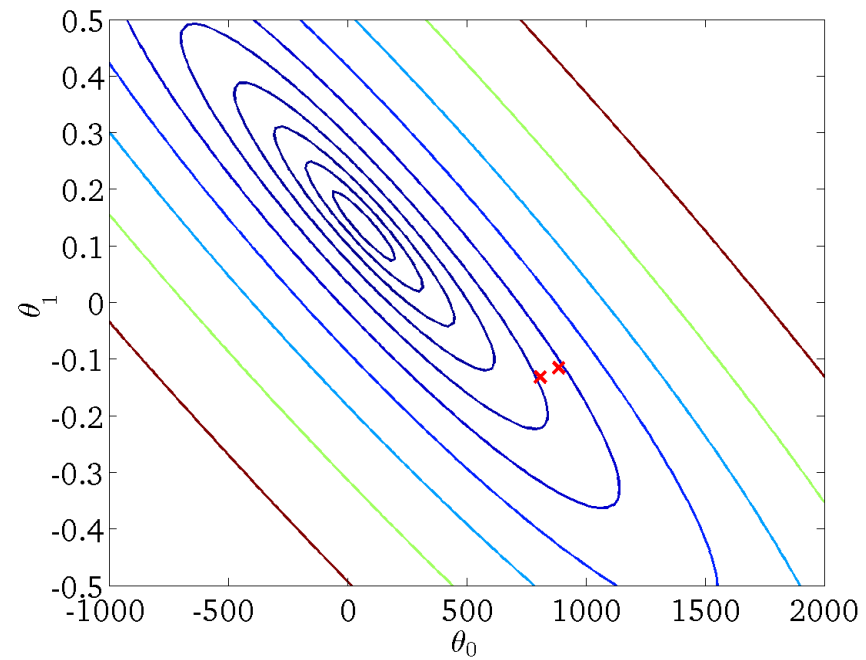
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



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

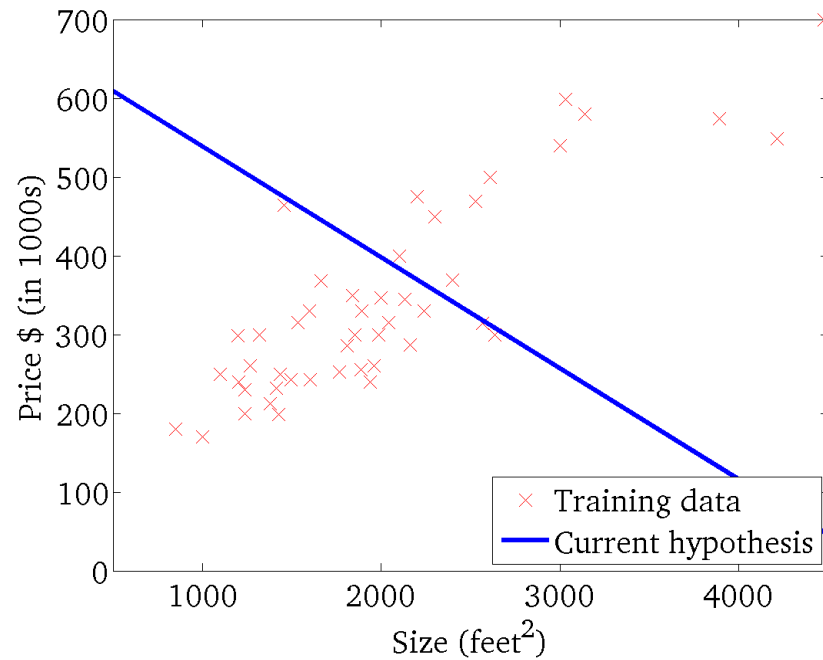
(function of the parameters θ_0, θ_1)



Gradient Descent

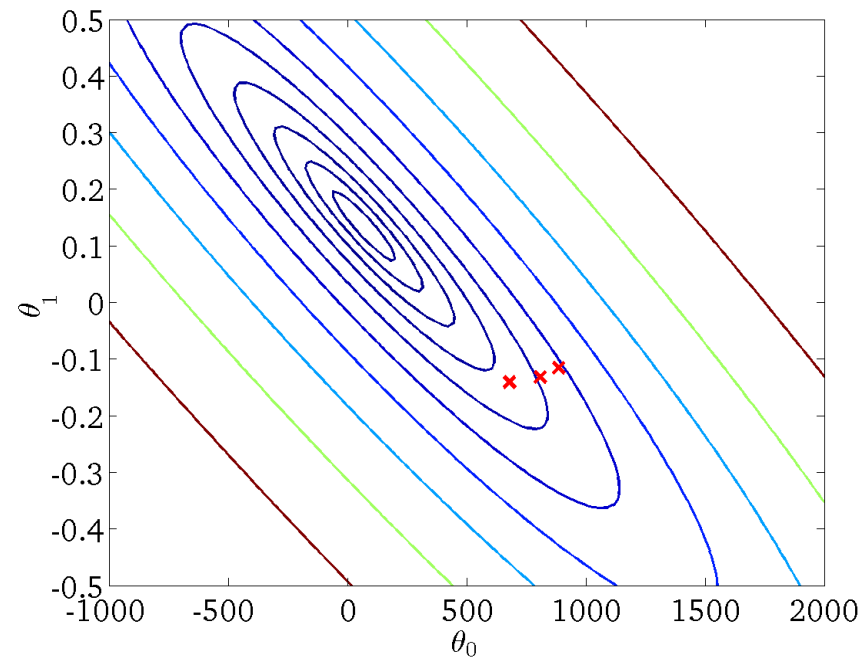
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



