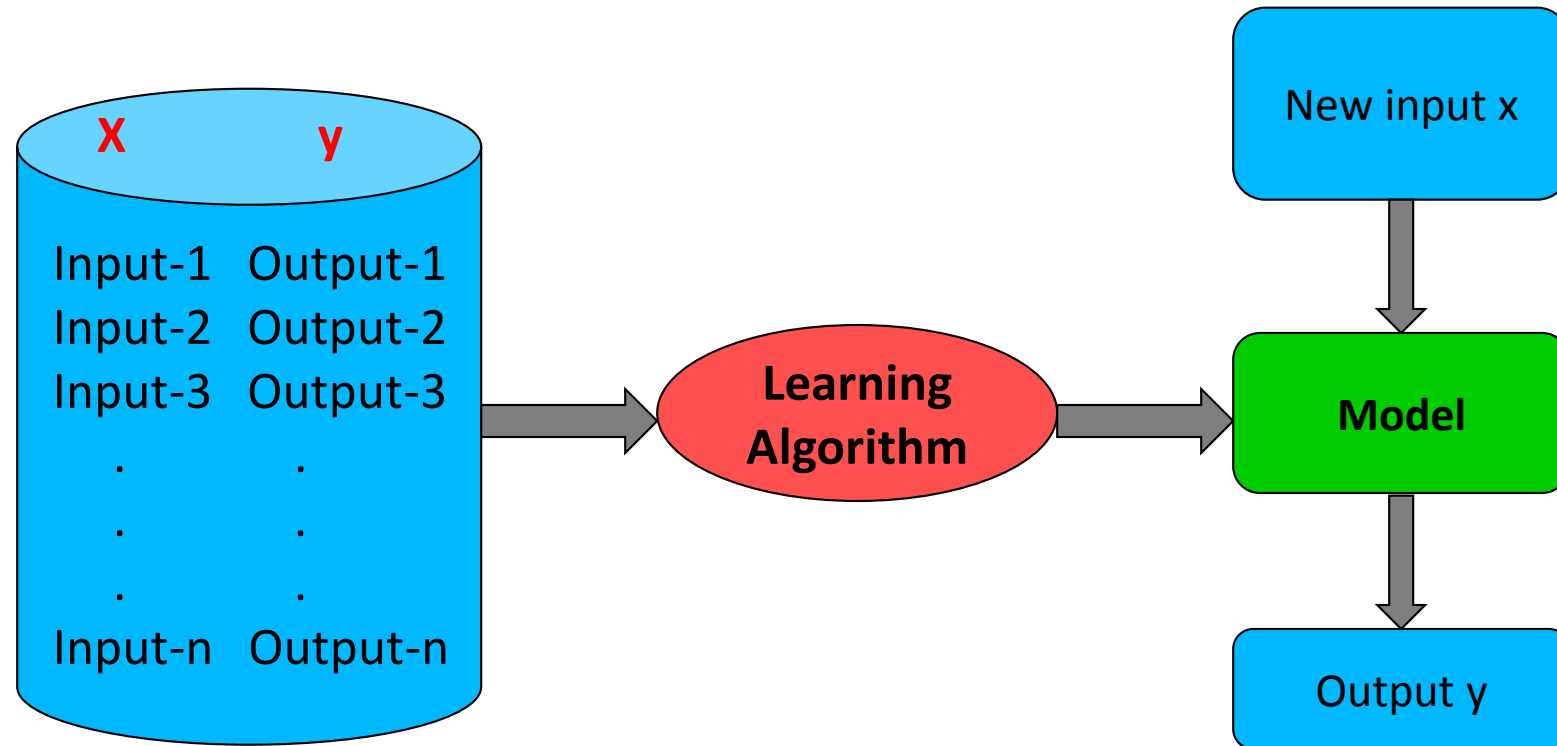


Supervised Learning

Regression



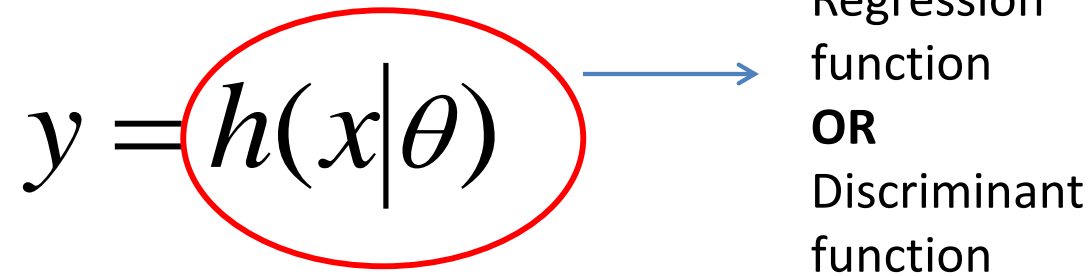
Classification

Supervised Learning

- A model defined with a set of parameters.

$$y = h(x|\theta)$$

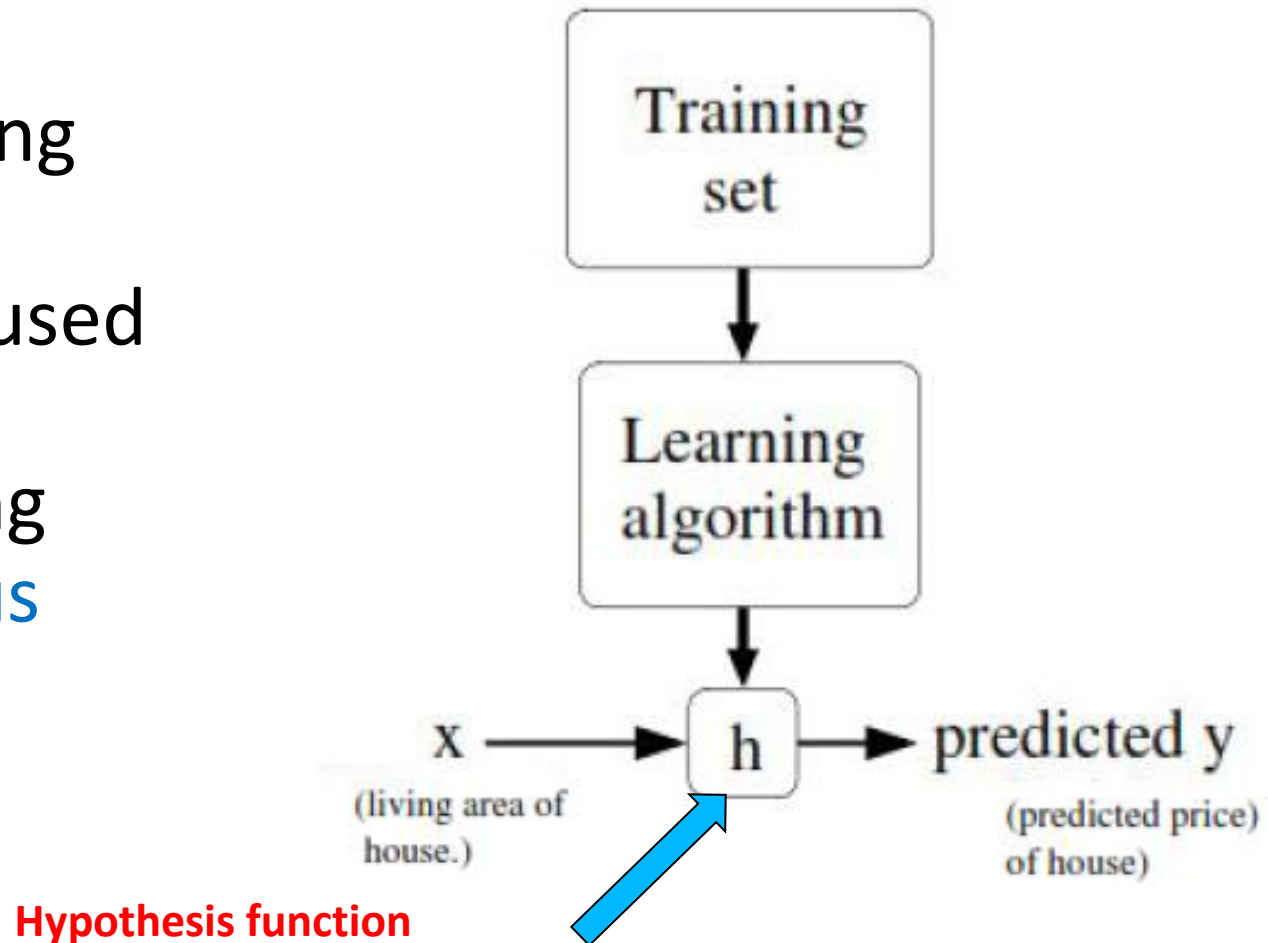
Regression
function
OR
Discriminant
function



- Where $h(\cdot)$ is the model and θ are its parameters
- **Regression** : y – number
- **Classification**: y – class code (e.g. 0/1)

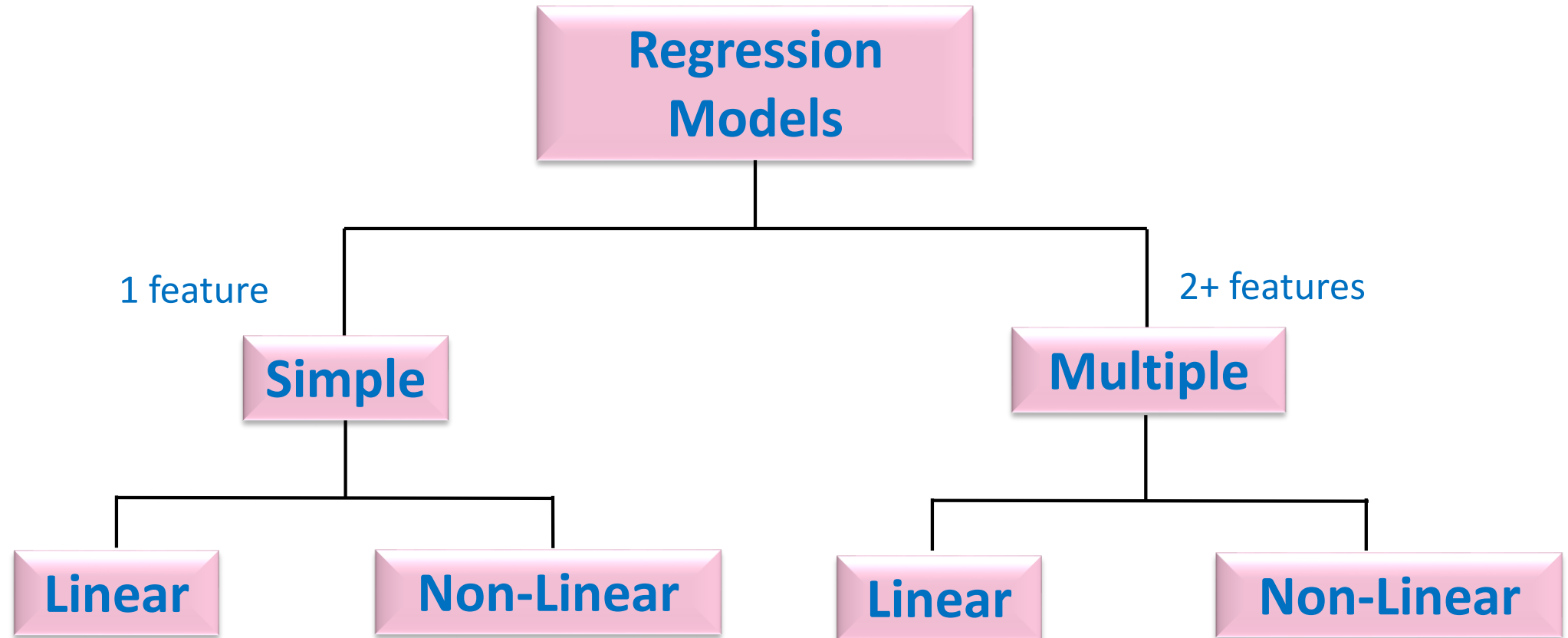
Regression

- Supervised Learning tasks where the datasets that are used for predictive / statistical modeling contain **continuous labels**.



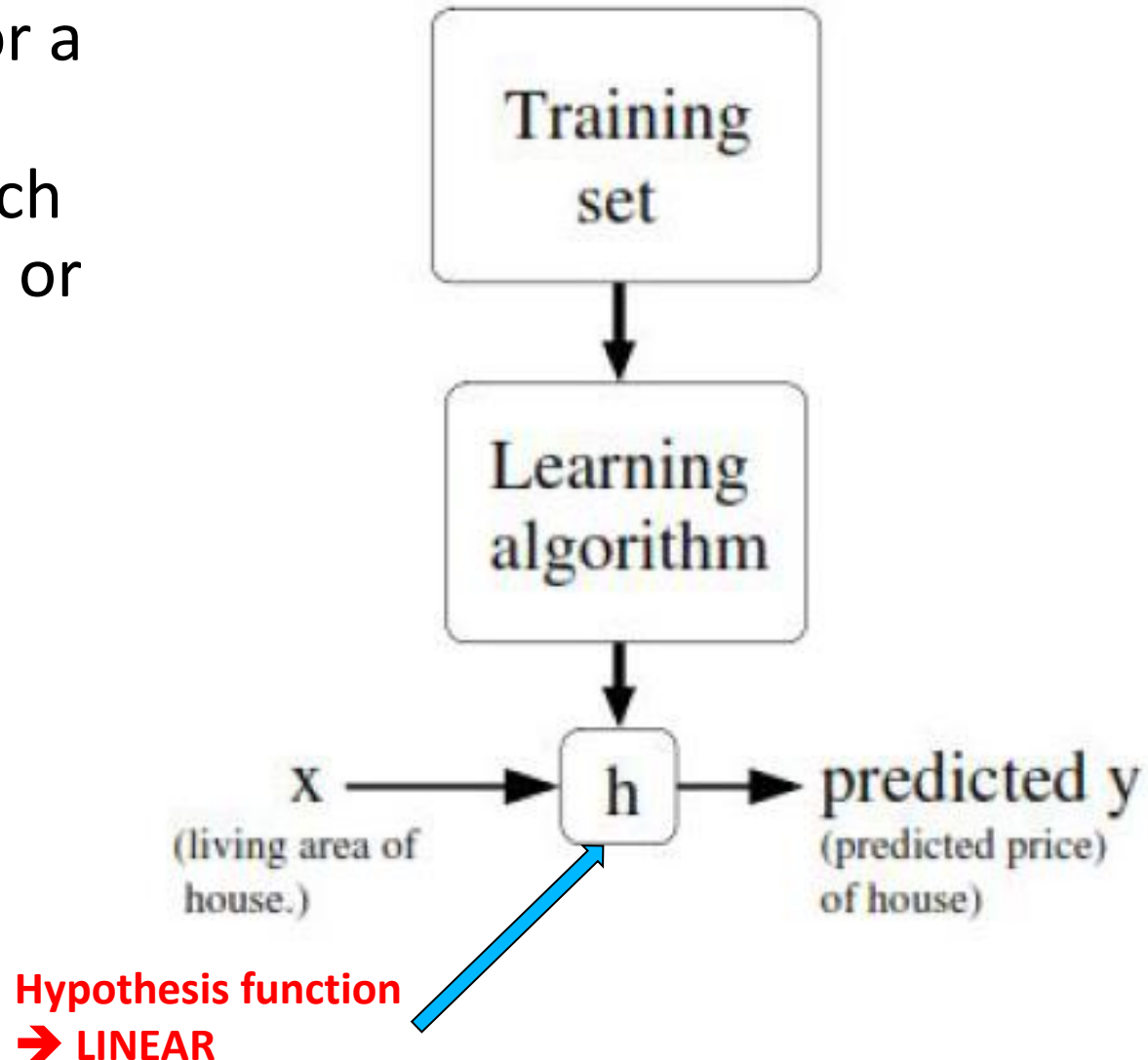
Learn a function $h(x)$, so that $h(x)$ is a good predictor for the corresponding value of y

Types of Regression Models



Linear Regression

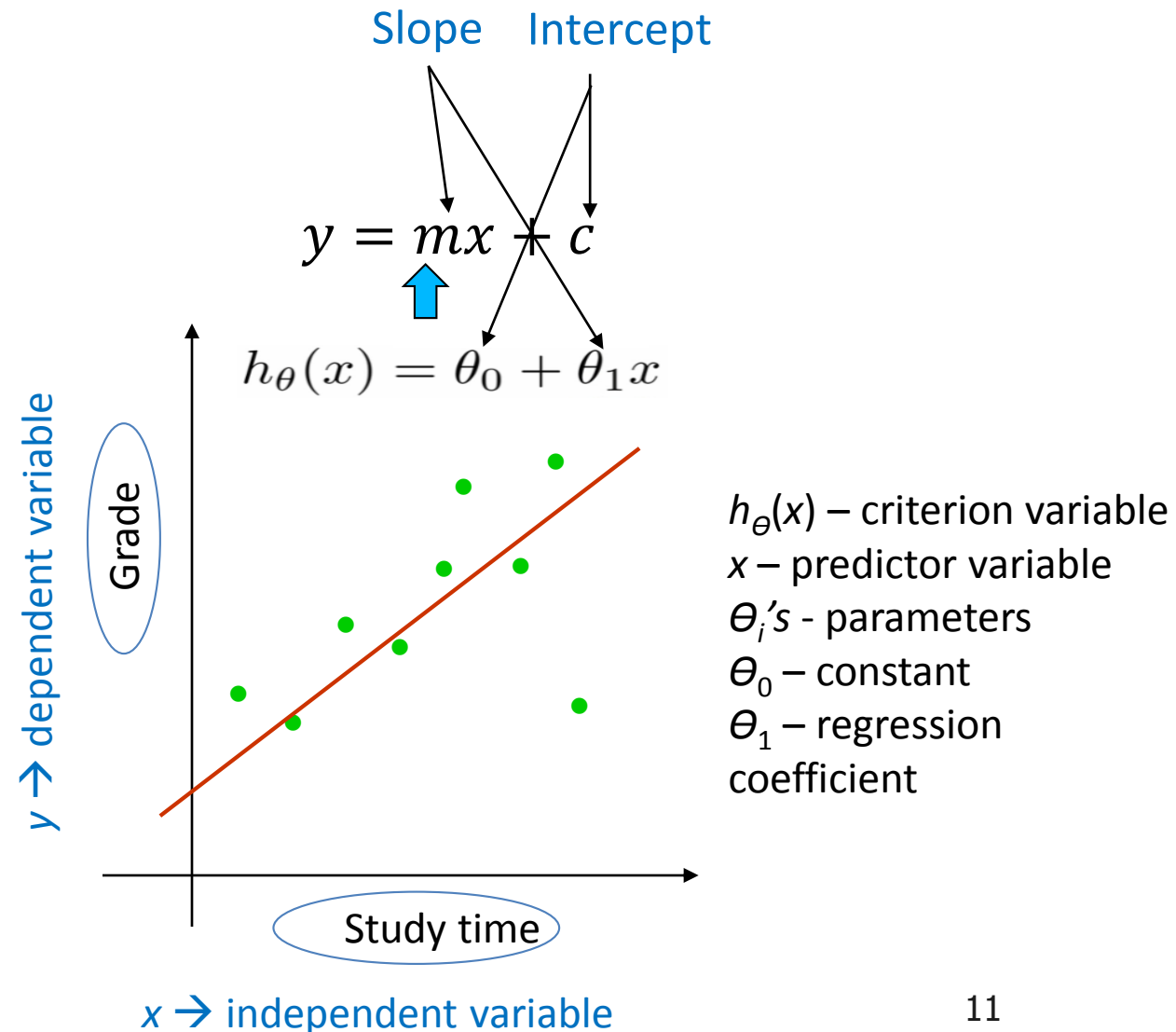
- Find a linear equation for a continuous response variable known as Y which will be a function of one or more variables (X)



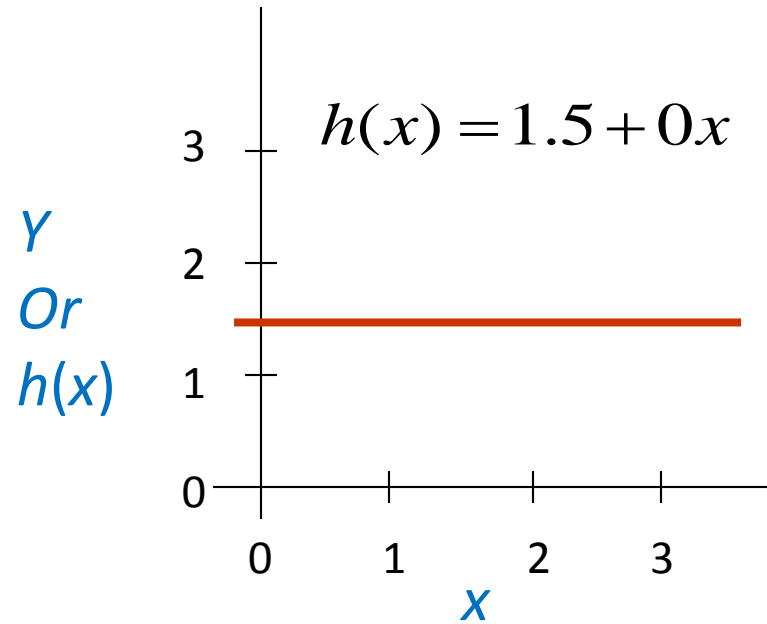
Linear Regression

- Given an input x , compute an output y
- For example –
 - Predict height from age
 - Predict house price from house area
 - Predict grade from study time

Linear regression with one variable / feature
Univariate Linear Regression

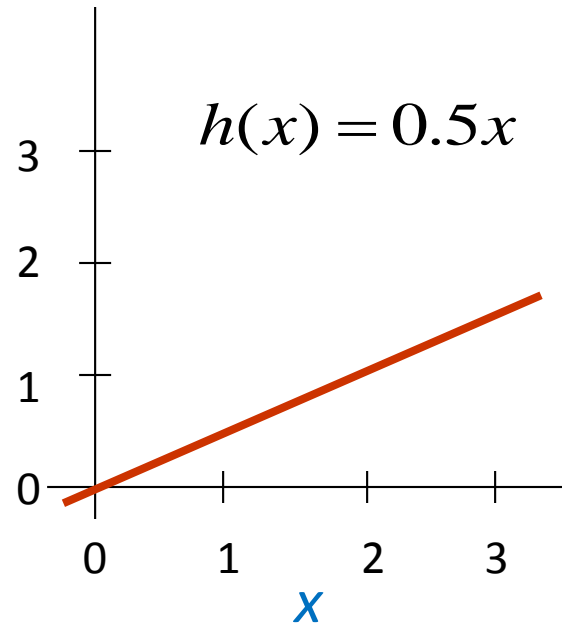


$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



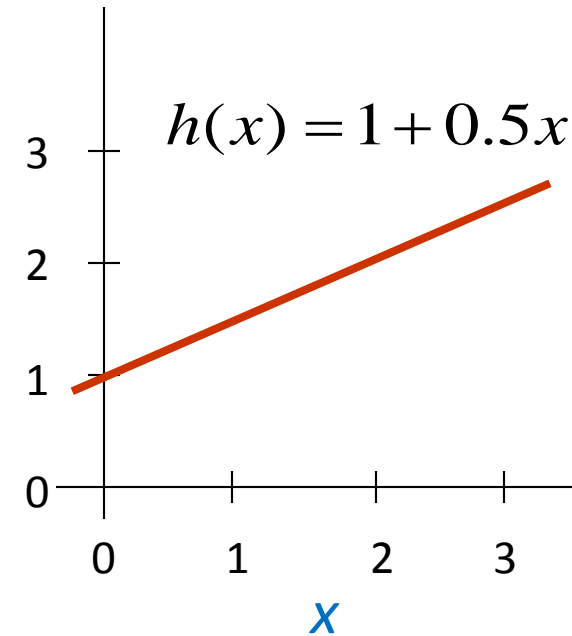
$$\theta_0 = 1.5$$

$$\theta_1 = 0$$



$$\theta_0 = 0$$

$$\theta_1 = 0.5$$

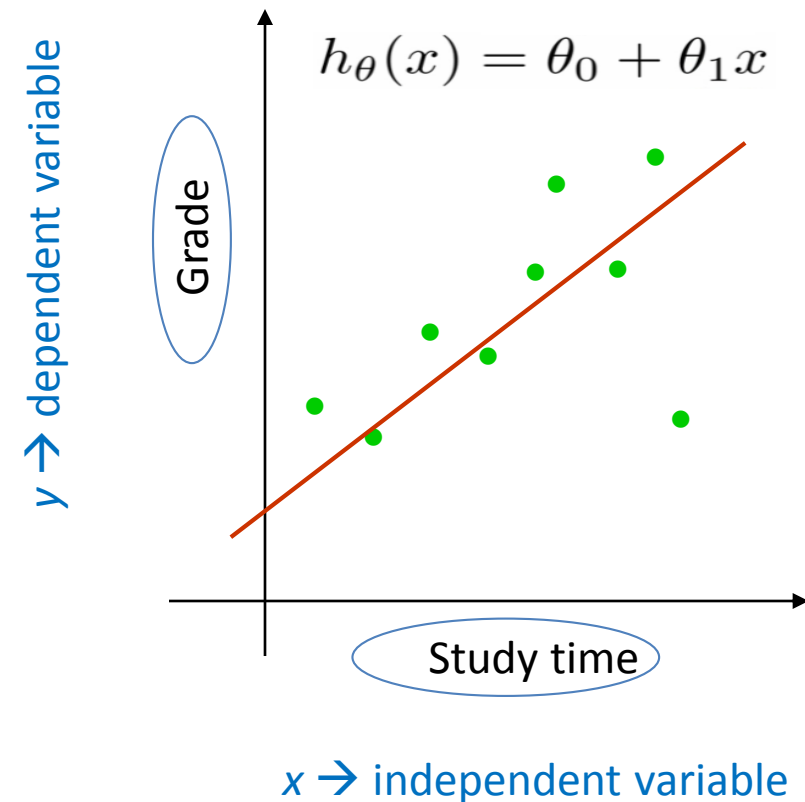


$$\theta_0 = 1$$

$$\theta_1 = 0.5$$

Linear Regression

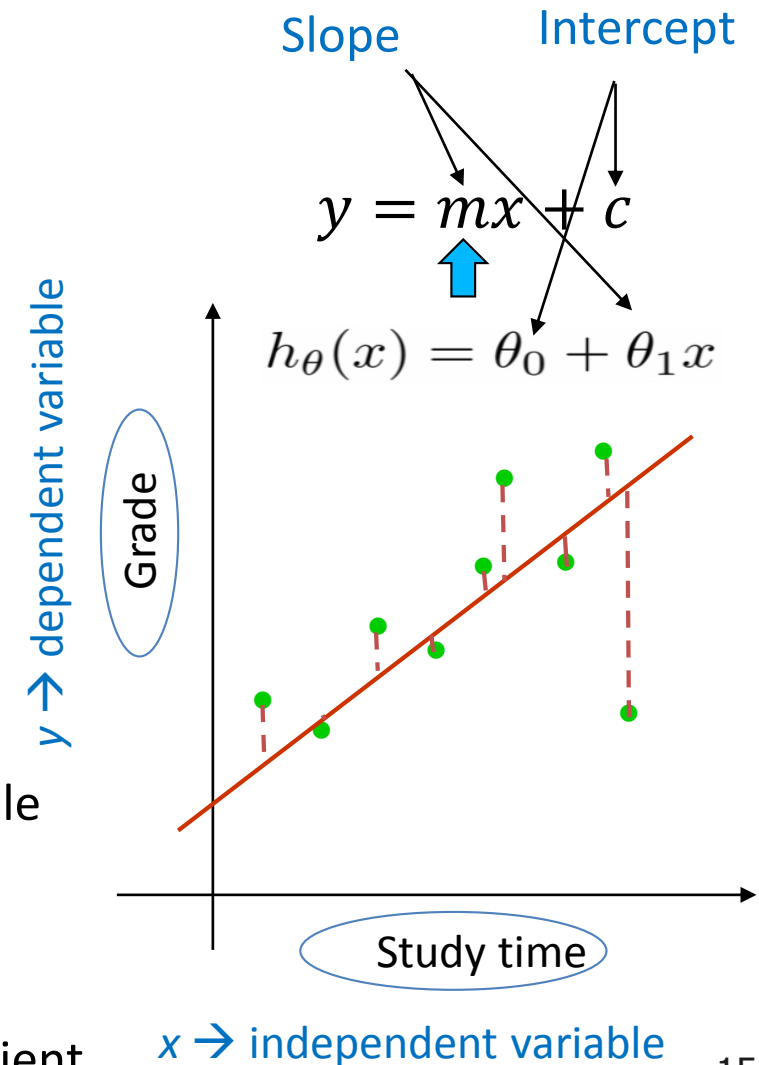
- Idea: Choose θ_0 , θ_1 so that $h_{\theta}(x)$ is close to y for our training example (X, y)



Linear Regression Line

- The **least-squares regression line** is the unique line such that the sum of squared vertical (y) distances between the data points and the line is the smallest possible.

$h_{\theta}(x)$ – criterion variable
 x – predictor variable
 θ_i 's – parameters
 θ_0 – constant
 θ_1 – regression coefficient

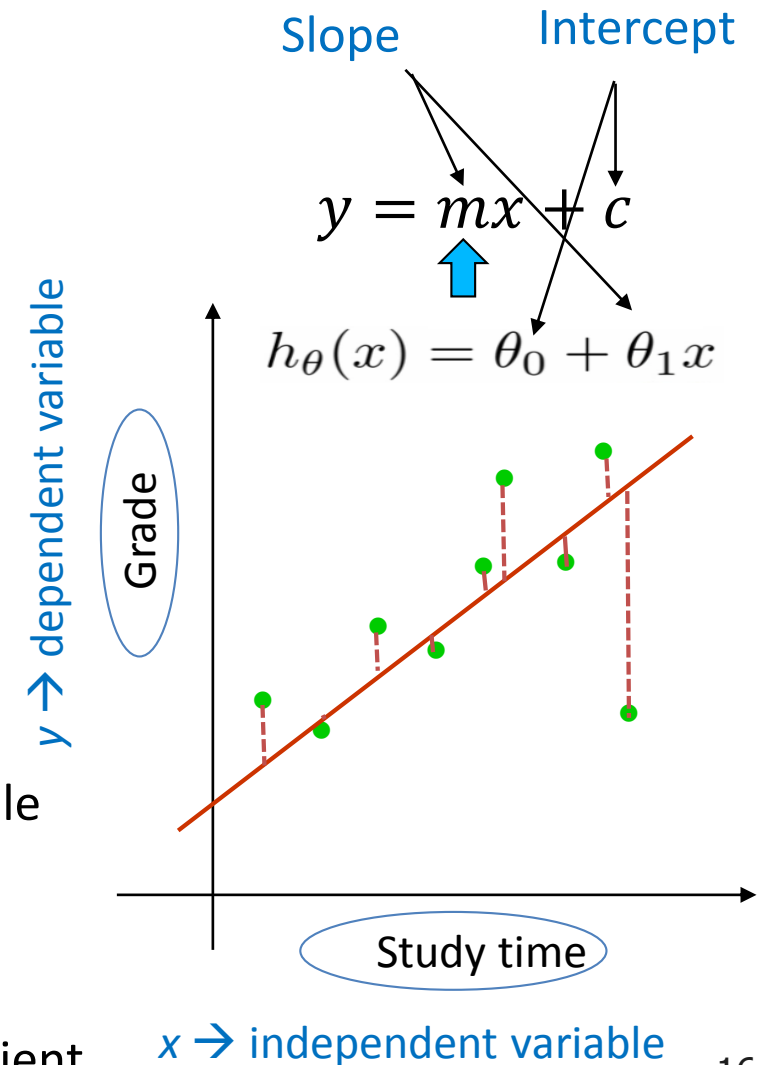


Linear Regression Line

- How to compute an error ?
- For given $x^{(i)}$, actual output is $y^{(i)}$ and using regression line (for some θ_0 and θ_1), it is $(\theta_0 + \theta_1 x^{(i)})$

- Error for $x^{(i)}$
= (predicted $y^{(i)}$ – actual $y^{(i)}$)
= $((\theta_0 + \theta_1 x^{(i)}) - y^{(i)})$
= $(h_{\theta}(x^{(i)}) - y^{(i)})$

$h_{\theta}(x)$ – criterion variable
 x – predictor variable
 θ_i 's – parameters
 θ_0 – constant
 θ_1 – regression coefficient



Linear Regression Line

- Sum of squared error

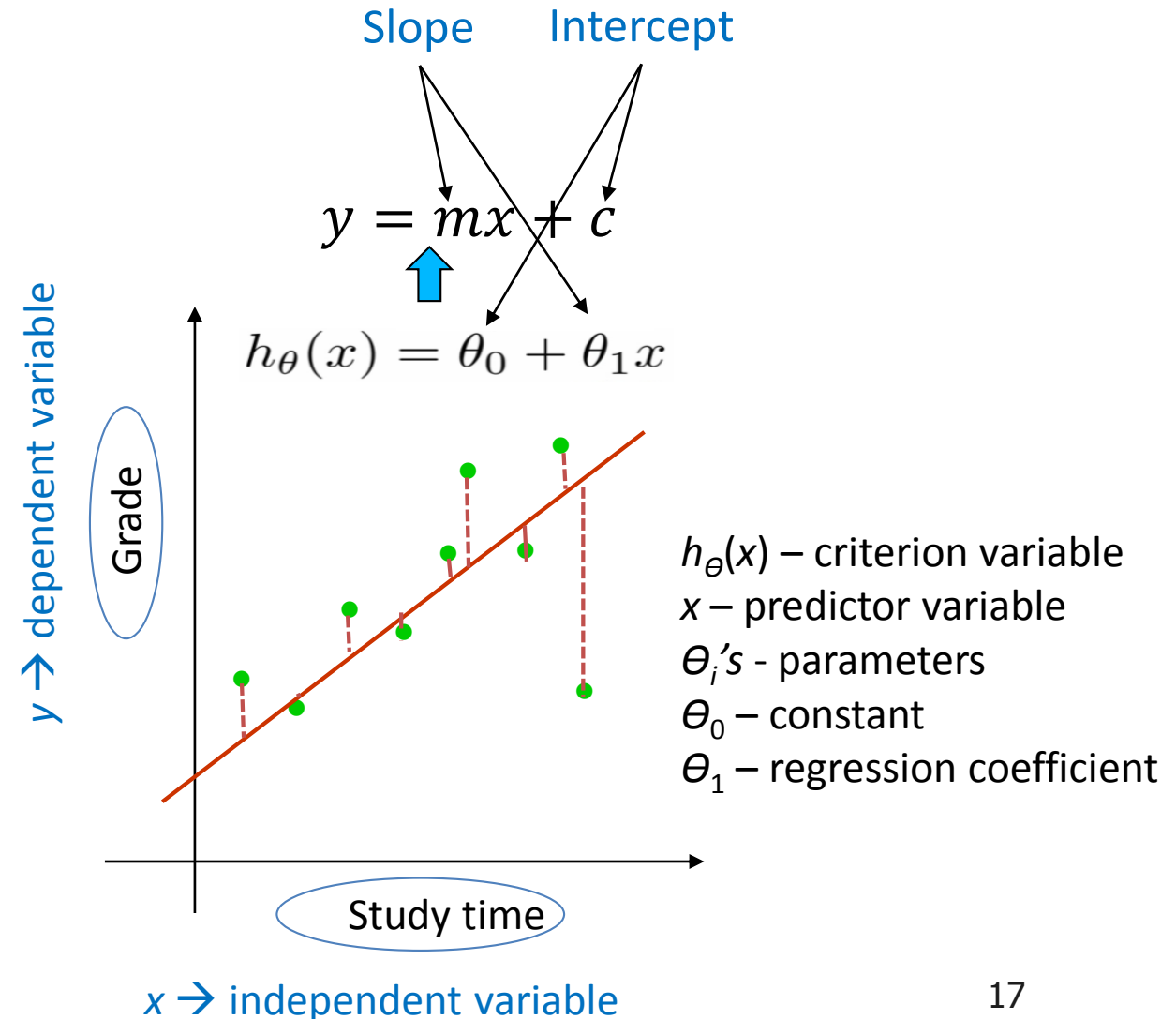
$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Objective Function / Cost Function

Where m = Number of data-points
(size of training set)

- **Goal** : Find θ 's such that Error is minimum (Find parameters θ_0 and θ_1 to minimize objective function)

→ **Learning**



Exercise

- Consider applying linear regression with the hypothesis as $h_{\theta}(x) = \theta_0 + \theta_1 x$
- The training data is given as

X	Y
6	7
5	4
10	8
3	4

- The cost function is

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

- What is value of $J(\theta)$, if $\theta = (2, 1)$?

- Hypothesis

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

- Parameters

$$\theta_0 \text{ and } \theta_1$$

- Cost Function

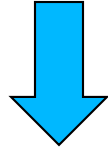
$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^i) - y^{(i)})^2$$

- Goal

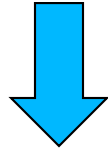
$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad J(\theta_0, \theta_1)$$

- Simplified Illustration with $\theta_0 = 0$

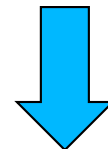
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



$$h_{\theta}(x) = \theta_1 x$$

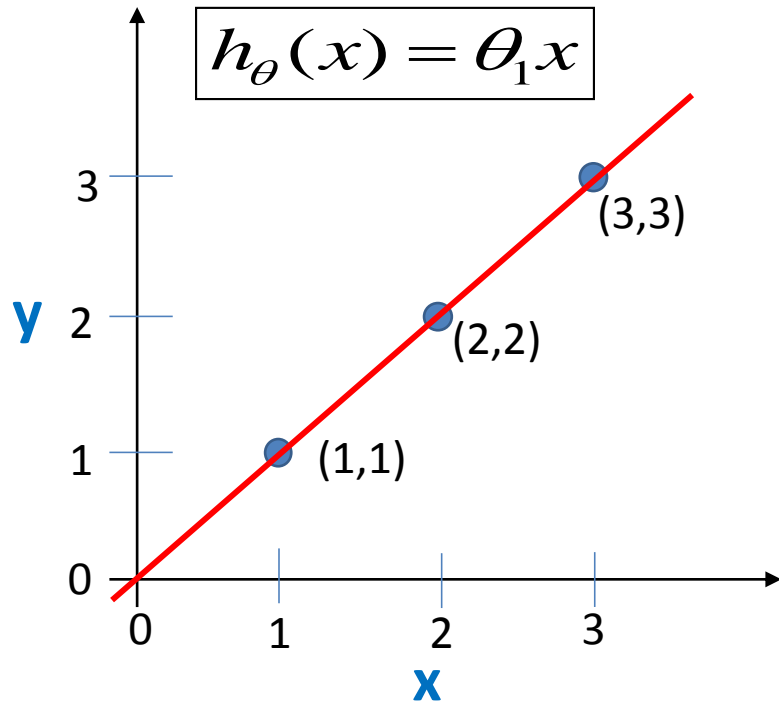


$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{2m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2$$



$$\min_{\theta_1} J(\theta_1)$$

Understanding Cost Function



$$\theta_0 = 0$$

$$\theta_1 = 1$$

$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$
$$= \frac{1}{2m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2$$

$$J(\theta_1) = J(1) = 0$$

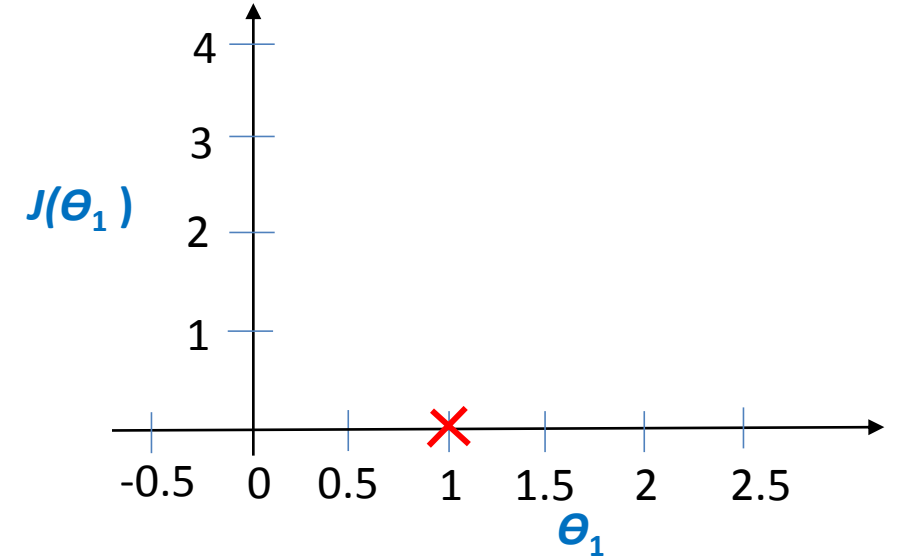


$$= \frac{1}{2 * 3} ((1 * 1 - 1)^2 + (1 * 2 - 2)^2 + (1 * 3 - 3)^2)$$
$$= 0$$

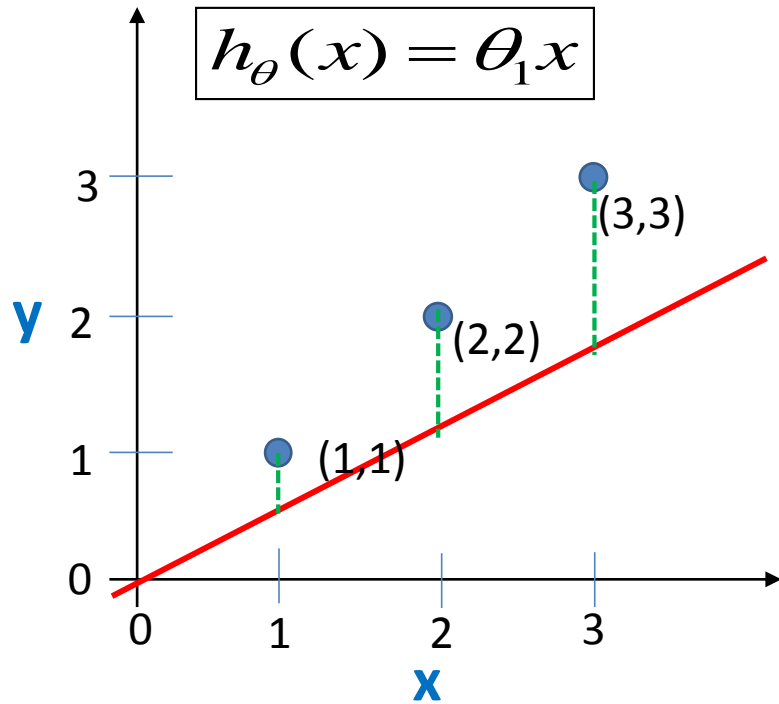
Cost Function

$$J(\theta_1)$$

Function of parameter θ_1



Understanding Cost Function



$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{2m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2$$