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

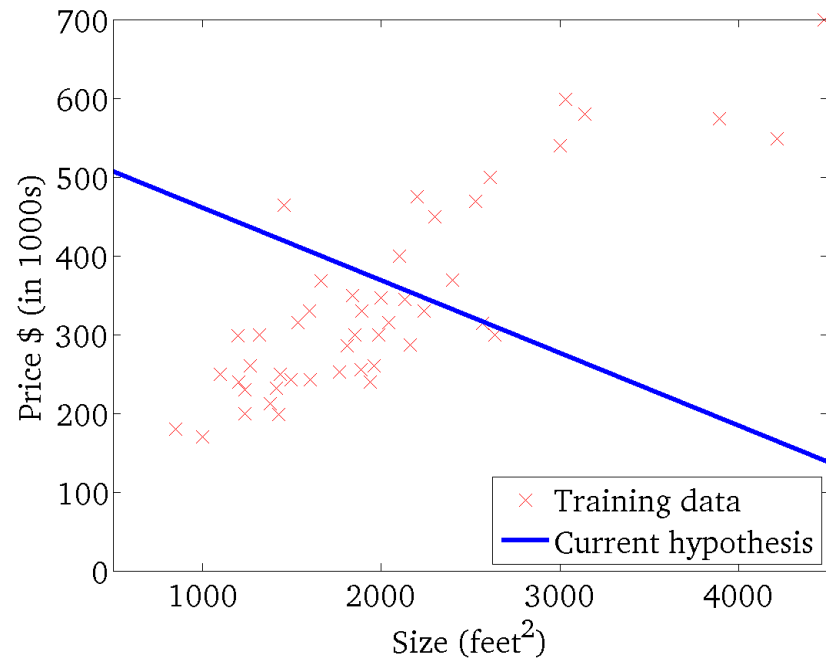
(function of the parameters θ_0, θ_1)



Gradient Descent

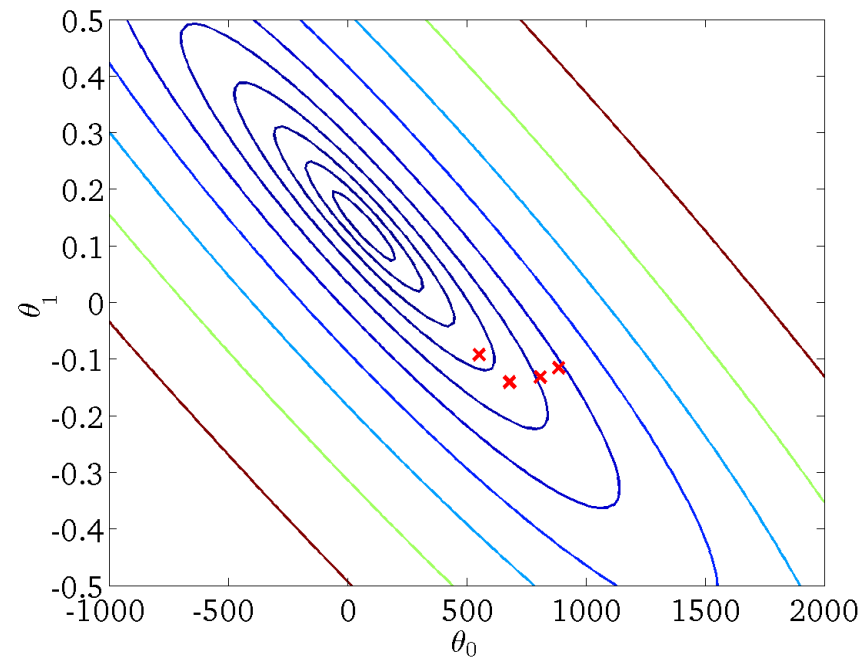
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



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

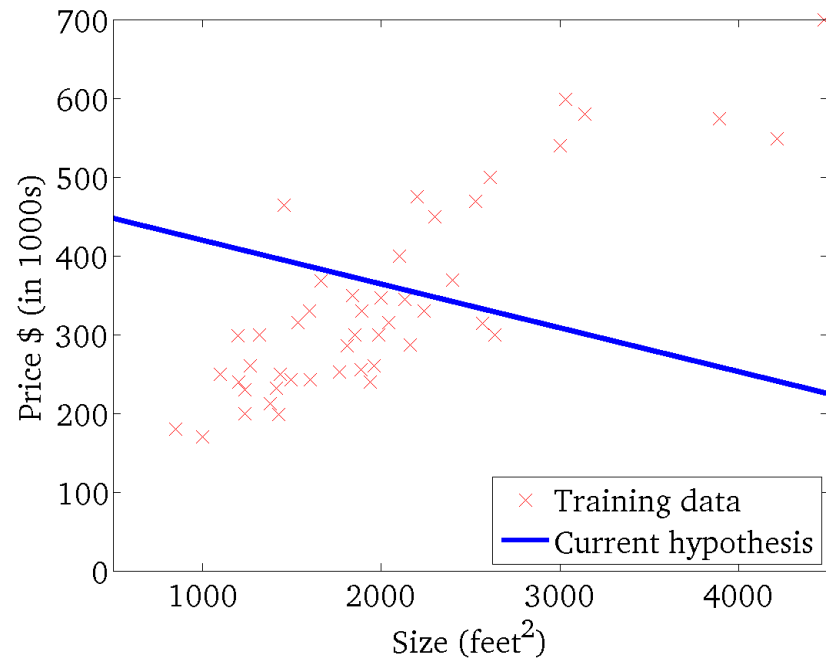
(function of the parameters θ_0, θ_1)



Gradient Descent

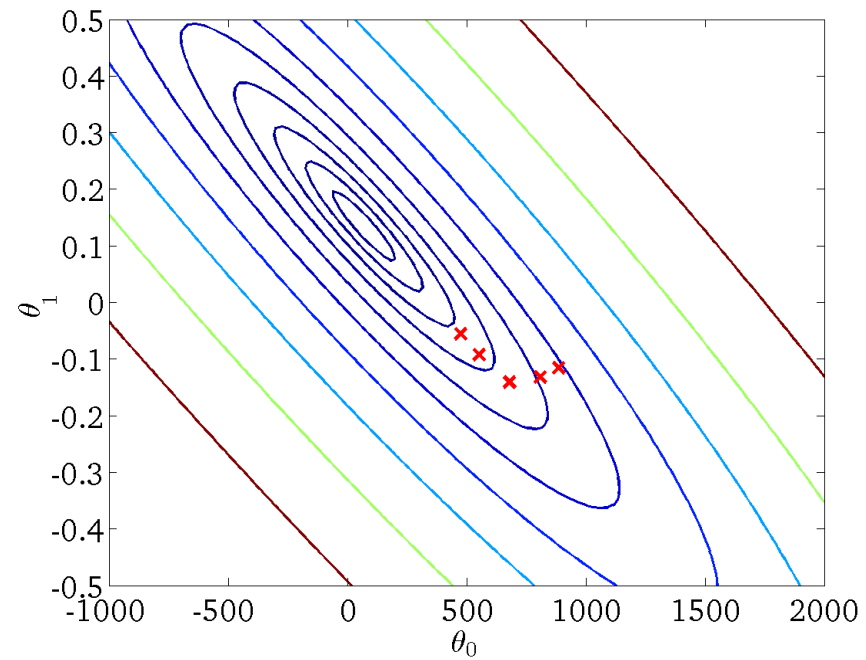
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