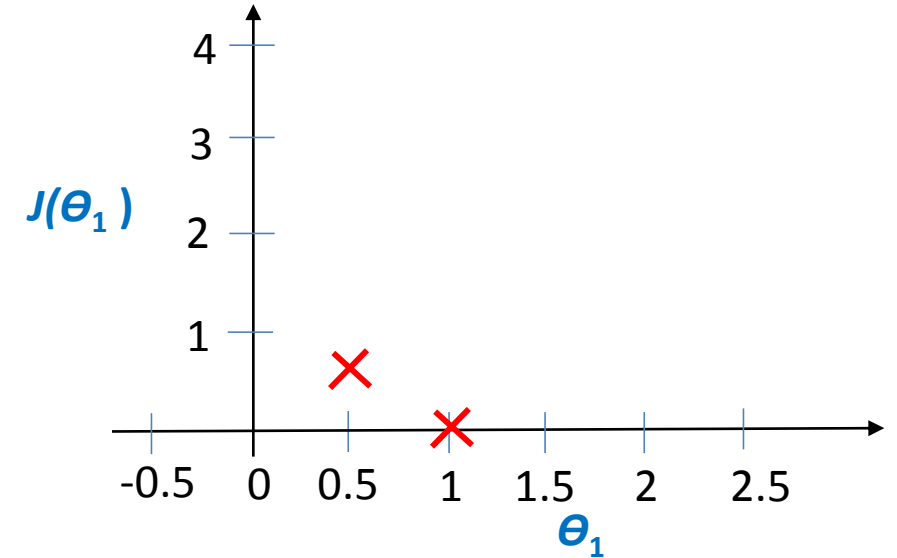
$$\begin{aligned} &= \frac{1}{2 \cdot 3} ((0.5 \cdot 1 - 1)^2 + (0.5 \cdot 2 - 2)^2 + (0.5 \cdot 3 - 3)^2) \\ &= \frac{1}{2 \cdot 3} ((-0.5)^2 + (-1)^2 + (-1.5)^2) = \frac{3.5}{6} \approx 0.58 \end{aligned}$$

$$\begin{aligned} J(\theta_1) &= \\ J(0.5) &= \\ 0.58 \end{aligned}$$

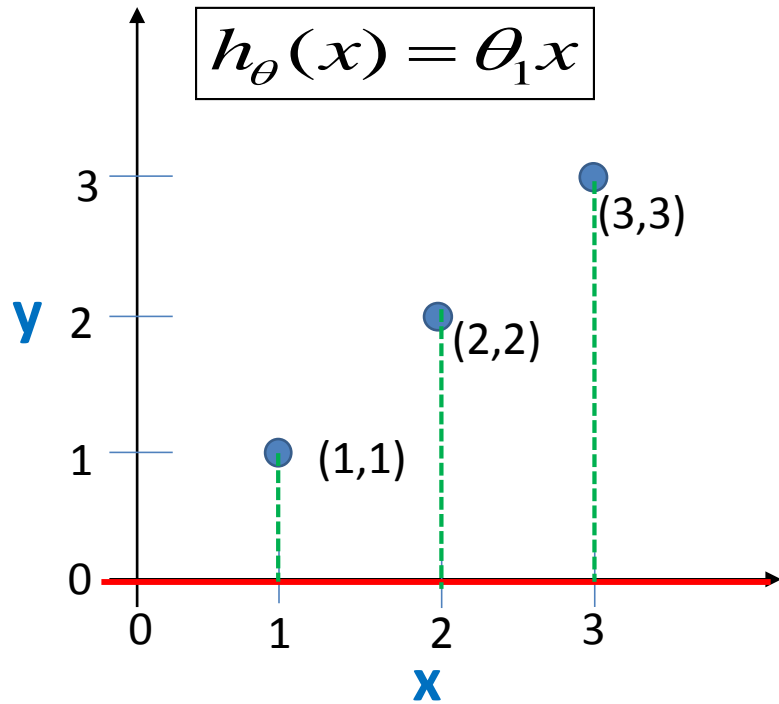


Cost Function $J(\theta_1)$

Function of parameter θ_1



Understanding Cost Function



$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{2m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2$$

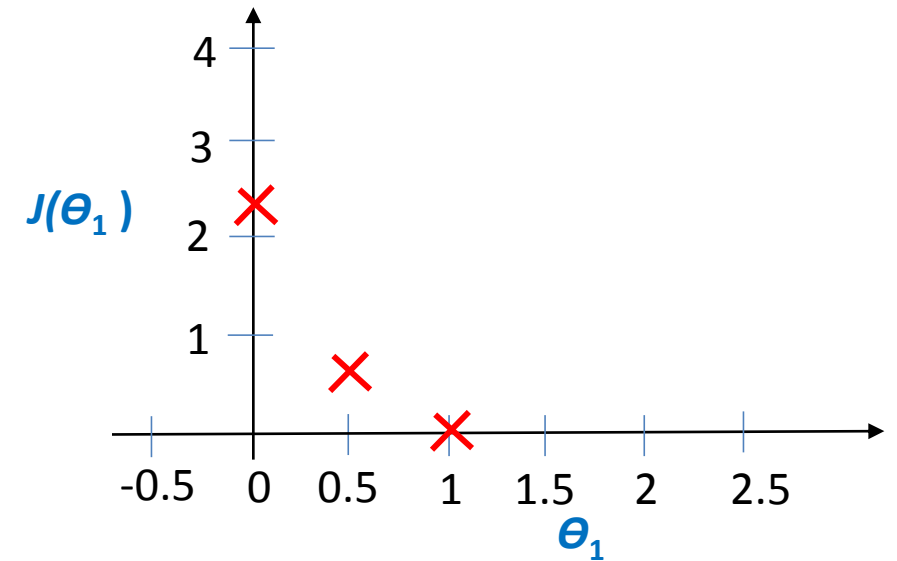
$$\begin{aligned} &= \frac{1}{2 \cdot 3} ((0 \cdot 1 - 1)^2 + (0 \cdot 2 - 2)^2 + (0 \cdot 3 - 3)^2) \\ &= \frac{1}{2 \cdot 3} ((-1)^2 + (-2)^2 + (-3)^2) = \frac{14}{6} \approx 2.33 \end{aligned}$$

$$\begin{aligned} J(\theta_1) &= \\ J(0) &= \\ 2.33 \end{aligned}$$

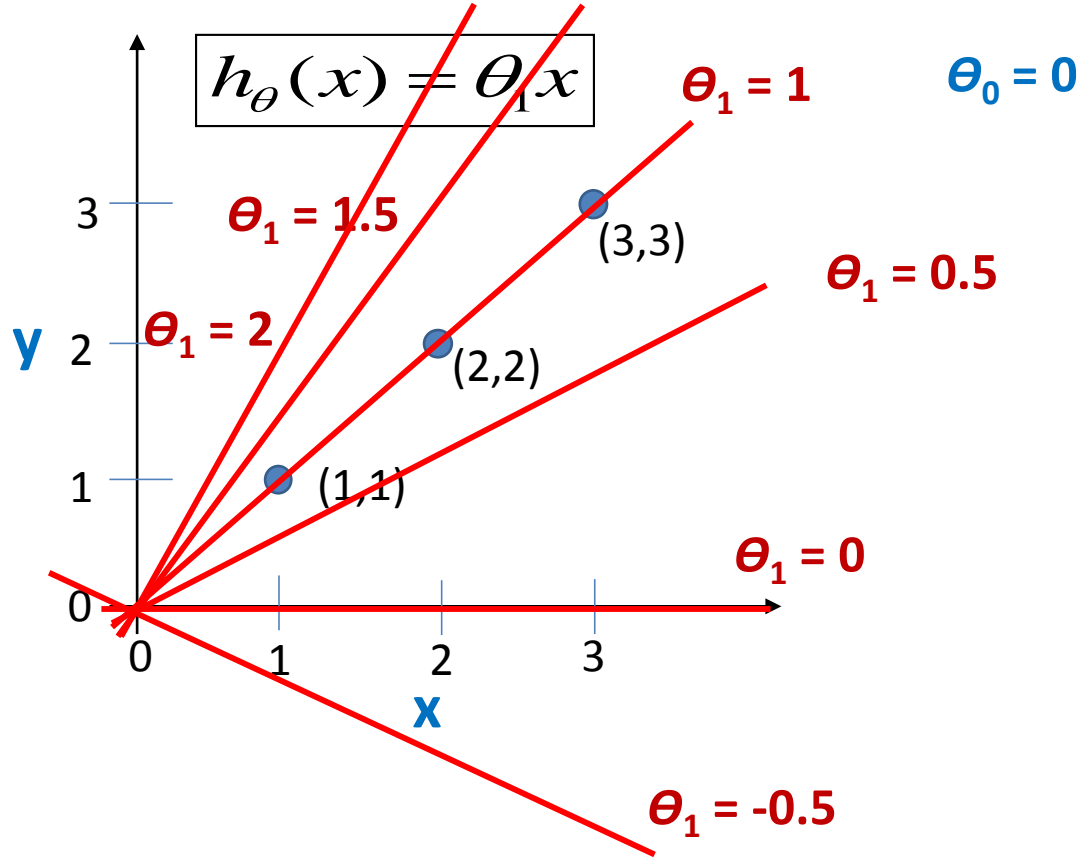


Cost Function $J(\theta_1)$

Function of parameter θ_1



Understanding Cost Function

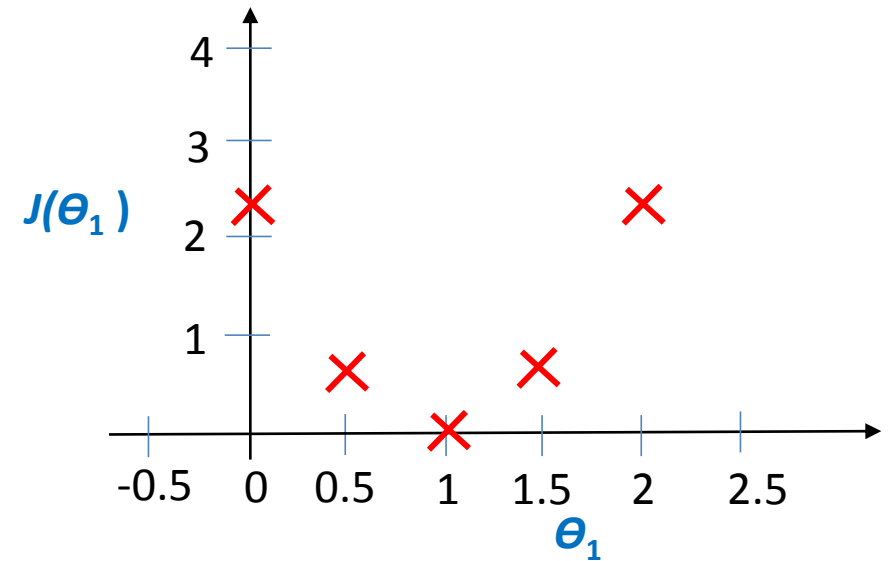


$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{2m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2$$

Cost Function

$$J(\theta_1)$$

Function of parameter θ_1



Cost Optimization

Finding θ_1 to minimize $J(\theta_1)$

$$\theta_1 = 1$$

- Hypothesis

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

- Parameters

$$\theta_0 \text{ and } \theta_1$$

- Cost Function

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

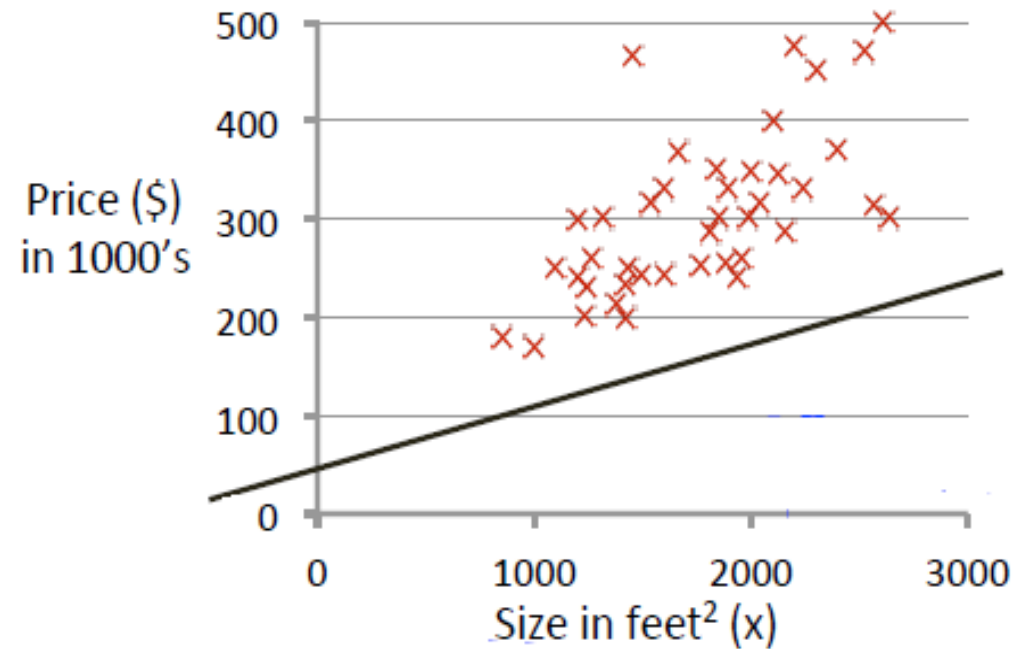
- Goal

$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad J(\theta_0, \theta_1)$$

Understanding Cost Function

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

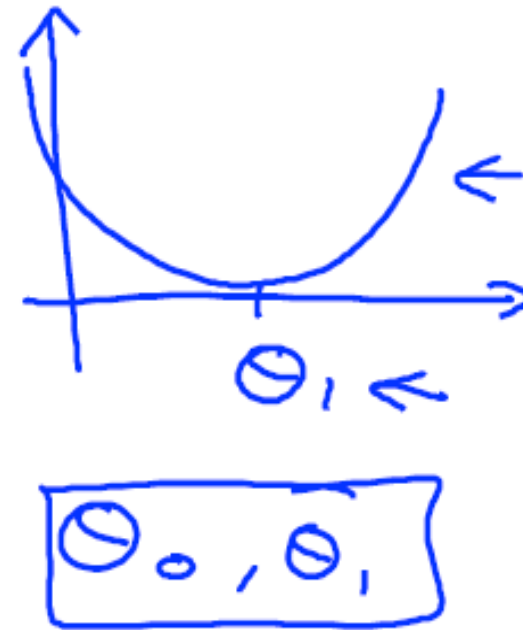
(for fixed θ_0, θ_1 , this is a function of x)



$$h_{\theta}(x) = 50 + 0.06x$$

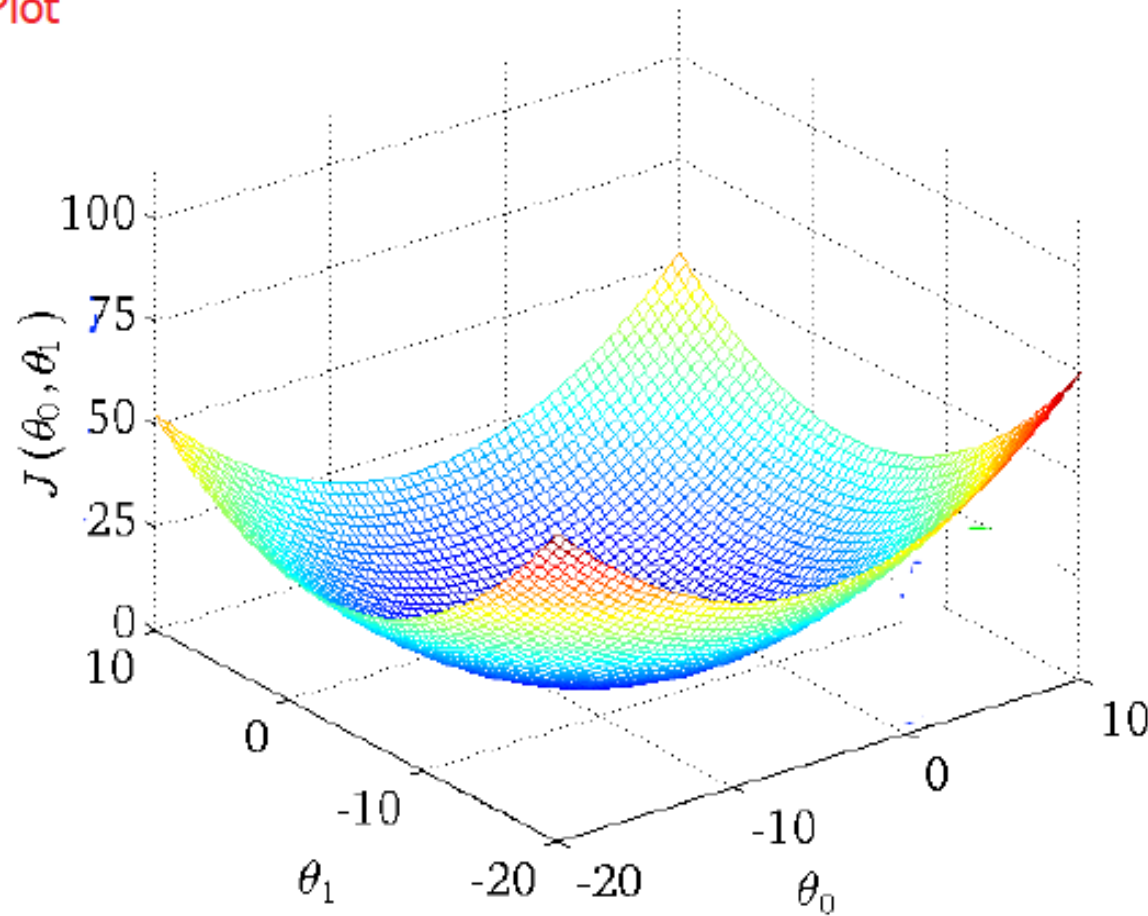
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



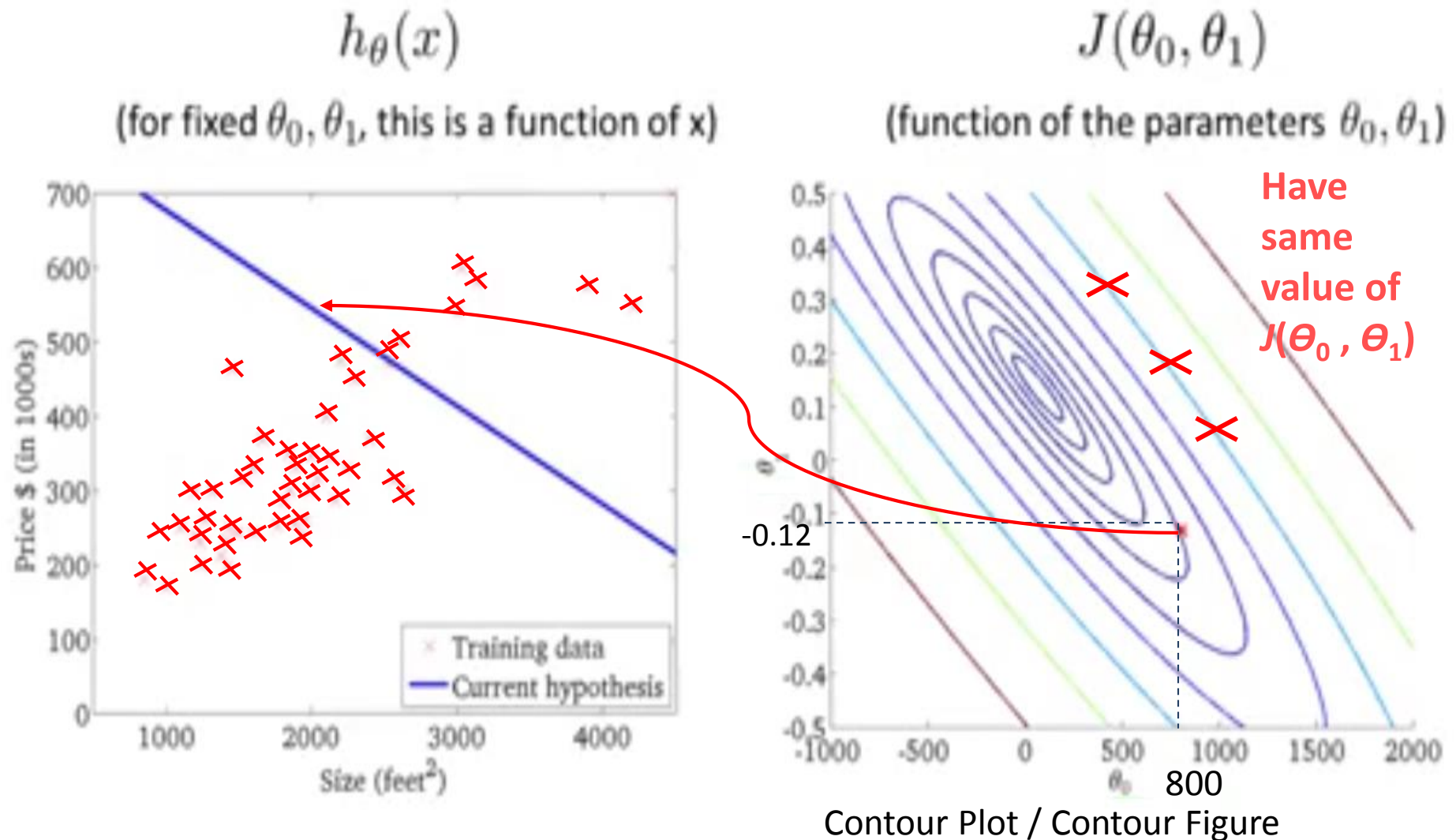
Understanding Cost Function

Surface Plot



Andrew Ng

Understanding Cost Function

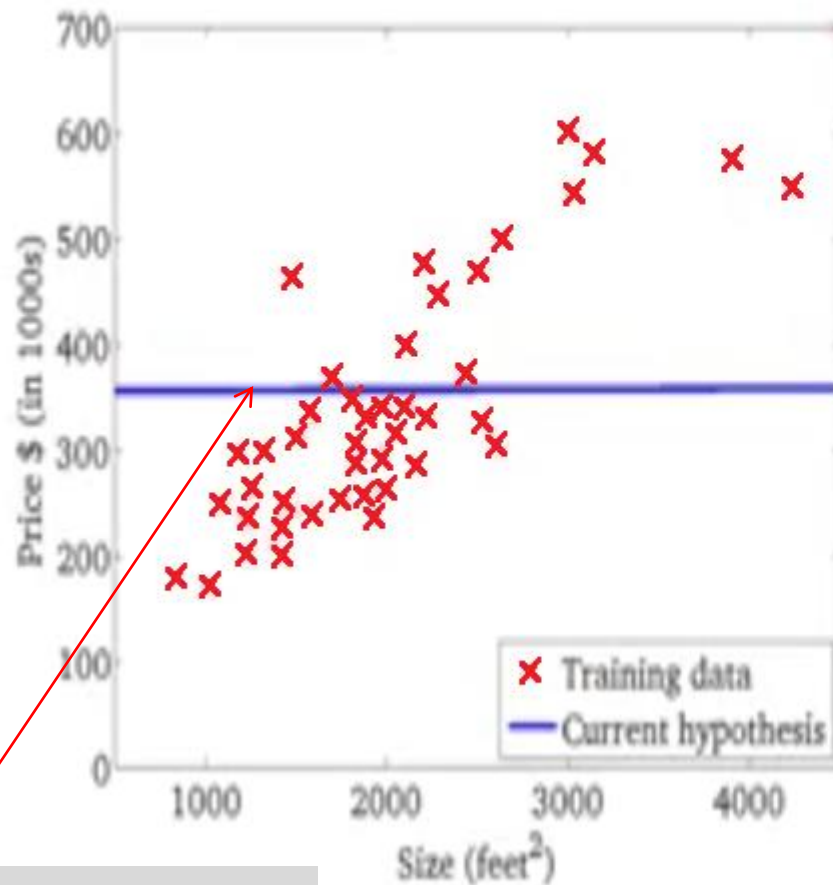


Elliptical contours show the set of points which takes on the same value of $J(\theta_0, \theta_1)$

Understanding Cost Function

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

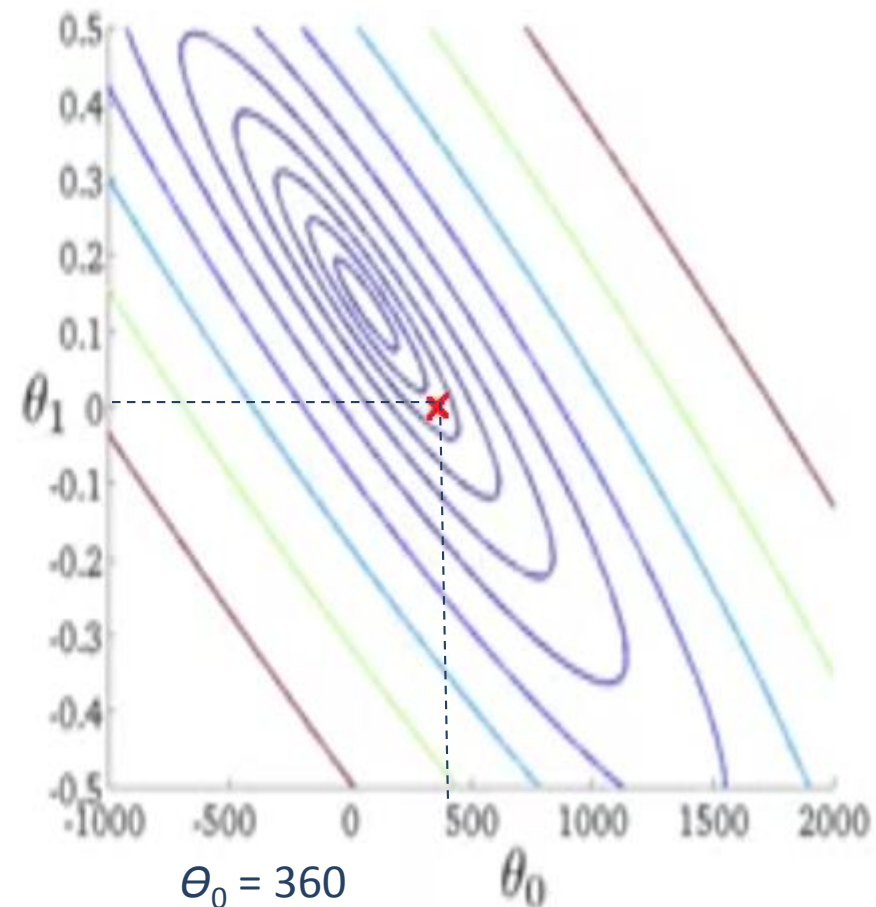
(for fixed θ_0, θ_1 , this is a function of x)



$$h_{\theta}(x) = 360 + 0 * x$$

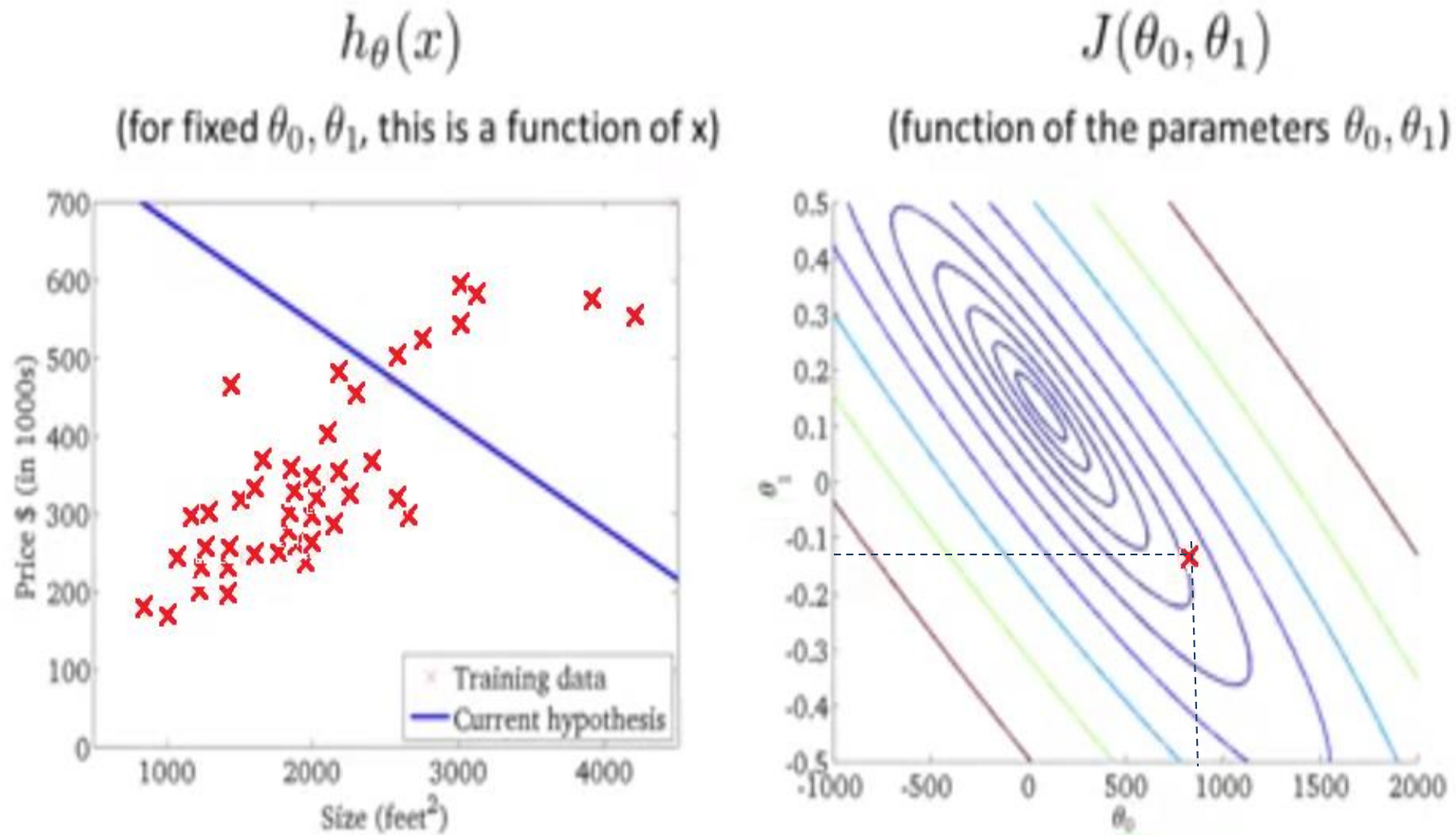
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)

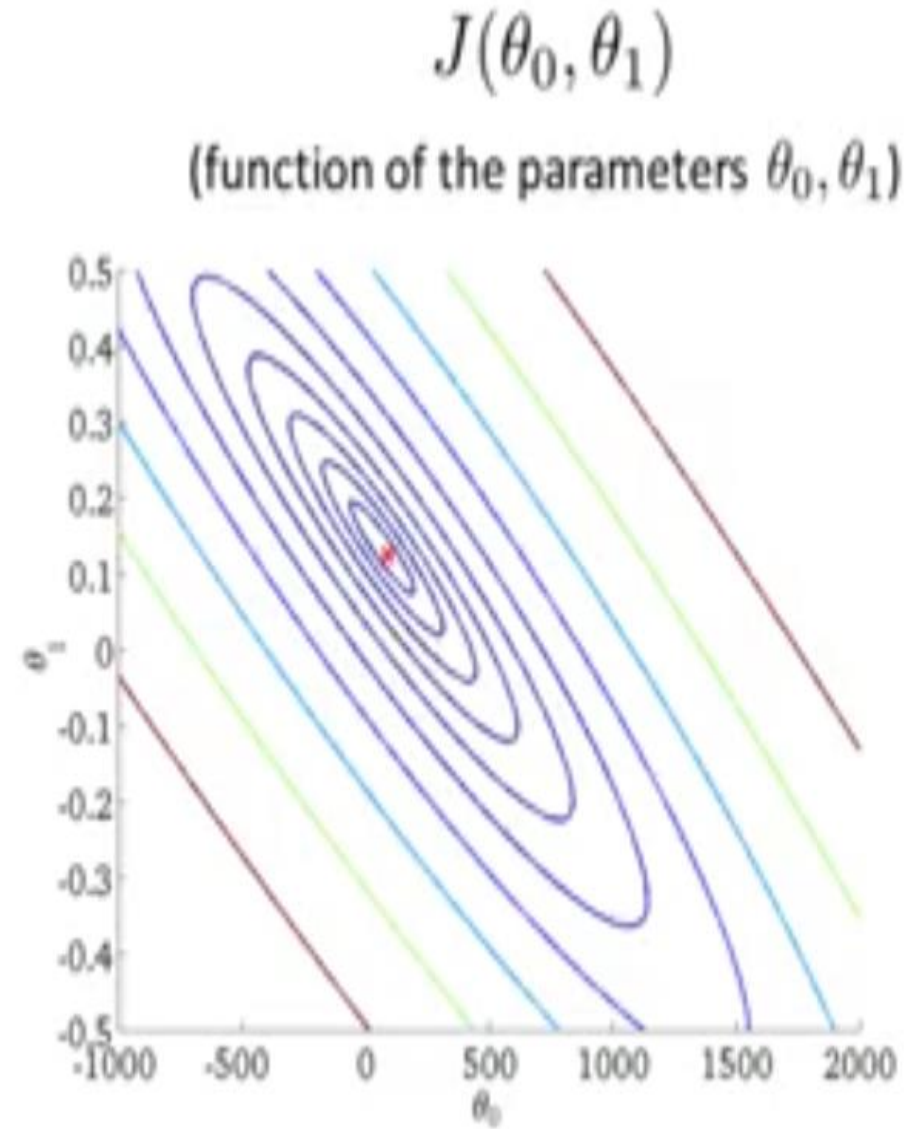
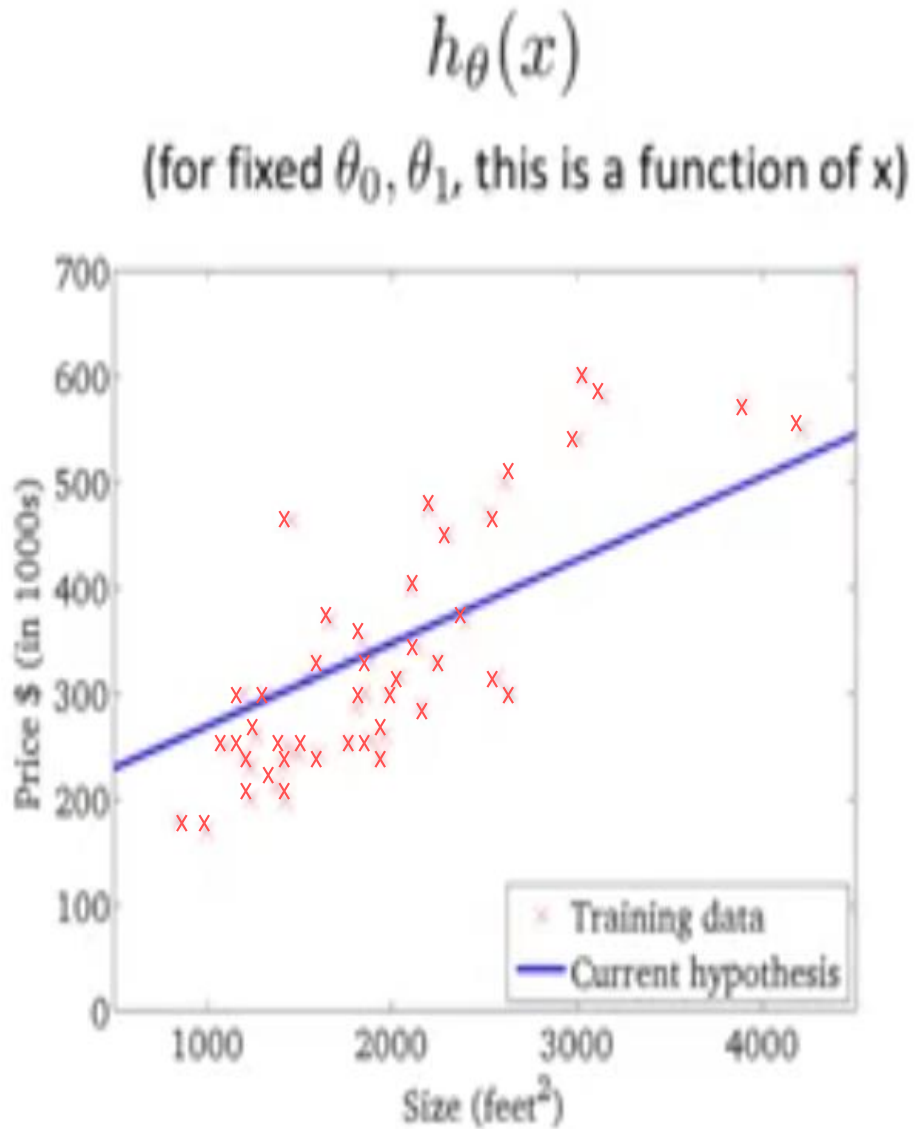