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

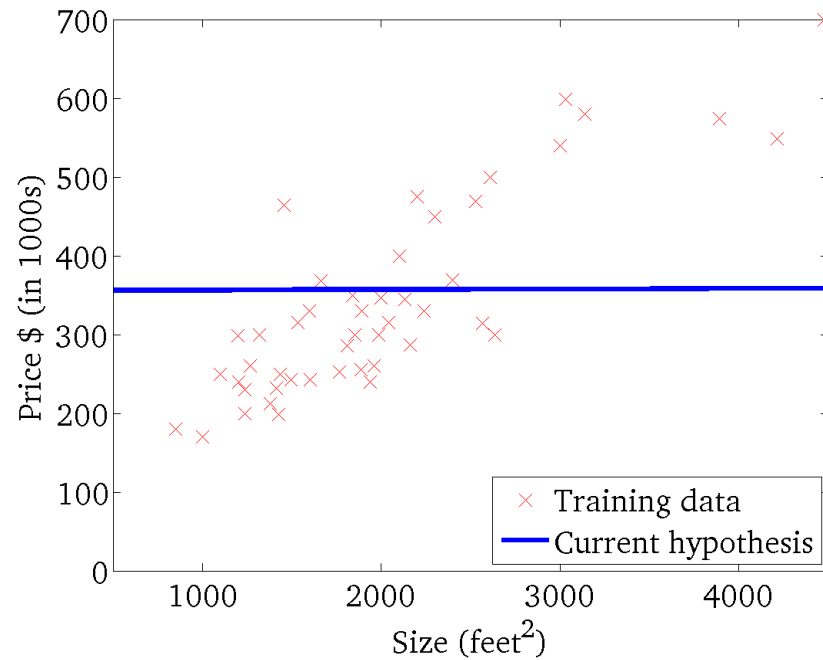
(function of the parameters θ_0, θ_1)



Gradient Descent

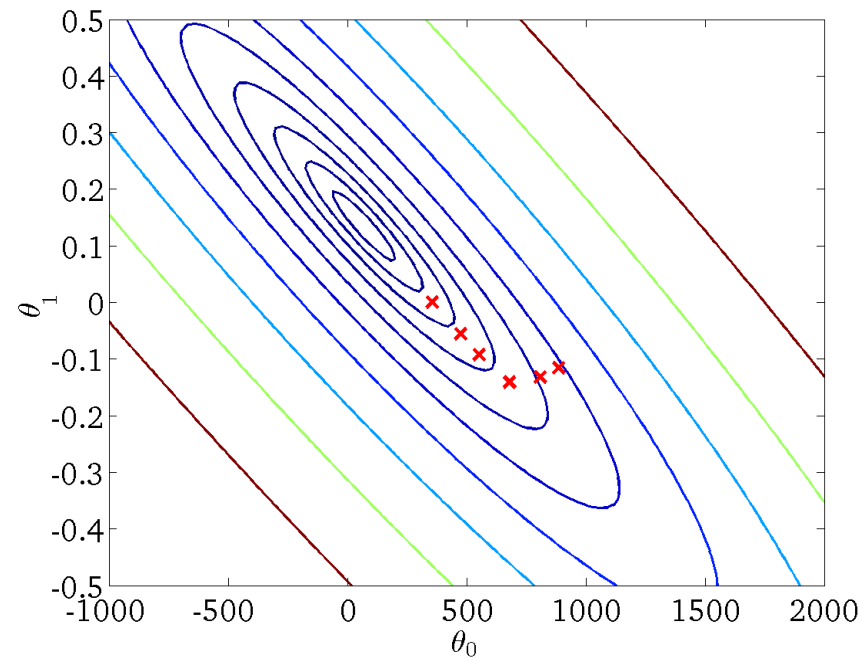
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



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

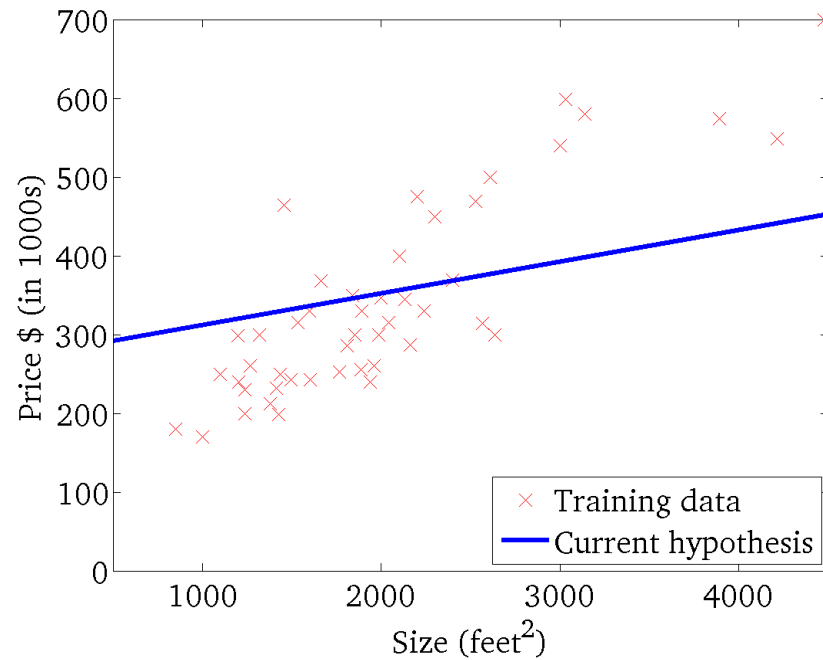
(function of the parameters θ_0, θ_1)