$$\theta_0 = 360$$
$$\theta_1 = 0$$

Understanding Cost Function



Understanding Cost Function



Unit 2

Gradient Descent

Machine Learning – Regression and Decision Trees

- Hypothesis

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

- Parameters

$$\theta_0 \text{ and } \theta_1$$

- Cost Function

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

- Goal

$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad J(\theta_0, \theta_1)$$

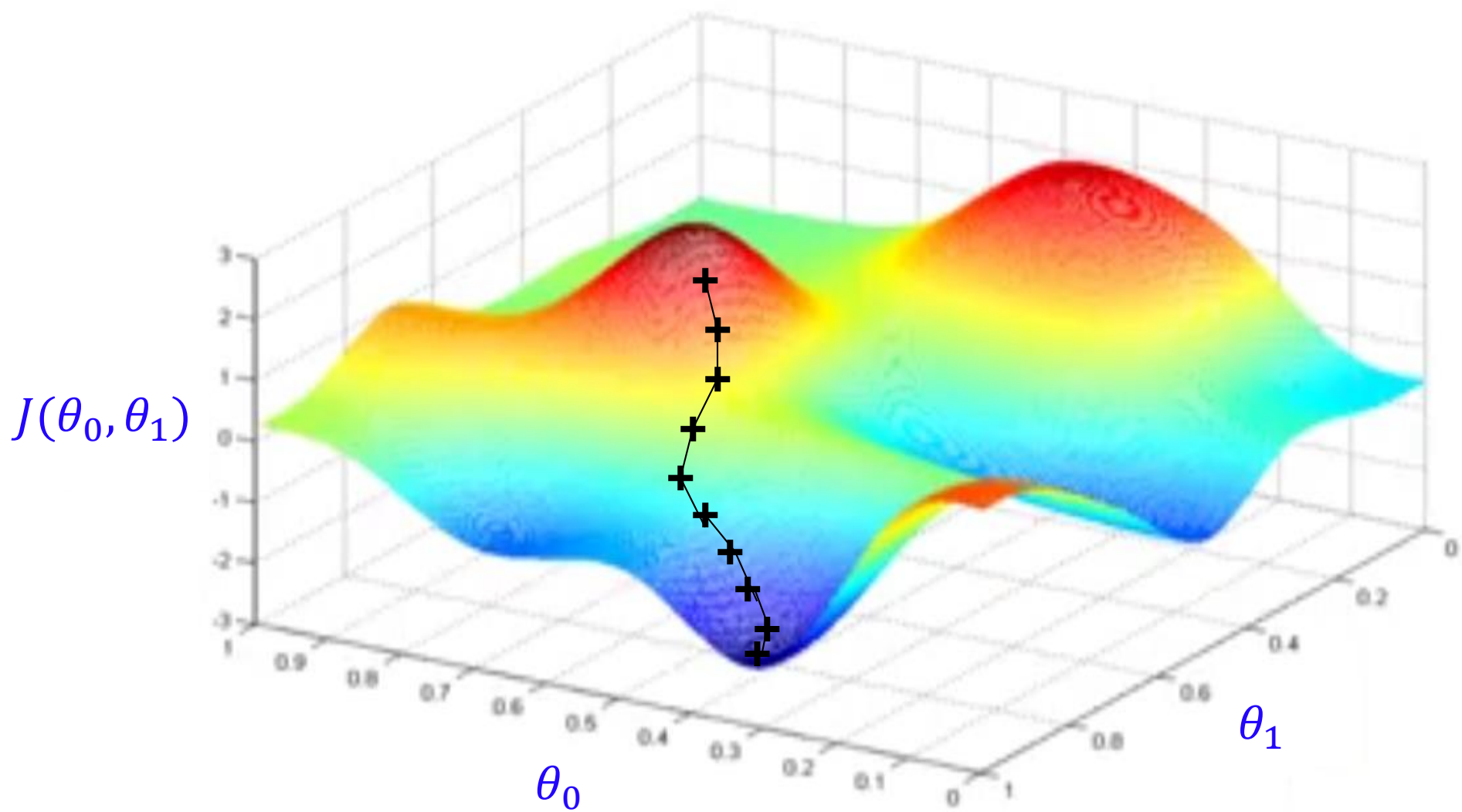
Gradient Descent

Have some function $J(\theta_0, \theta_1)$

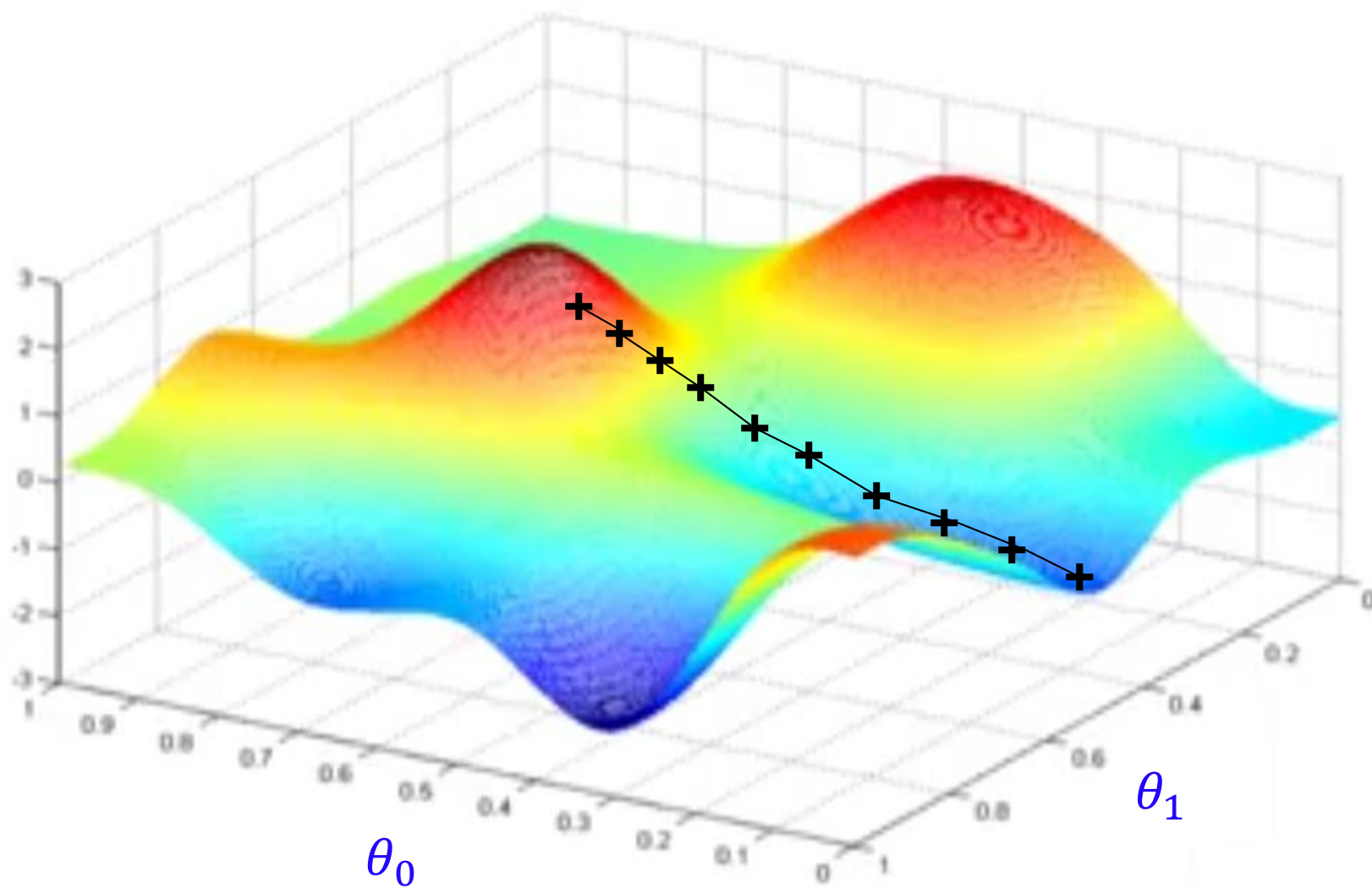
Want $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

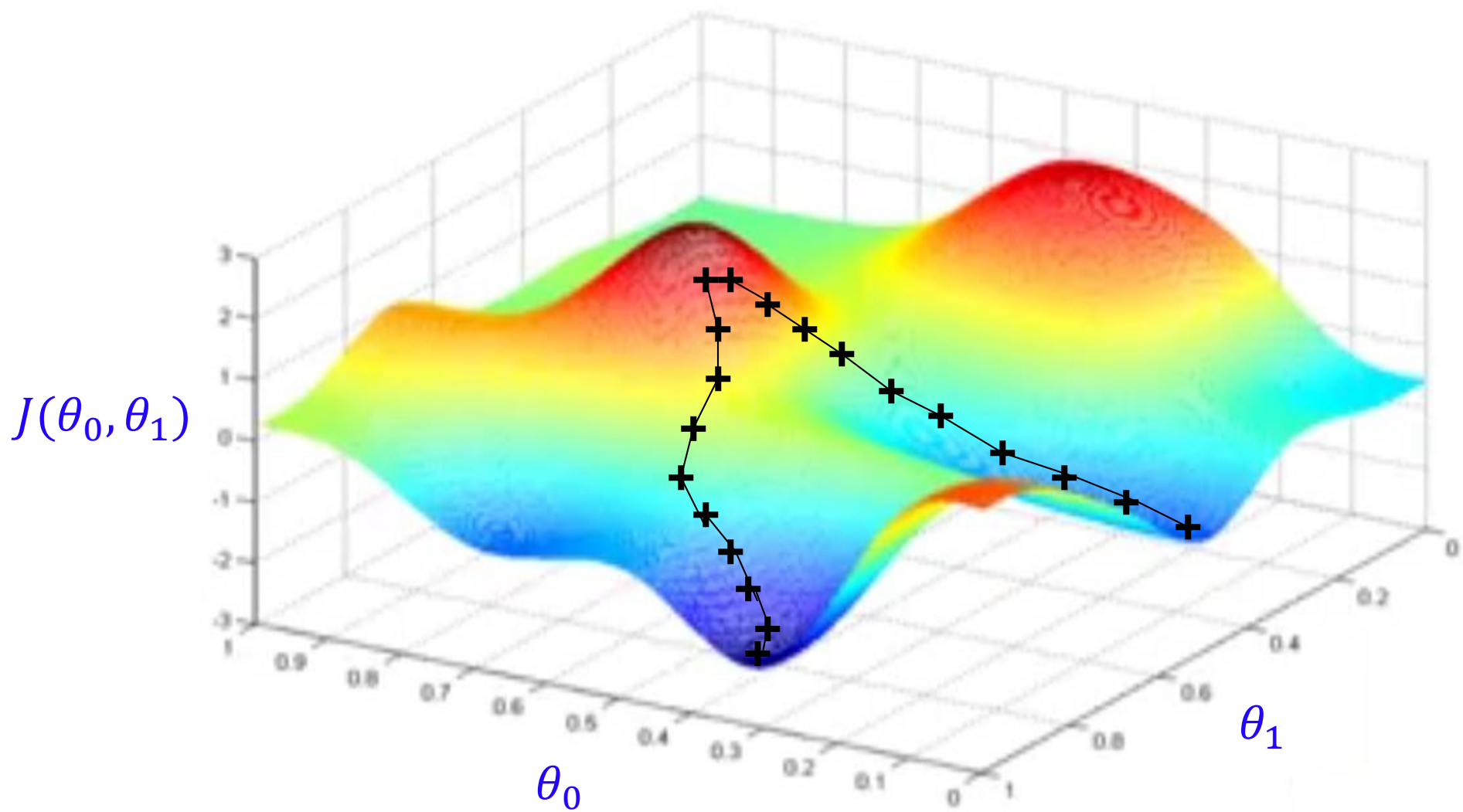
Outline:

- Start with some θ_0, θ_1
- Keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$
until we hopefully end up at a minimum



$J(\theta_0, \theta_1)$





Unit 2

19

Gradient Descent for Linear Regression

Machine Learning – Regression and Decision Trees

Gradient Descent Algorithm

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad \text{for } (j = 0 \text{ and } j = 1)$$

Simultaneous update
of θ_0, θ_1

}

Learning Rate

Derivative

Correct Implementation: Simultaneous Update

$$temp0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$temp1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := temp0$$

$$\theta_1 := temp1$$

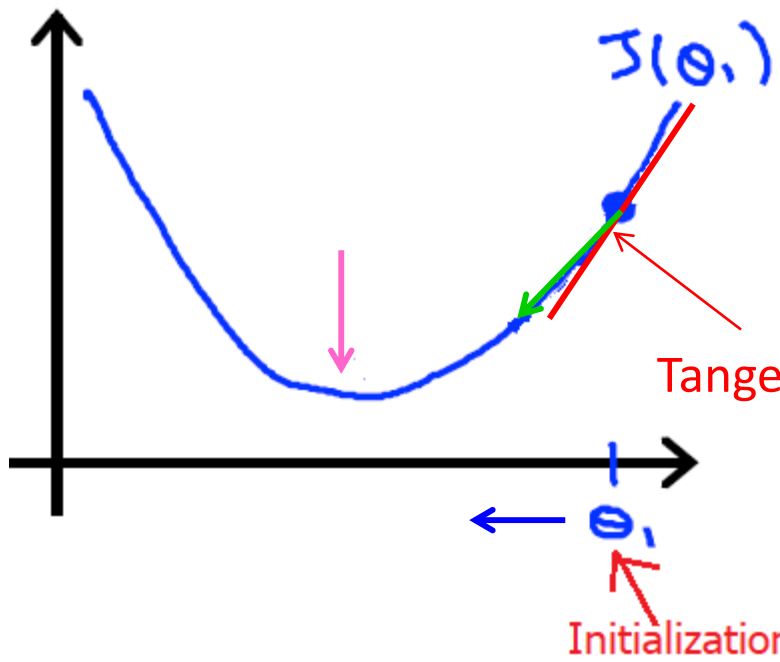
Incorrect Implementation

$$temp0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\theta_0 := temp0$$

$$temp1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

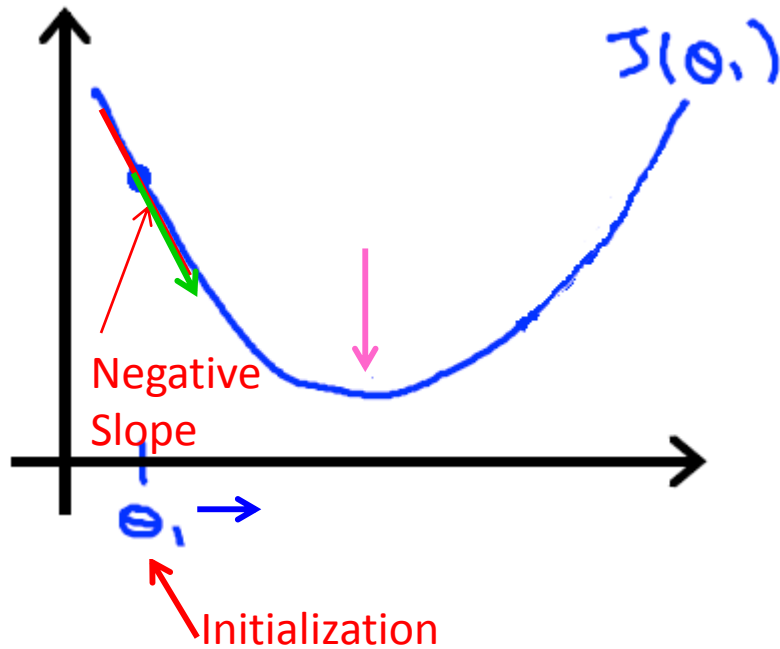
$$\theta_1 := temp1$$



$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

Slope of Tangent at θ_1
 ≥ 0

$$\theta_1 := \theta_1 - \alpha * (\text{positive number})$$



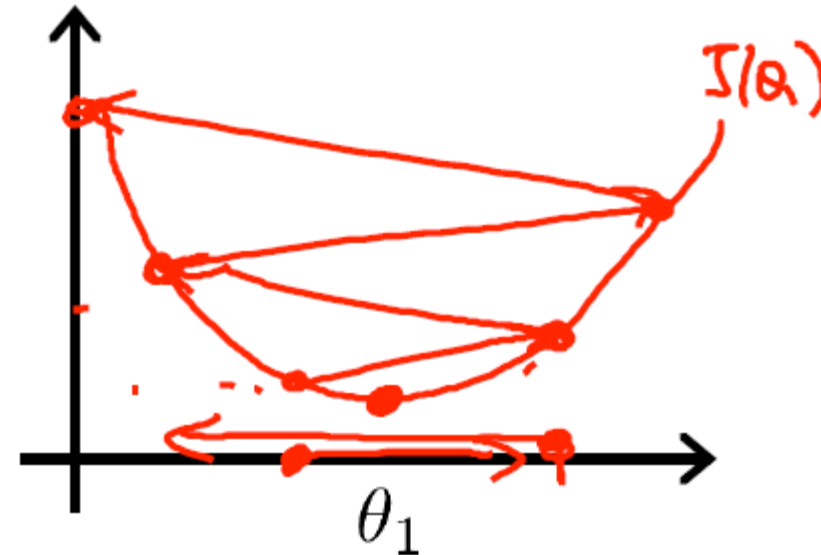
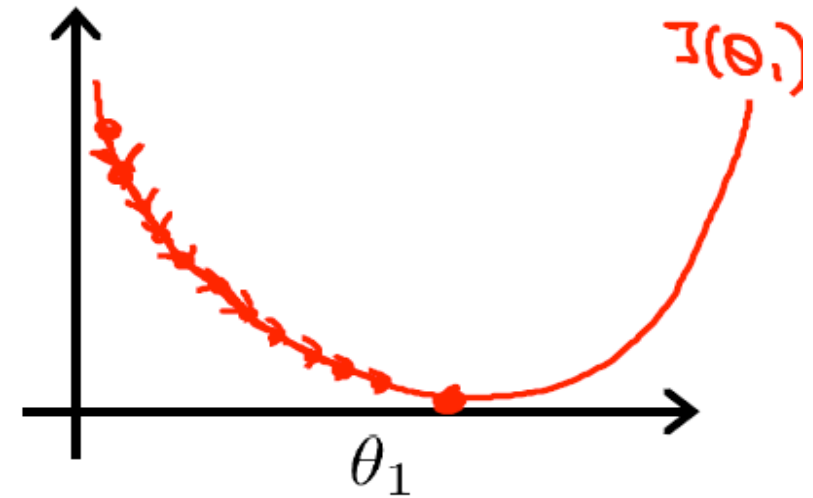
$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

$$\theta_1 := \theta_1 - \alpha * (\text{negative number})$$

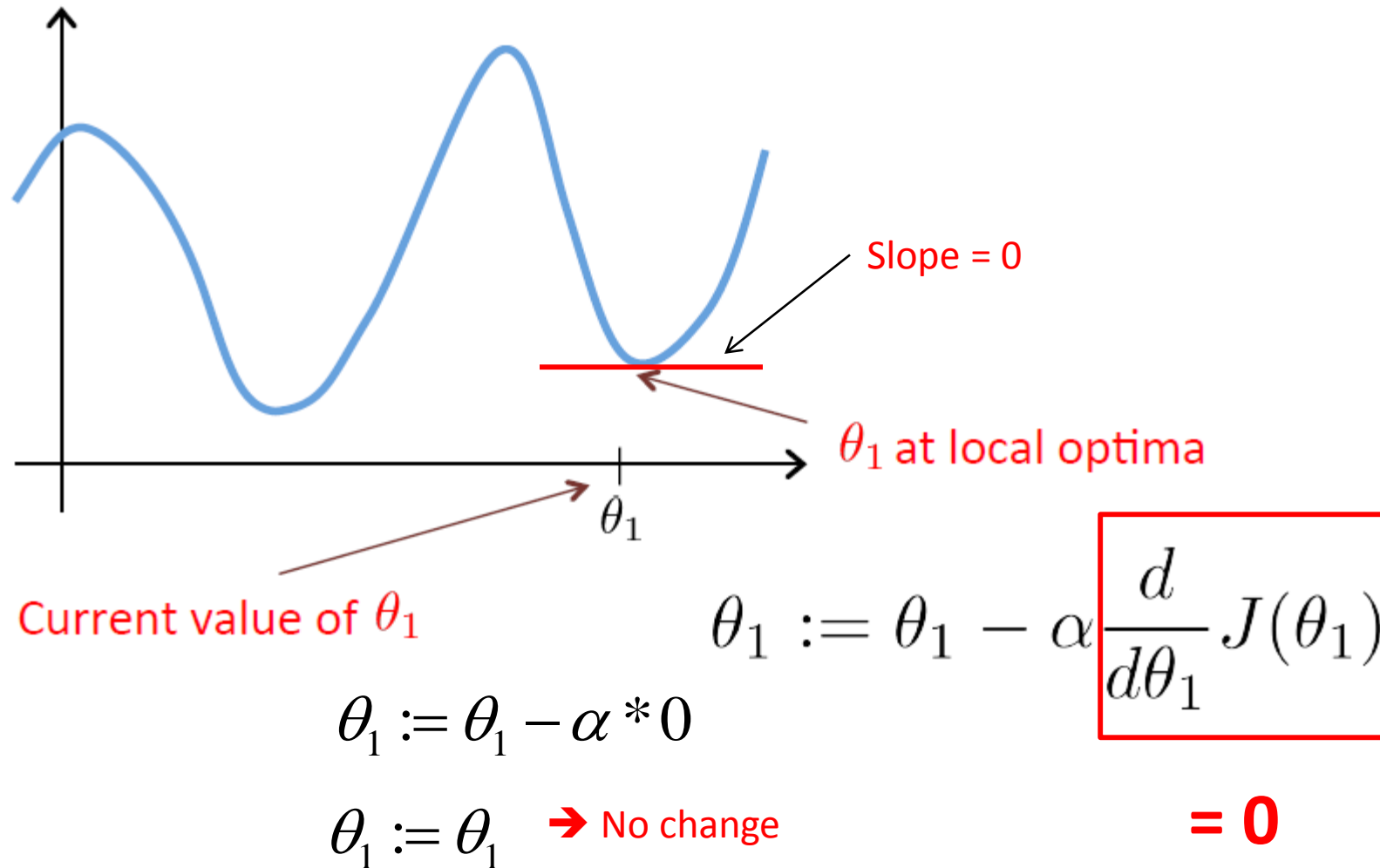
The Learning Rate

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta} J(\theta_1)$$

- If α is **too small**, gradient descent can be slow
- If α is **too large**, gradient descent can overshoot the minimum. It may fail to converge, or even diverge



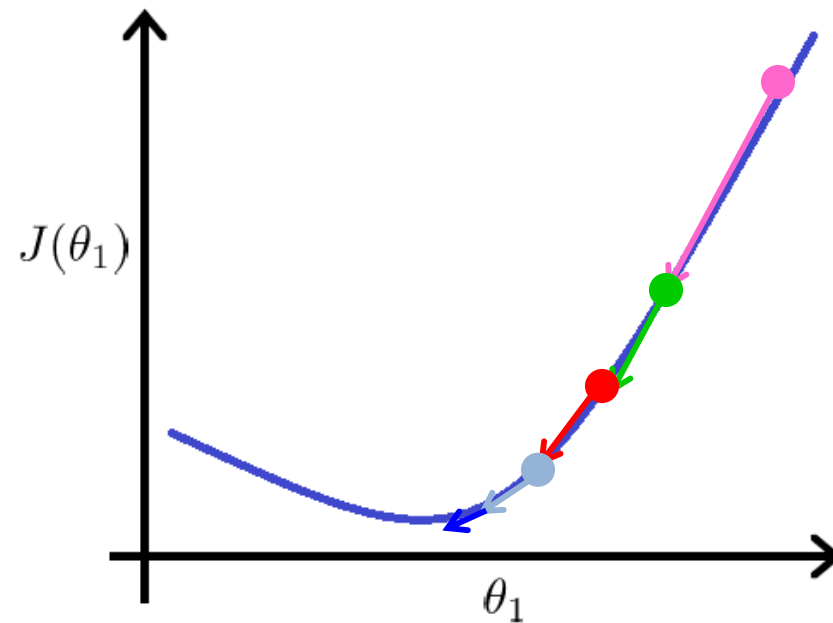
- What if θ_1 is already at local minimum ?