Gradient Descent

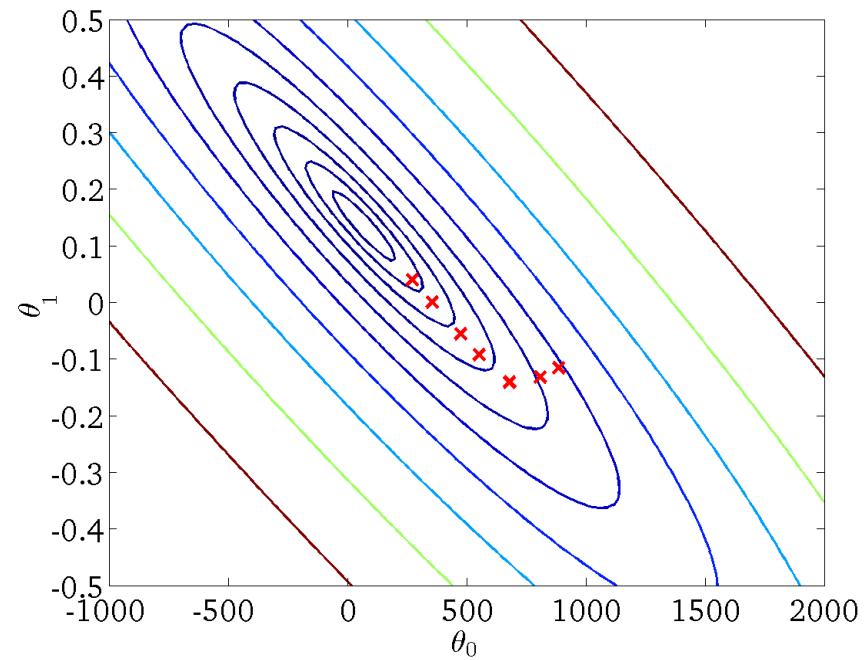
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



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

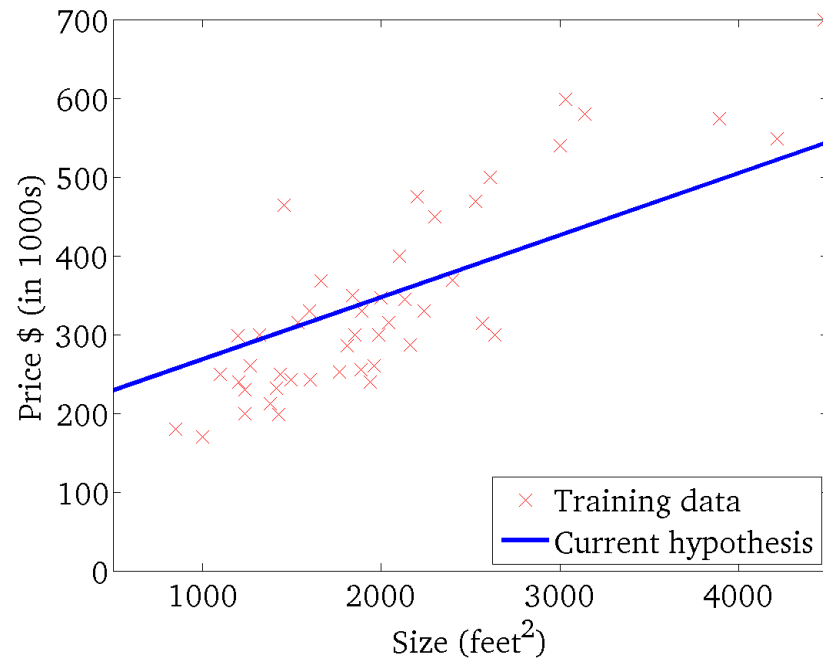
(function of the parameters θ_0, θ_1)



Gradient Descent

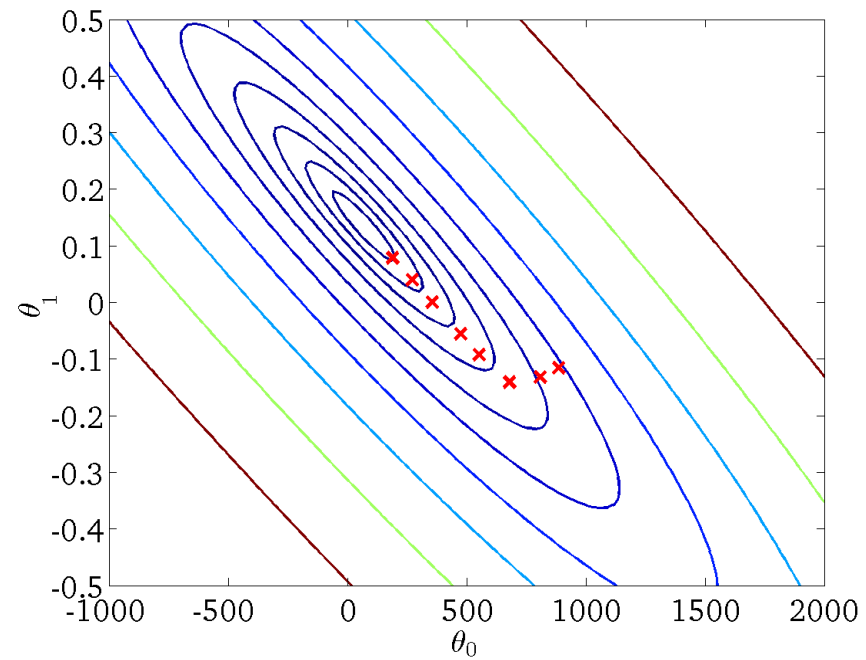
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



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

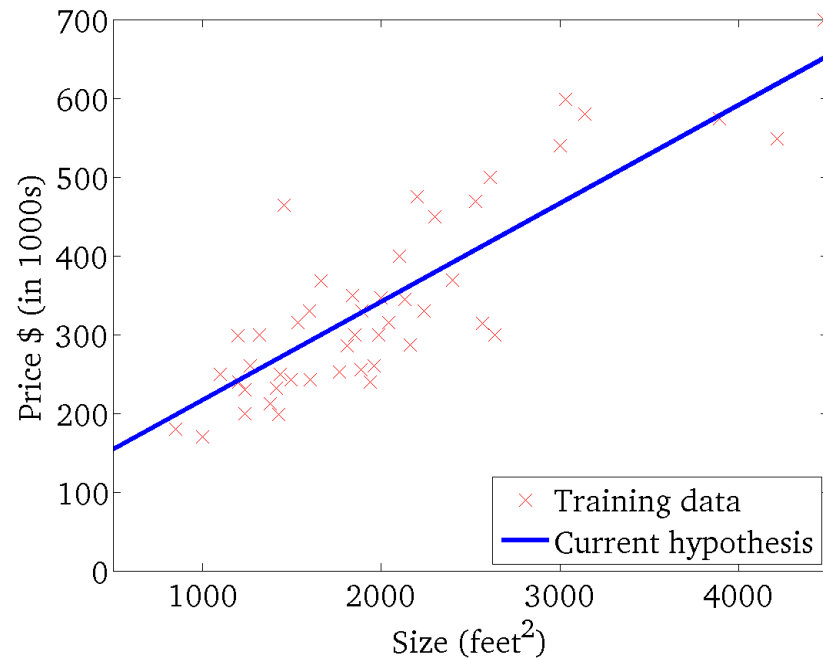
(function of the parameters θ_0, θ_1)



Gradient Descent

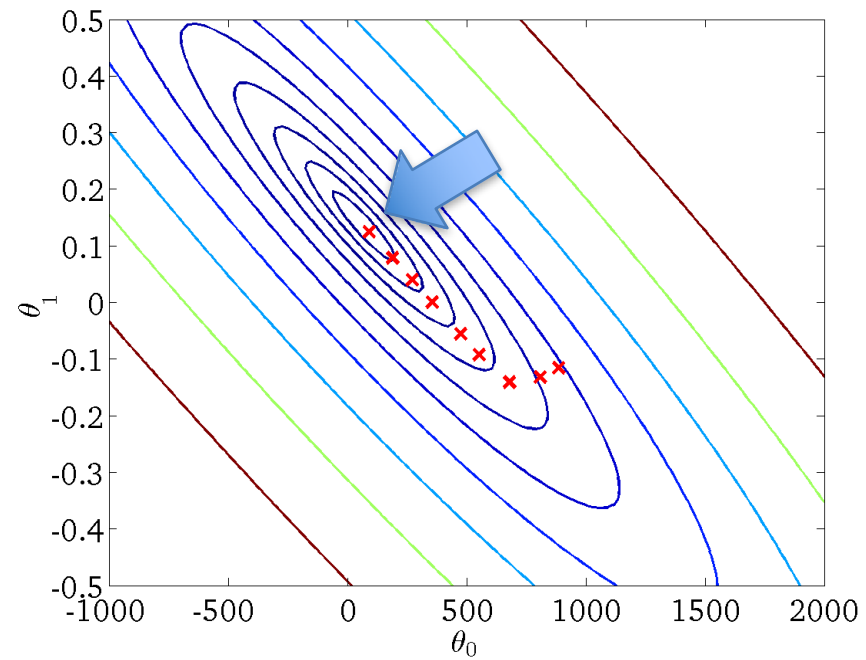
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



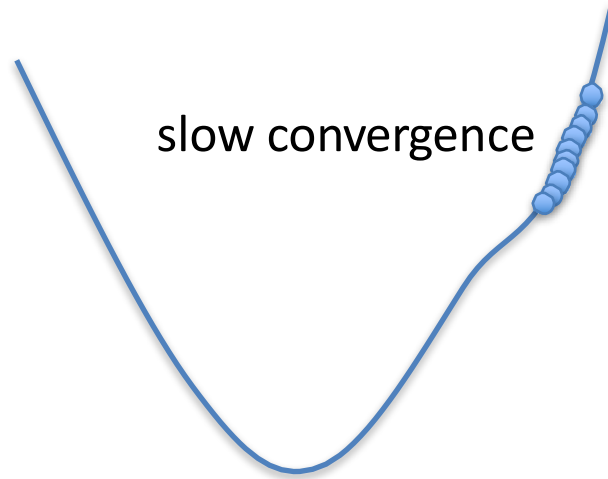
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)

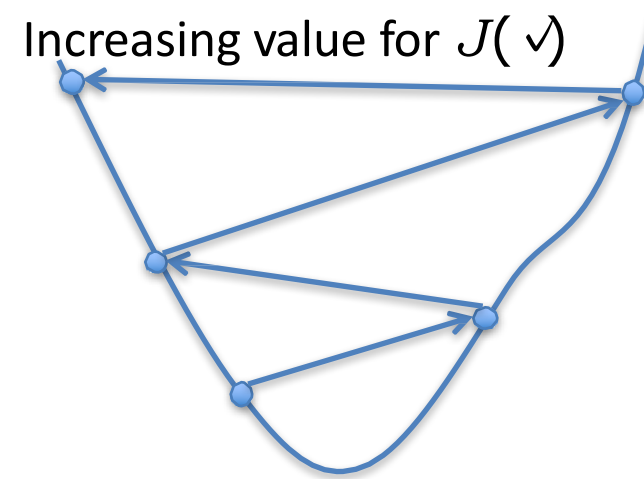


Choosing α

α too small



α 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 α

Extending Linear Regression to More Complex Models

- The inputs \mathbf{X} for linear regression can be:
 - Transformation of quantitative inputs
 - e.g. log, exp, square root, square, etc.
 - Interactions between variables
 - example: $x_3 = x_1 \times x_2$
 - Polynomial transformation
 - example: $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$
 - Basis expansions
 - Take the input x , and add “expand” it to make whatever method you are using more powerful.

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

Linear Basis Function Models

- Generally,

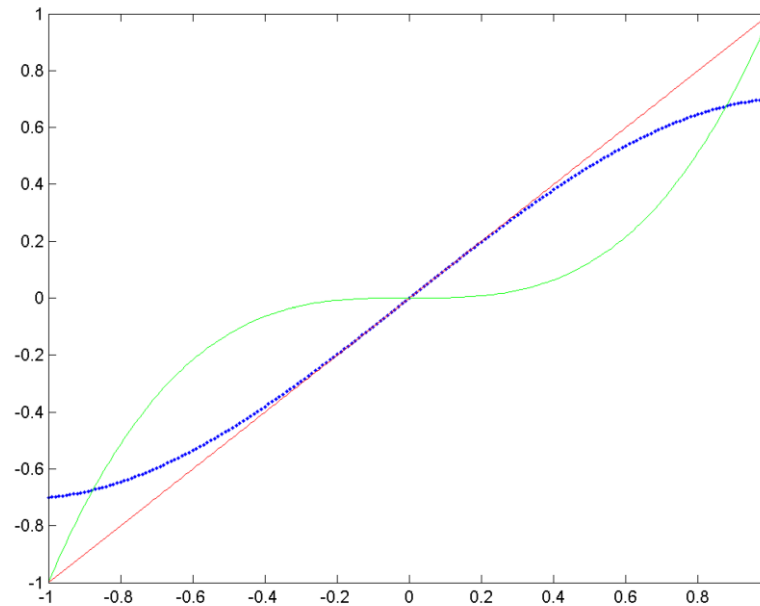
$$h_{\theta}(x) = \sum_{j=0}^d \theta_j \underbrace{\phi_j(x)}_{\text{basis function}}$$