Gradient descent can converge to a local minimum, even with the learning rate α fixed.

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

As we approach a local minimum, gradient descent will automatically take smaller steps. So, no need to decrease α over time.



Gradient Descent for Linear Regression

Gradient Descent Algorithm

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

for ($j = 0$ and $j = 1$)

}

Linear Regression Model

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Apply Gradient Descent to minimize cost function of linear regression

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$

$$j = 0: \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

$$j = 1:$$

Gradient Descent Algorithm for Linear Regression

repeat until convergence {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

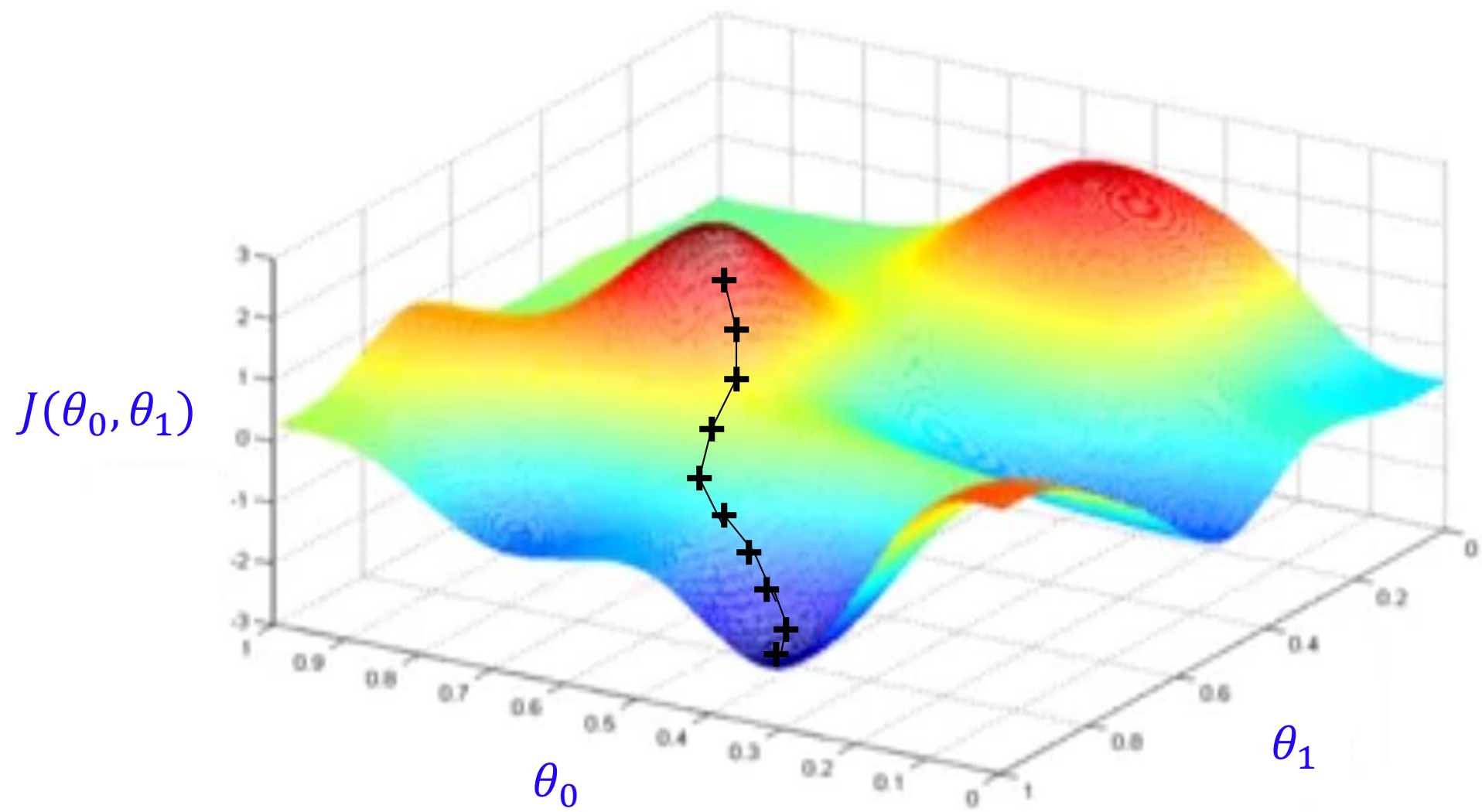
$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

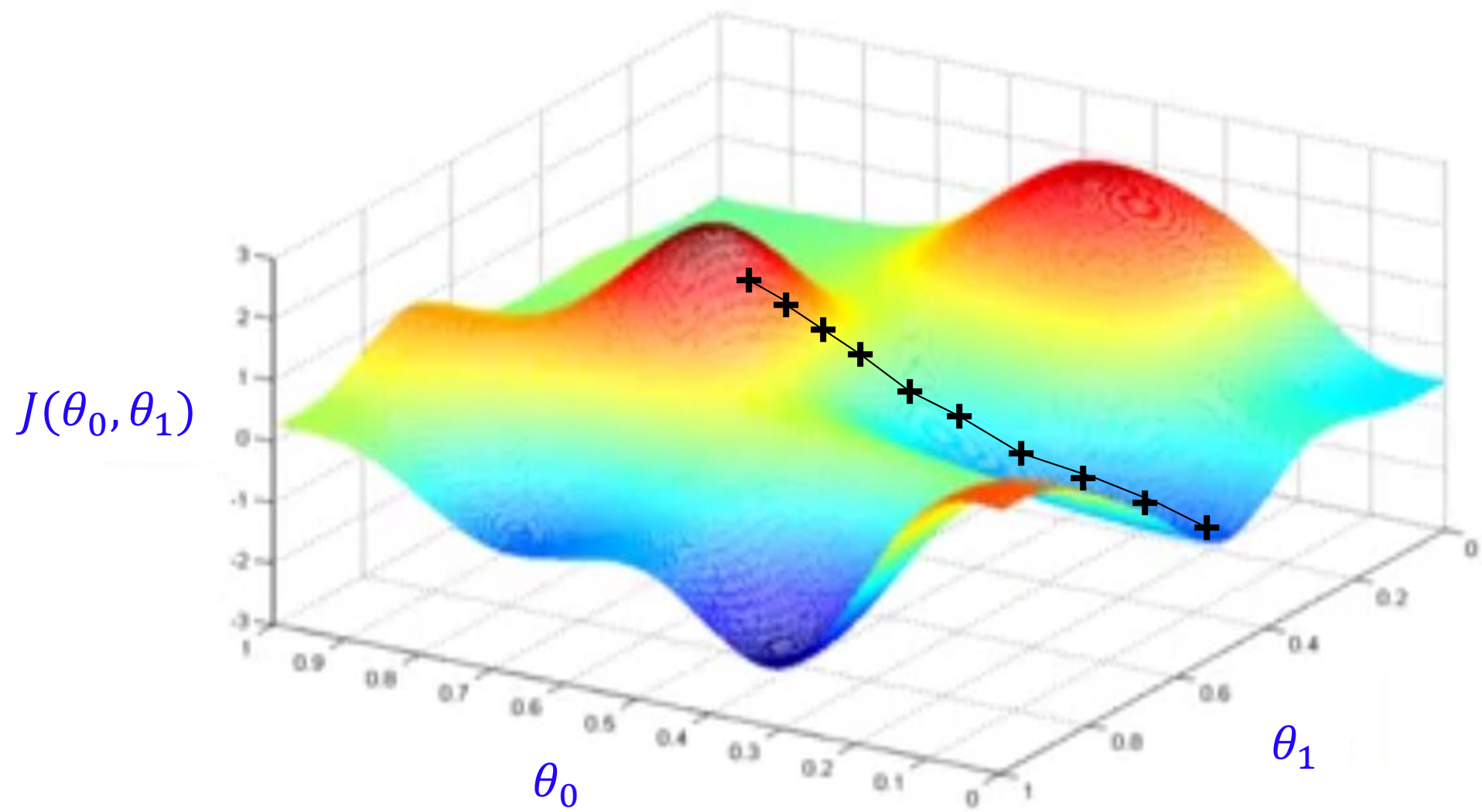
Update θ_0 and θ_1
simultaneously

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

}

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

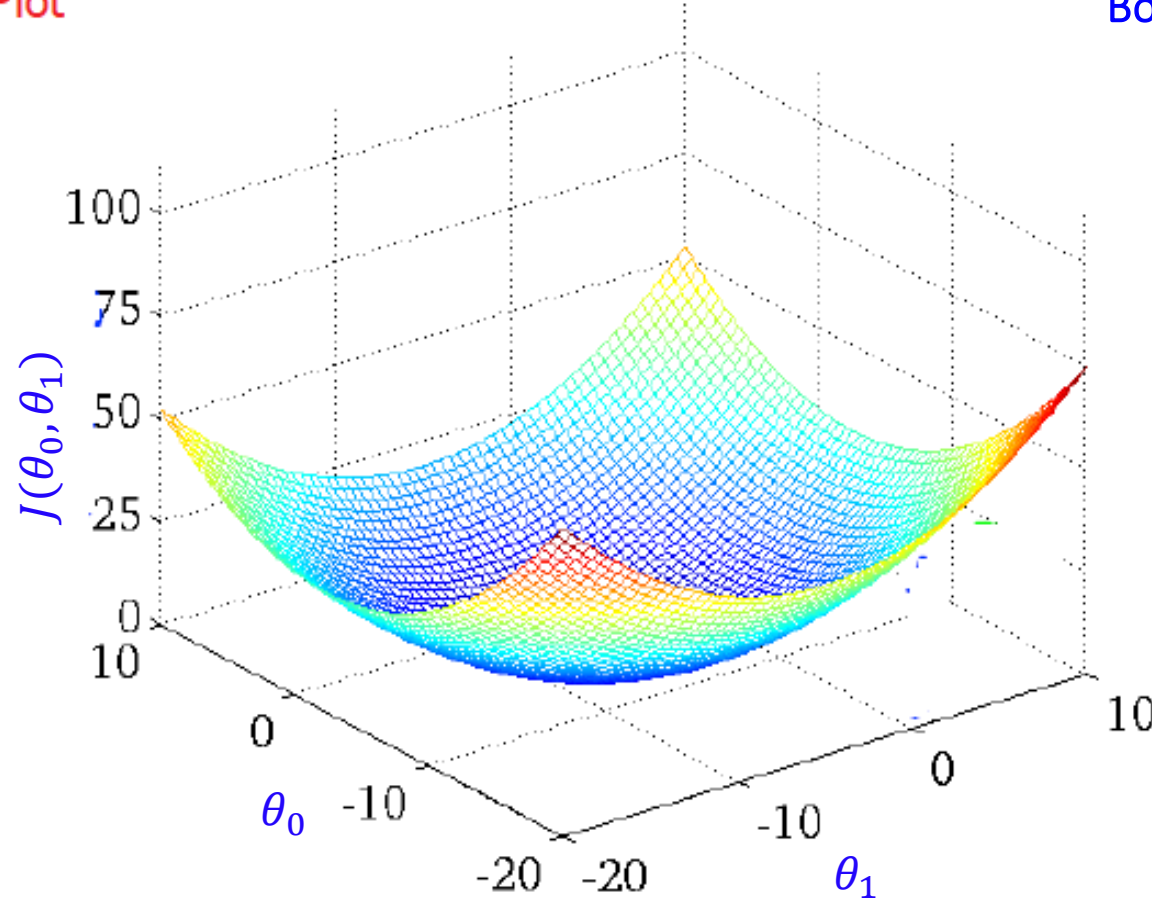




Surface Plot

Convex Function

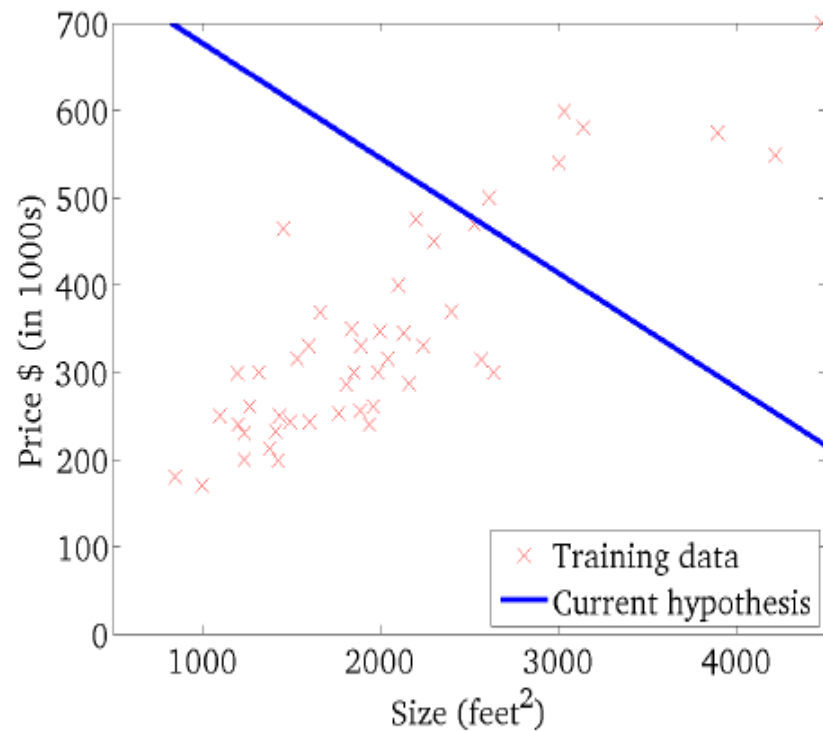
Bowl-shaped Function



- Cost function of linear regression has always one global minima and not the local minima.

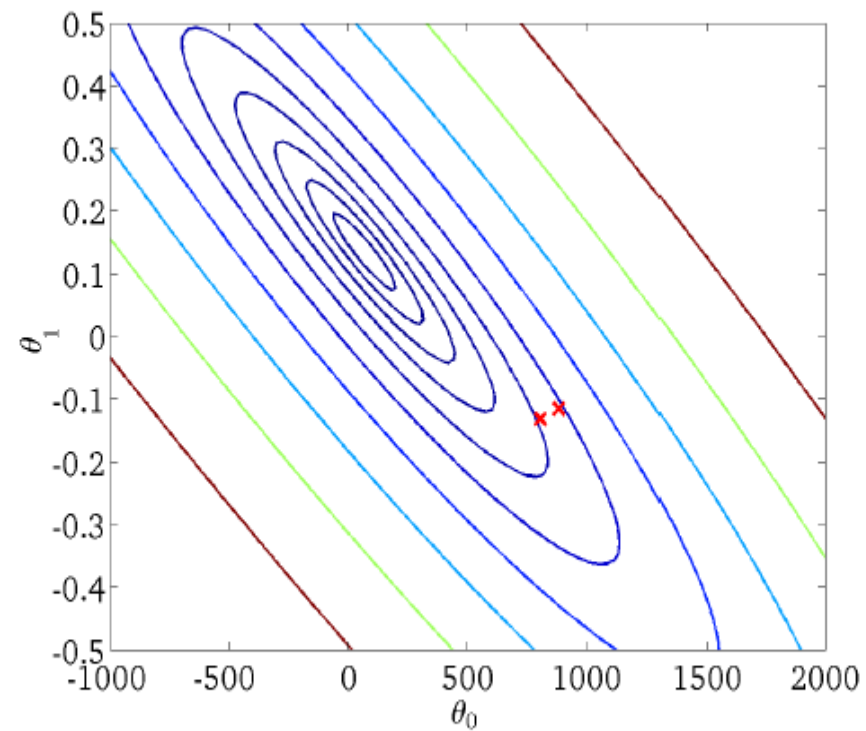
$$h_{\theta}(x)$$

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



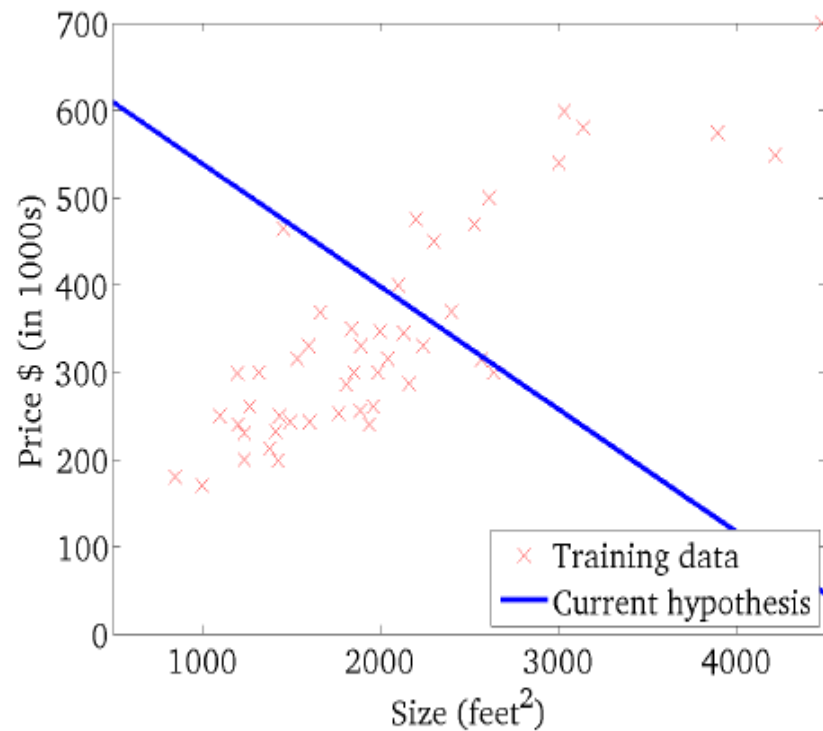
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