- Typically, $\phi_0(x) = 1$ so that θ_0 acts as a bias
- In the simplest case, we use linear basis functions:

$$\phi_j(x) = x_j$$

Example of Basis Function Models

$$f(x) = x - 0.3x^3$$



Basis functions $\phi_1(x) = 1, \phi_2(x) = x, \phi_3(x) = x^2, \phi_4(x) = x^3$ und $\theta = (0, 1, 0, -0.3)$

Basic Idea of Basis Function Models

- The simple idea: in addition to the original inputs, we add inputs that are calculated as deterministic functions of the existing inputs and treat them as additional inputs

- Example: Polynomial Basis Functions

$$\{1, x_1, x_2, x_3, x_1x_2, x_1x_3, x_2x_3, x_1^2, x_2^2, x_3^2\}$$

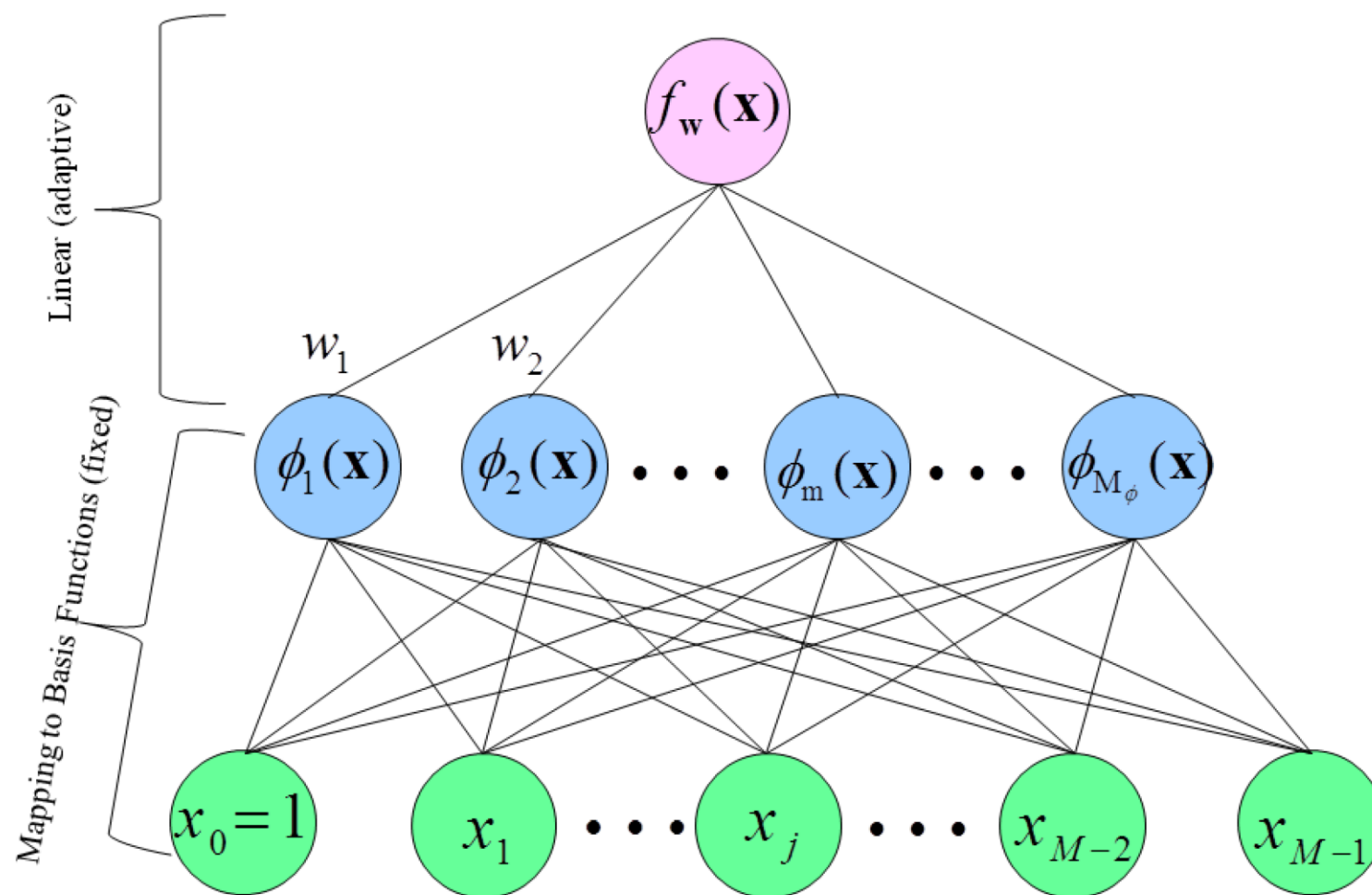
- Basis functions $\{\phi_m(\mathbf{x})\}_{m=1}^{M_\phi}$

- In the example:

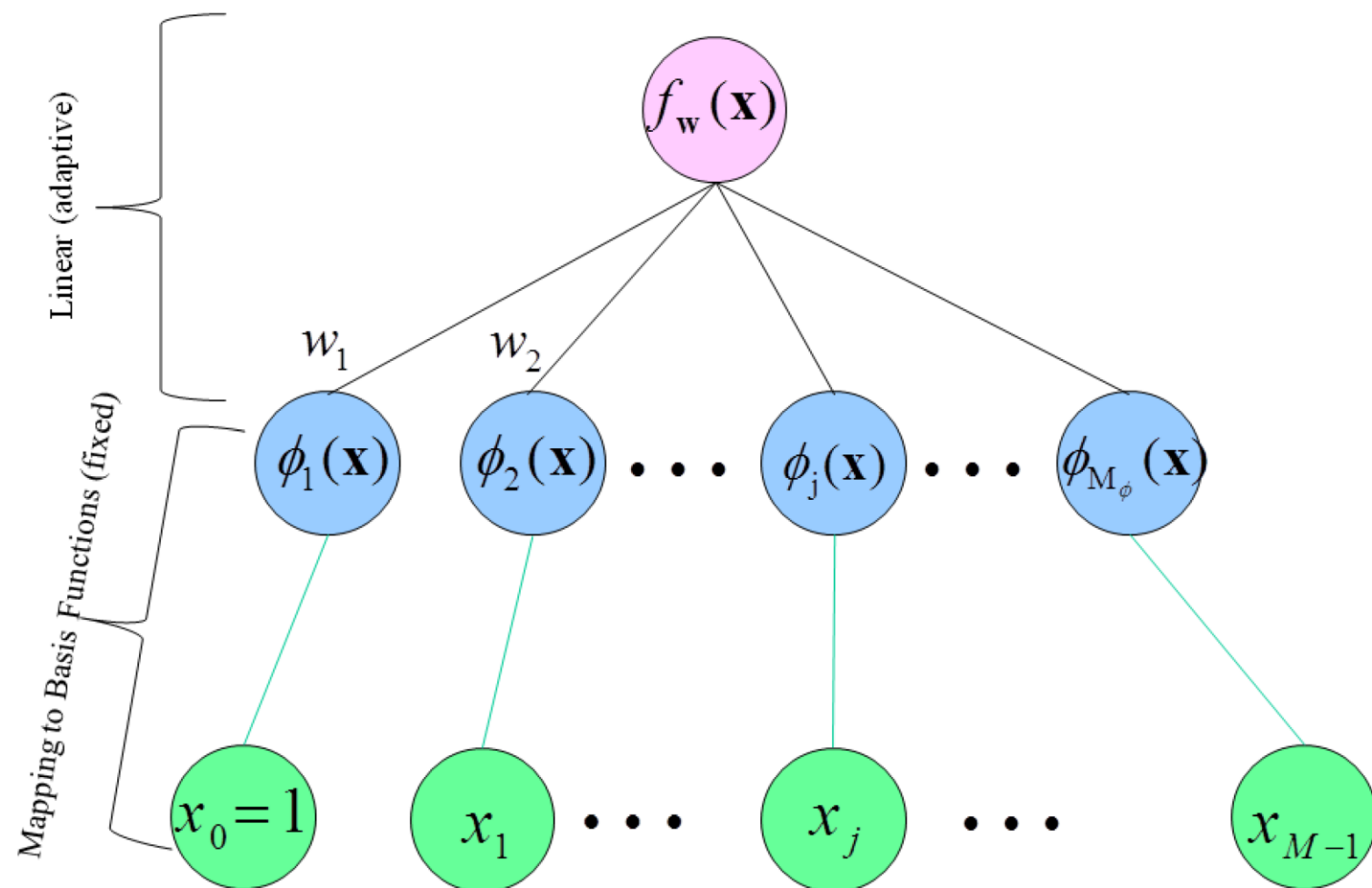
$$\phi_1(\mathbf{x}) = 1 \quad \phi_2(\mathbf{x}) = x_1 \quad \phi_6(\mathbf{x}) = x_1x_3 \quad \dots$$

- Independent of the choice of basis functions, the regression parameters are calculated using the well-known equations for linear regression

Network of Basis Functions



Network of Linear Basis Functions

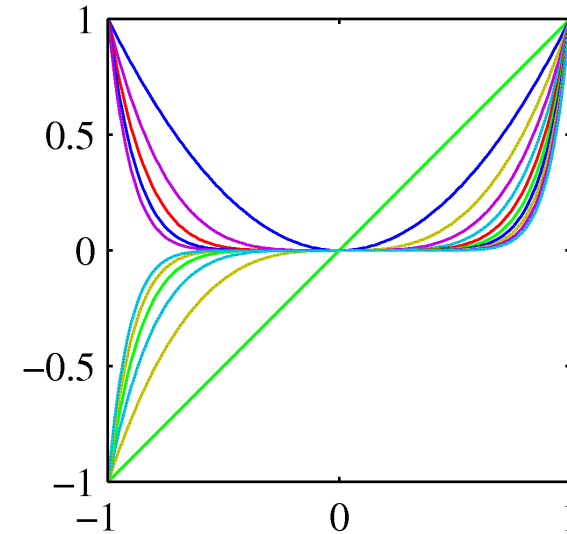


Linear Basis Function Models

- Polynomial basis functions:

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

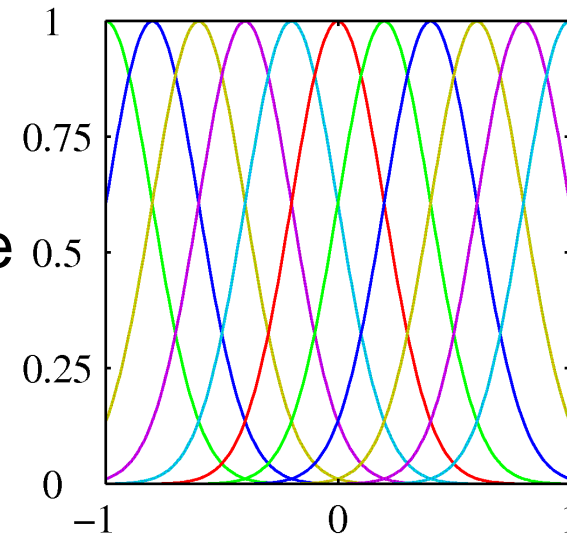
- These are global; a small change in x affects all basis functions



- 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. μ_j and s control location and scale (width).