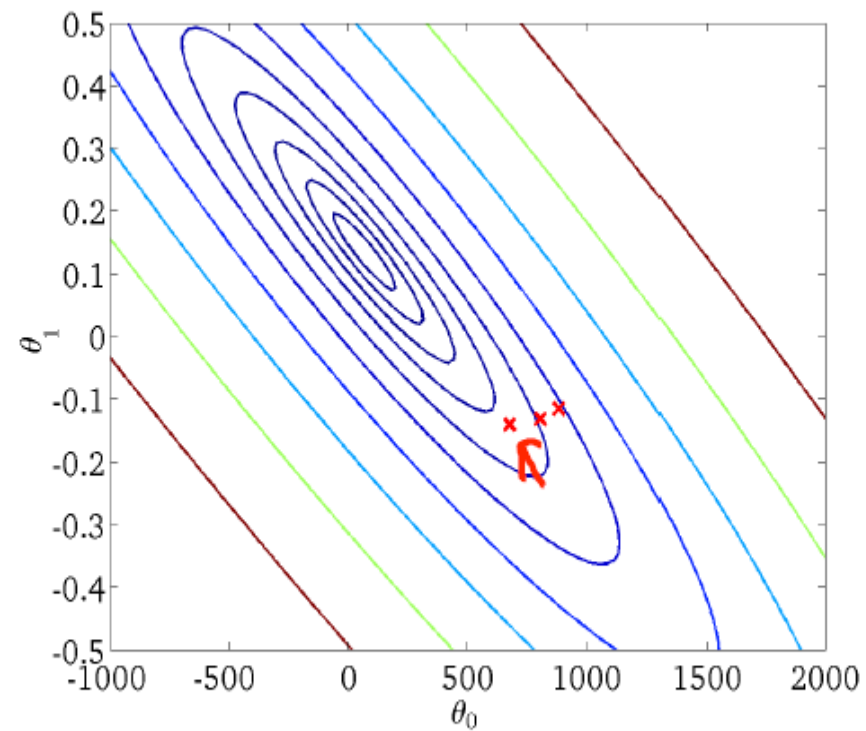
$$h_{\theta}(x)$$

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



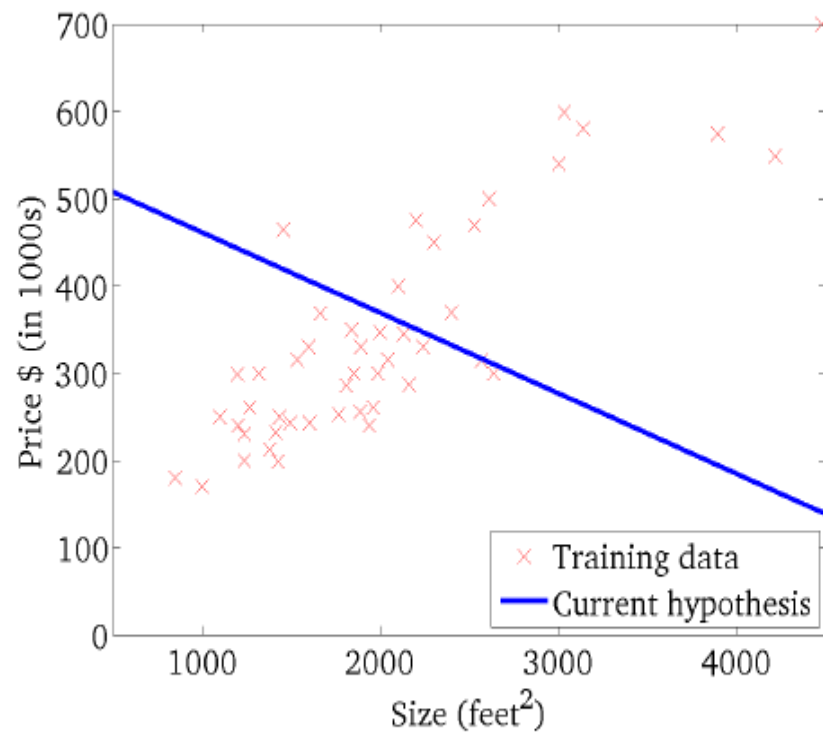
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



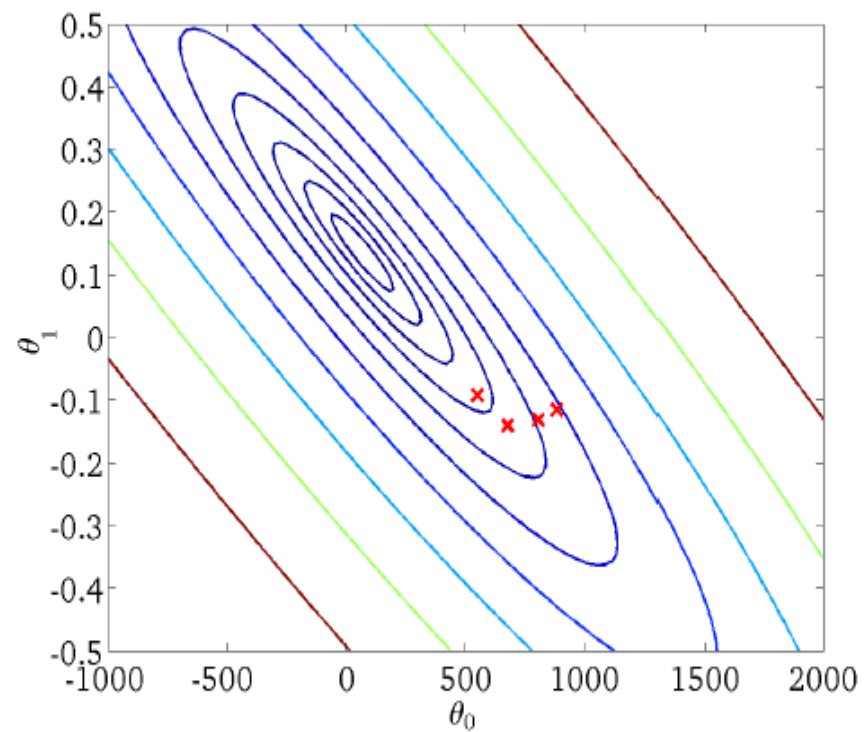
$$h_{\theta}(x)$$

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



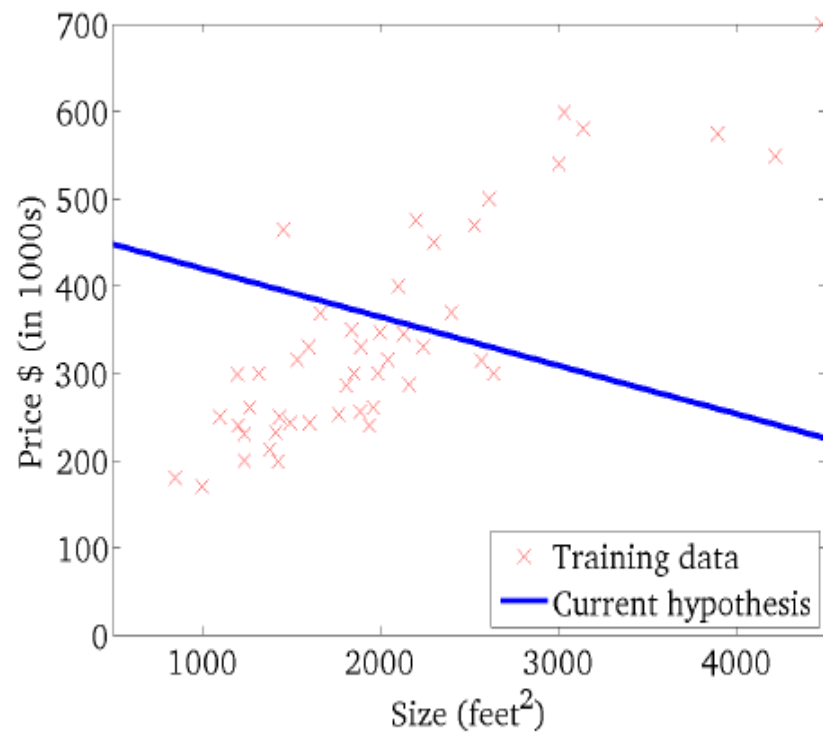
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



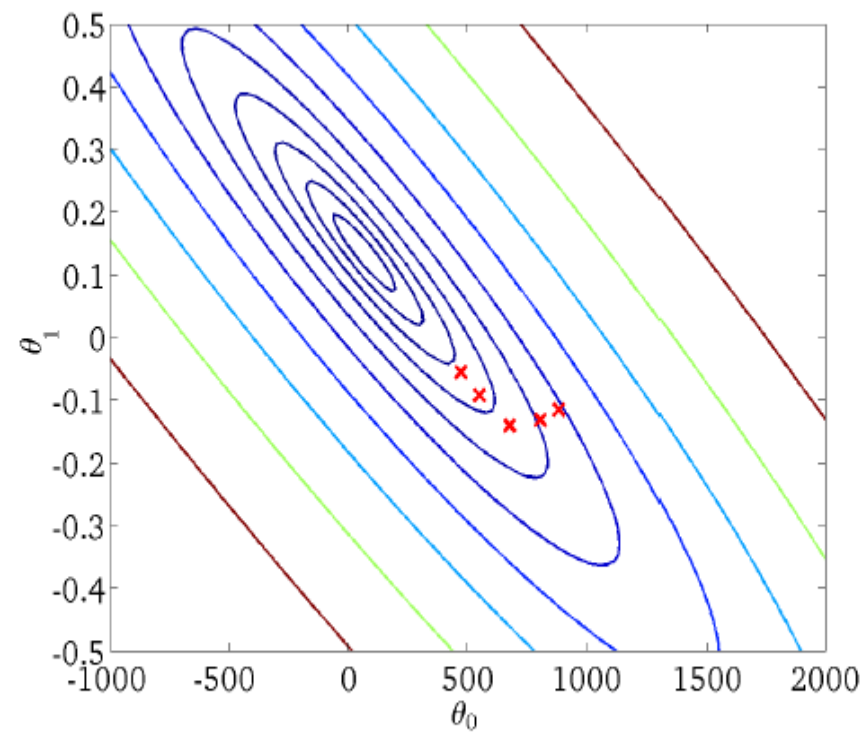
$$h_{\theta}(x)$$

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



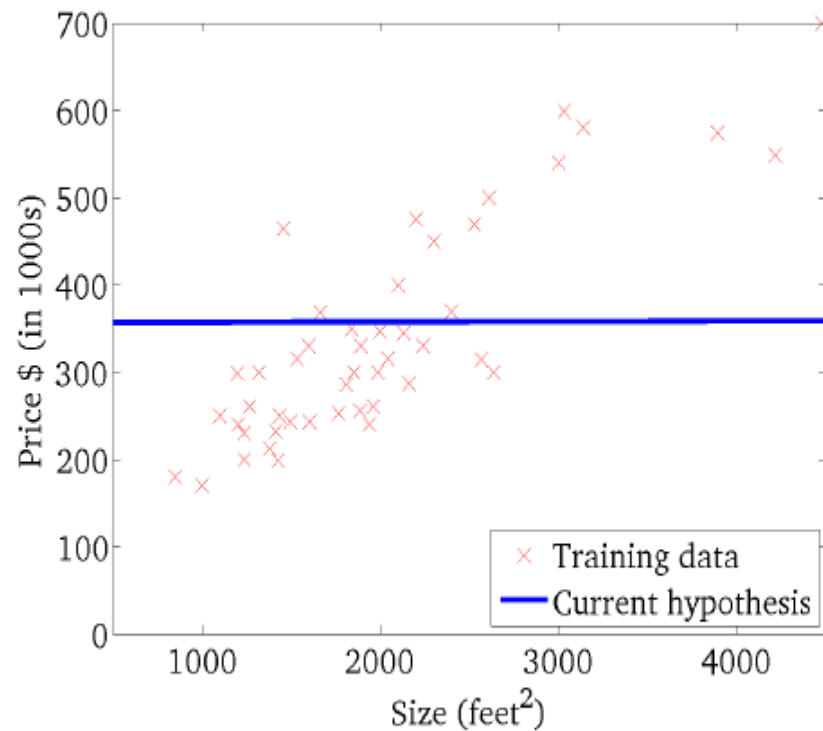
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



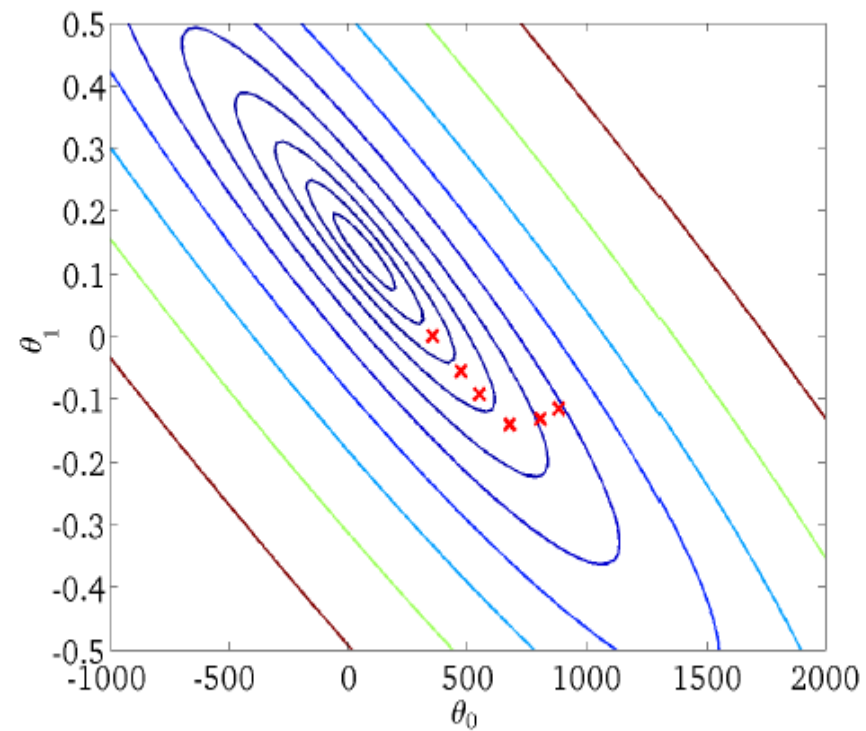
$$h_{\theta}(x)$$

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



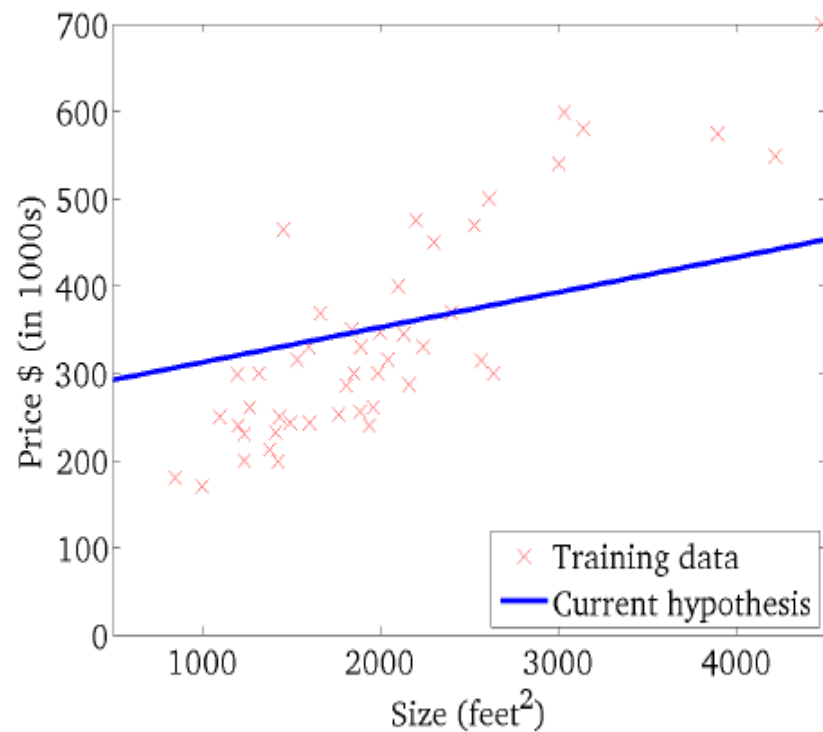
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



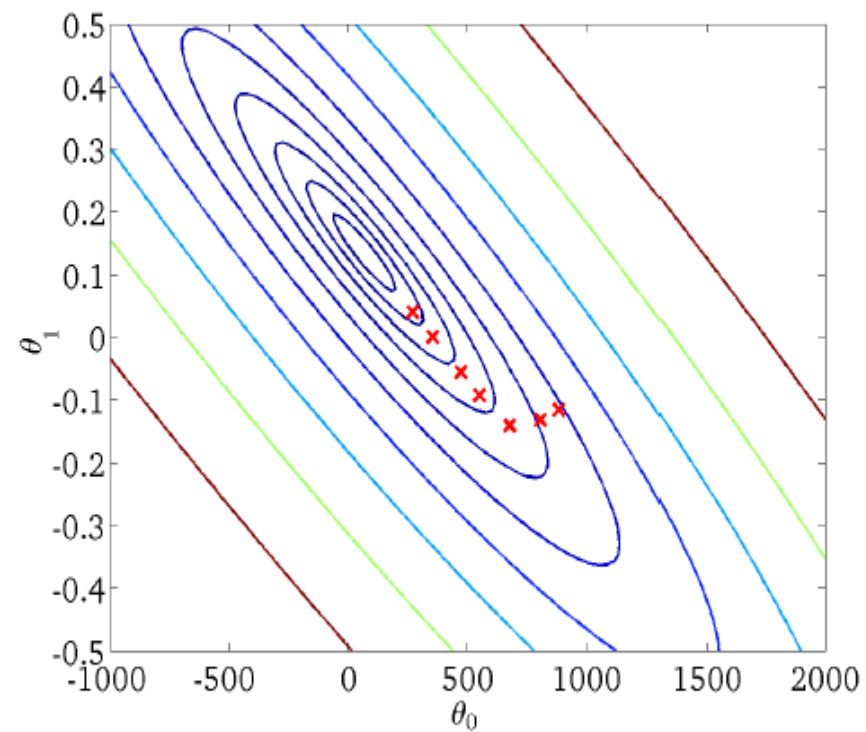
$$h_{\theta}(x)$$

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



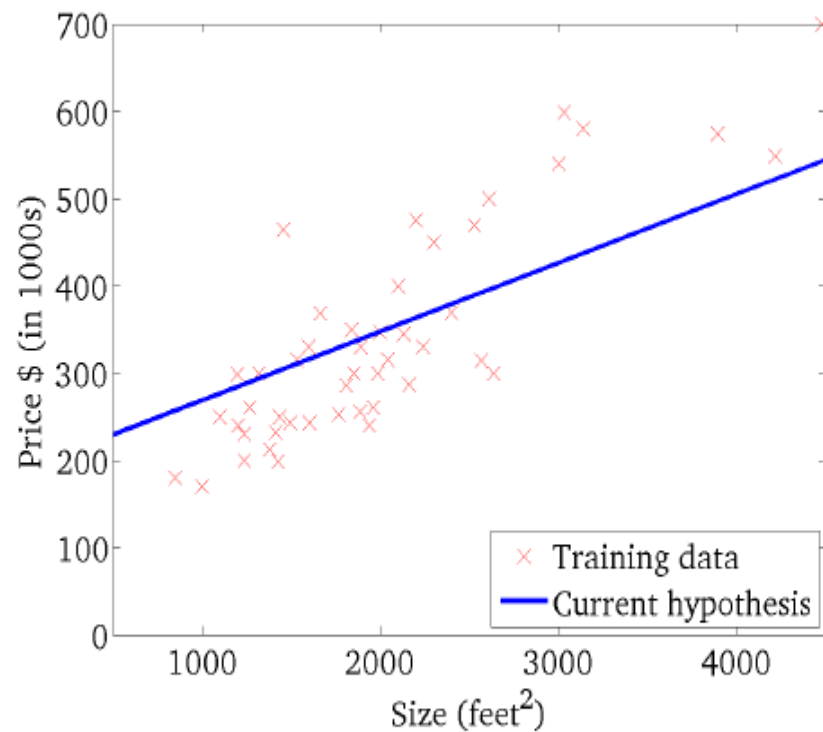
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