Linear Basis Function Models

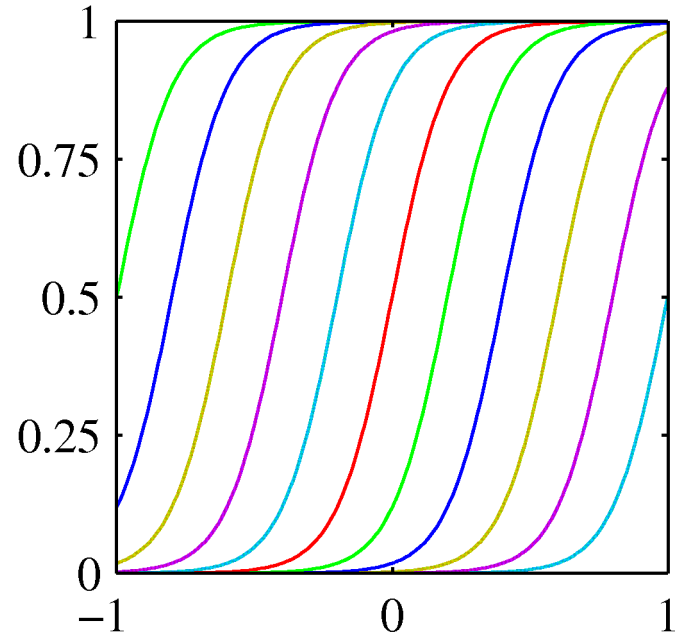
- Sigmoidal basis functions

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

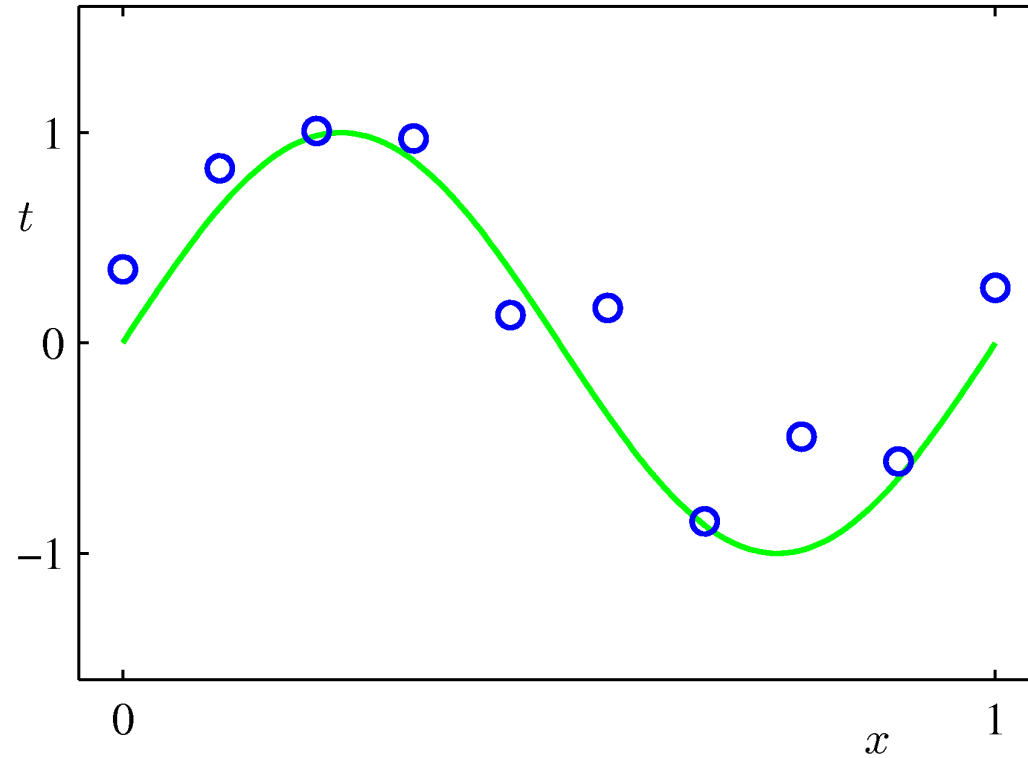
where

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

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



Example of Fitting a Polynomial Curve with a Linear Model



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

Linear Basis Function Models

- Basic Linear Model:

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

- Generalized Linear Model:

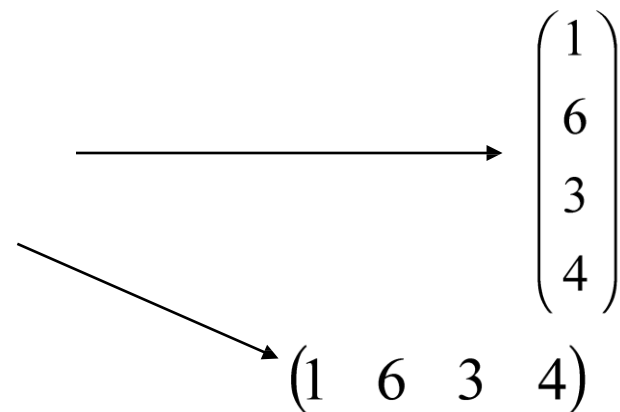
$$h_{\theta}(x) = \sum_{j=0}^d \theta_j \phi_j(x)$$

- 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

Linear Algebra Concepts

- Vector in \mathbb{R}^d is an ordered set of d real numbers

- e.g., $v = [1,6,3,4]$ is in \mathbb{R}^4
- “[1,6,3,4]” is a column vector:
- as opposed to a row vector:



- An m -by- n matrix is an object with m rows and n columns, where each entry is a real number:

$$\begin{pmatrix} 1 & 2 & 8 \\ 4 & 78 & 6 \\ 9 & 3 & 2 \end{pmatrix}$$

Linear Algebra Concepts

- Transpose: flips a matrix over its diagonal

$$\begin{pmatrix} a \\ b \end{pmatrix}^T = (a \quad b)$$

$$\begin{pmatrix} a & b \\ c & d \end{pmatrix}^T = \begin{pmatrix} a & c \\ b & d \end{pmatrix}$$

- Note: $(Ax)^T = x^T A^T$

(We'll define multiplication soon...)

- Vector norms:

- L_p norm of $\mathbf{v} = (v_1, \dots, v_k)$ is

- Common norms: L_1 , L_2

- $L_{\text{infinity}} = \max_i |v_i|$

$$\left(\sum_i |v_i|^p \right)^{\frac{1}{p}}$$

- Length of a vector \mathbf{v} is $L_2(\mathbf{v})$

Linear Algebra Concepts

- Vector dot product: $u \bullet v = \begin{pmatrix} u_1 & u_2 \end{pmatrix} \bullet \begin{pmatrix} v_1 & v_2 \end{pmatrix} = u_1 v_1 + u_2 v_2$
 - Note: dot product of u with itself = $\text{length}(u)^2 = \|u\|_2^2$

- Matrix product:

$$A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}, B = \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{pmatrix}$$
$$AB = \begin{pmatrix} a_{11}b_{11} + a_{12}b_{21} & a_{11}b_{12} + a_{12}b_{22} \\ a_{21}b_{11} + a_{22}b_{21} & a_{21}b_{12} + a_{22}b_{22} \end{pmatrix}$$

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

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

Vectorization

- Consider our model for n instances:

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

- Let

$$\boldsymbol{\theta}=\left[\begin{array}{c} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{array}\right] \quad \mathbf{X}=\left[\begin{array}{cccc} 1 & x_1^{(1)} & \cdots & x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & \cdots & x_d^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \cdots & x_d^{(n)} \end{array}\right]$$

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

- Can write the model in vectorized form as

$$h_{\boldsymbol{\theta}}(\mathbf{x})=\mathbf{X} \boldsymbol{\theta}$$

Vectorization

- For the linear regression cost function:

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

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

Let:

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

HW2

- <https://canvas.wpi.edu/courses/58900/assignments/35566>
6