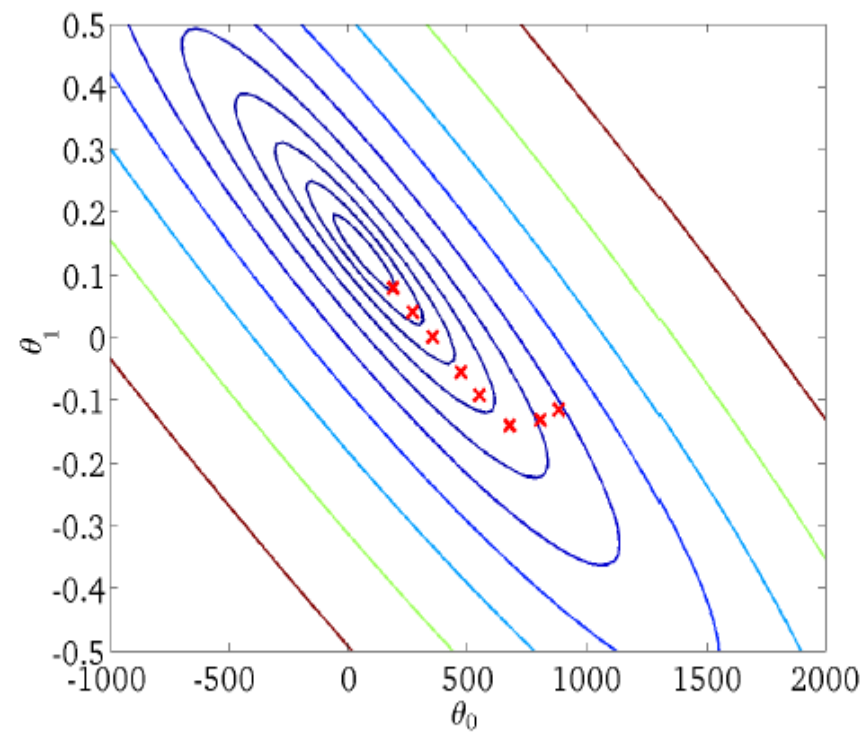
$$h_{\theta}(x)$$

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



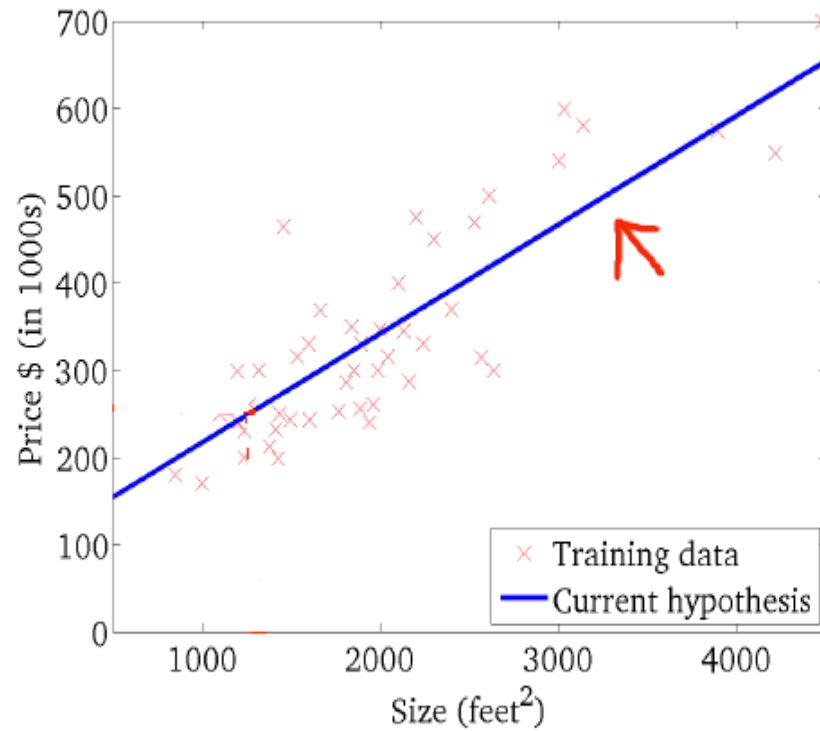
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



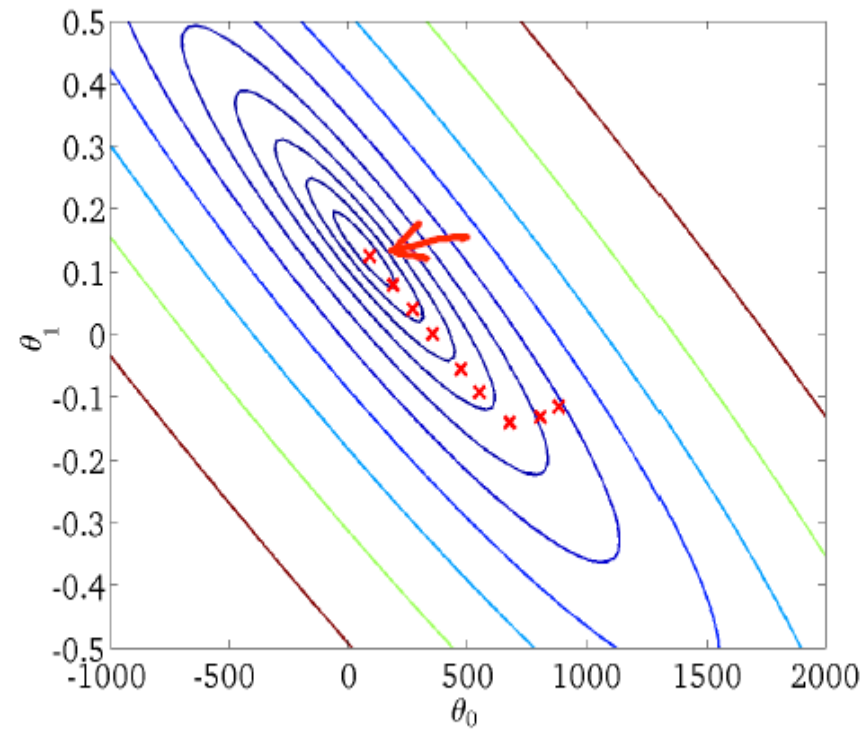
$$h_{\theta}(x)$$

(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)



Batch and Stochastic Gradient Descent

- GD is an iterative algorithm, that starts from a random point on the function and travels down its slope in steps until it reaches the lowest point of that function
- Batch GD
 - “Batch”: Each step of gradient descent uses all the training examples.
- Stochastic GD
 - Randomly pick one data point from whole dataset at each iteration → reduce computations
- Mini-batch GD
 - Sample small number of data points at each step

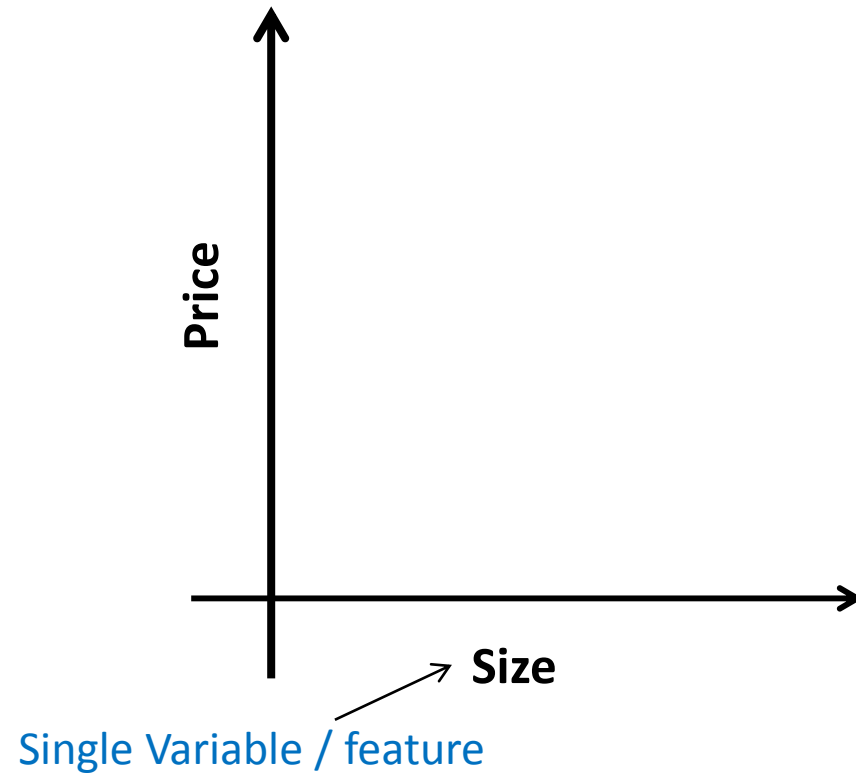
Multiple (Multivariable) Linear Regression

- Simple Regression
 - Only one feature or independent variable
- Multiple Regression
 - Multiple (More than 1) features or independent variables

Example – Simple Regression

Size (feet ²)	Price (\$1000)
x	y
2104	460
1416	232
1534	315
852	178
...	...

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



Example – Multiple Regression

- m = Number of Examples
- n = Number of features
- X = Feature Vector for training example
- $x_j^{(i)}$ = value of feature j in i^{th} training example

	Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
	2104	5	1	45	460
$\mathbf{X}^{(2)} \rightarrow$	1416 $x_1^{(2)}$	3 $x_2^{(2)}$	2 $x_3^{(2)}$	40 $x_4^{(2)}$	232
	1534	3	2	30	315
	852	2	1	36	178

$$\mathbf{X}^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$$

- Simple Regression

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

- Multiple Regression

$$h_{\theta}(\mathbf{X}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots$$

$x_1, x_2, x_3, x_4, \dots$ are features of a single example under consideration

- In general,

$$h_{\theta}(\mathbf{X}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

- For convenience of notations, let $x_0=1$

$$h_{\theta}(\mathbf{X}) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

$$\mathbf{X} = \begin{bmatrix} x_0 \\ x_1 \\ \cdot \\ \cdot \\ x_n \end{bmatrix} \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \cdot \\ \cdot \\ \theta_n \end{bmatrix}$$

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

$n+1$ column vectors

- Hypothesis

$$h_{\theta}(\mathbf{X}) = \boldsymbol{\theta}^T \mathbf{X} = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

$n+1$ dimensional vector

- Parameters

$$\theta_0, \theta_1, \theta_2, \dots, \theta_n$$

$$\boldsymbol{\theta}$$

$n+1$ dimensional vector

- Cost Function

$$J(\theta_0, \theta_1, \dots, \theta_n) = J(\boldsymbol{\theta}) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(\mathbf{X}^{(i)}) - y^{(i)})^2$$

- Goal

$$\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

- Gradient Descent

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1, \dots, \theta_n)$$

}

Simultaneous update
for every $j = 0, 1, 2, \dots, n$

Gradient Descent Algorithm

Simple Linear Regression

repeat until convergence {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

}

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

Update θ_0 and θ_1 simultaneously

Multiple Linear Regression

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(\mathbf{X}^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

}

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

Simultaneously update θ_j for every $j = 0, 1, 2, \dots, n$

Unit 2

Practical Tips for Linear Regression

Machine Learning – Regression and Decision Trees

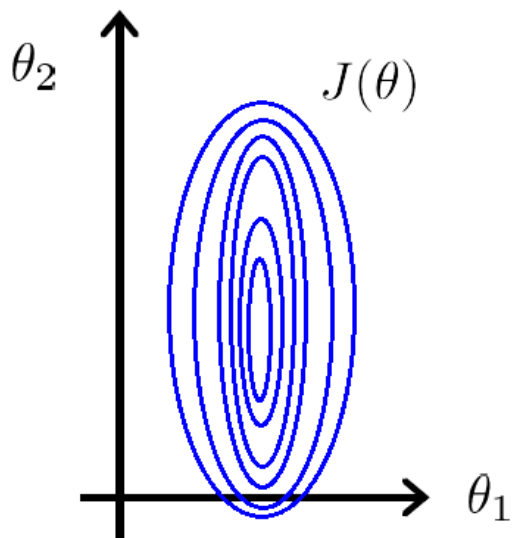
Some Practical Tips

- **Feature Scaling**

- Features should be in similar scale (i.e. they should have similar range of values) → For quick convergence

E.g. $x_1 = \text{size (0 – 2000 feet}^2\text{)}$

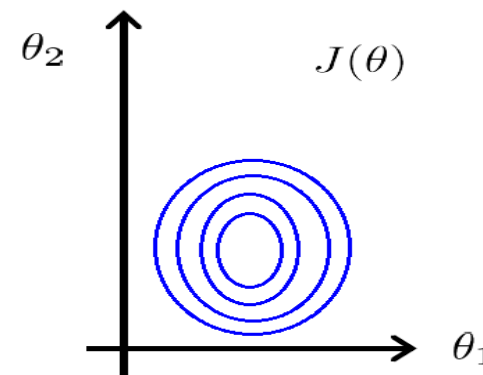
$x_2 = \text{number of bedrooms (1-5)}$



Method 1: Divide the feature value by maximum in the range

$$x_1 = \frac{\text{size}(\text{feet}^2)}{2000} \quad 0 \leq x_1 \leq 1$$

$$x_2 = \frac{\text{number of bedrooms}}{5} \quad 0 \leq x_2 \leq 1$$



Some Practical Tips

- Generally, get every feature approximately in the range $-1 \leq x_i \leq 1$
- Need not be exactly in the same range but at least closer to one another

Some Practical Tips

- Feature Scaling

- **Method 2: Mean Normalization**

Replace x_i with $(x_i - \mu_i)$ → make the features to have approximately zero mean
(Do not apply to $x_0 = 1$)

E.g.

Average value of x_i in training set

$$x_1 = \frac{\text{size} - 1000}{2000} \quad \text{Assuming average size of house} = 1000 \text{ feet}^2$$

$$x_2 = \frac{\text{number of bedrooms} - 2}{5} \quad -0.5 \leq x_1, x_2 \leq 0.5 \text{ (approx.)}$$

Generalization:

Average value of x_i in training set

$$x_i = \frac{x_i - \mu_i}{s_i}$$

Range of values i.e. $\max(x_i) - \min(x_i)$ in training set OR standard deviation

Some Practical Tips

- Learning Rate (α)

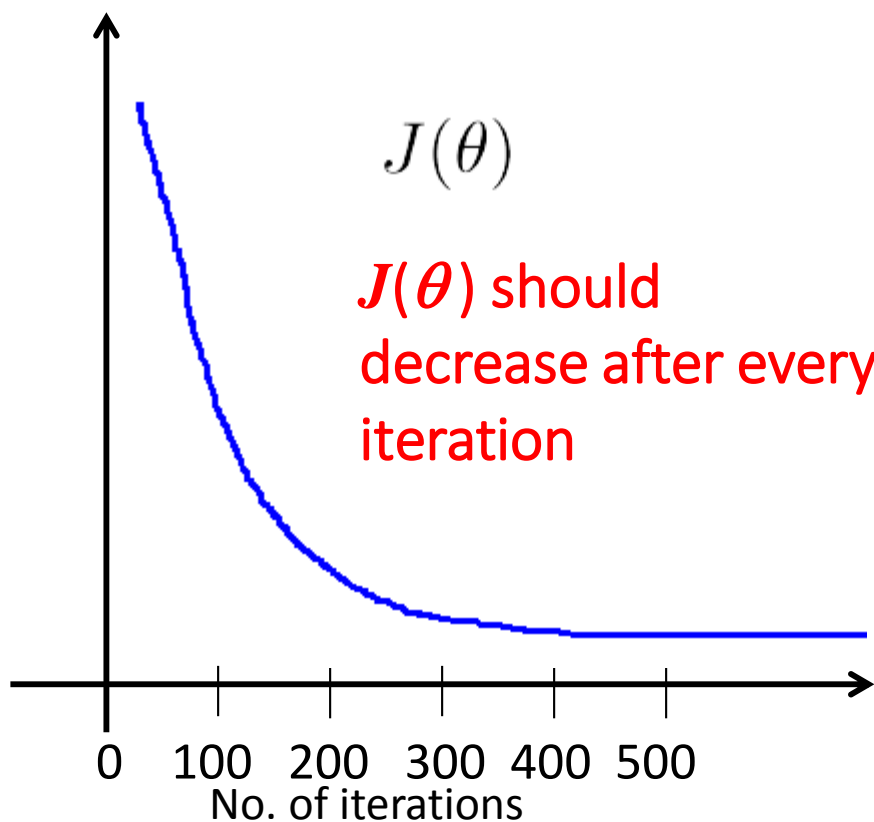
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$$

- How to make sure that gradient descent is working correctly ?
- How to choose learning rate ?

Some Practical Tips

- How to make sure that gradient descent is working correctly ?

$$\min_{\theta} J(\theta)$$

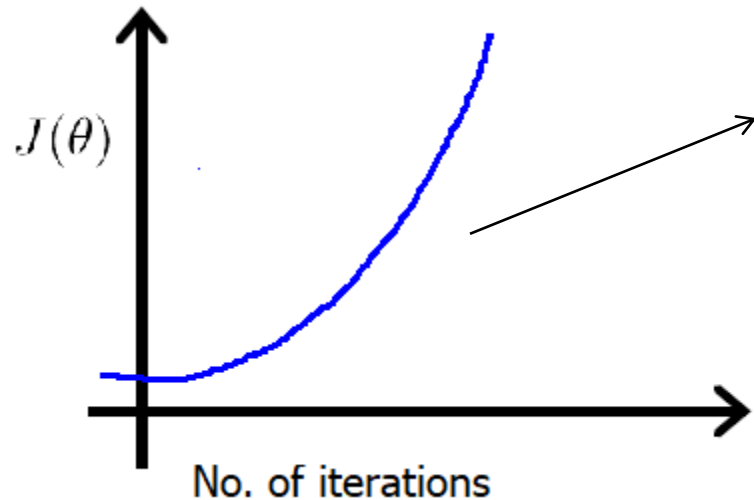


Automatic Convergence Test

Declare convergence if $J(\theta)$ decreases by less than 10^{-3} in one iteration. (ϵ)

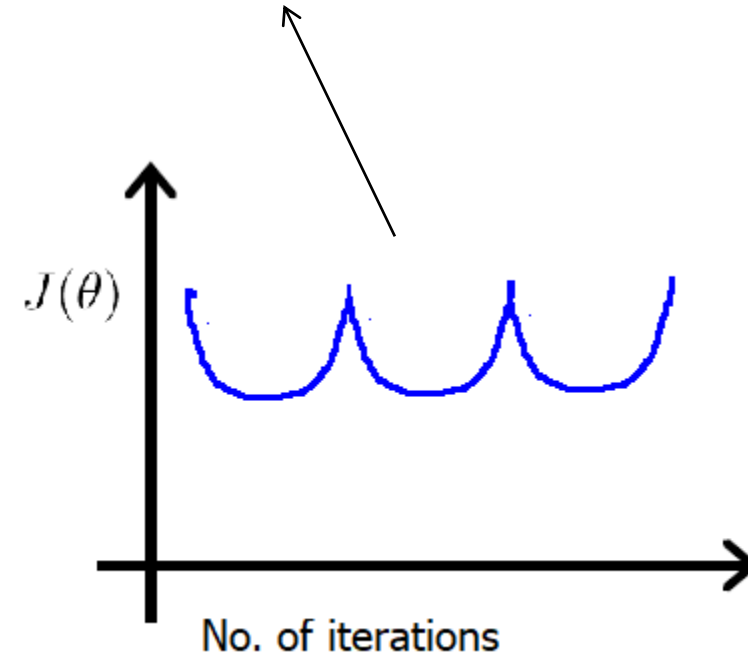
Some Practical Tips

- How to make sure that gradient descent is not working correctly ?



Gradient Descent not working →
Use smaller α

- For sufficiently small α , $J(\theta)$ should decrease on every iteration
- But if α is too small, gradient descent can be slow to converge



Some Practical Tips

- **Learning Rate (α) : Summary**

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

- If α is too small \rightarrow Slow convergence
- If α is too large : $J(\theta)$ may not decrease on every iteration; may not converge
- To choose α try the following
0.001, 0.01, 0.1, ,1

Unit 2

Polynomial Regression, Normal Equation

Machine Learning – Regression and Decision Trees

Creating our own features

- Example – Housing price prediction

$$h_{\theta}(\mathbf{X}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

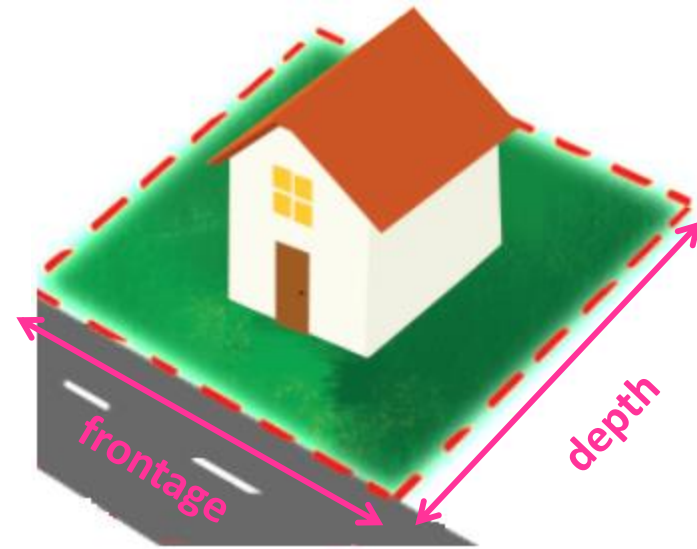
frontage

depth

Create New feature of our own

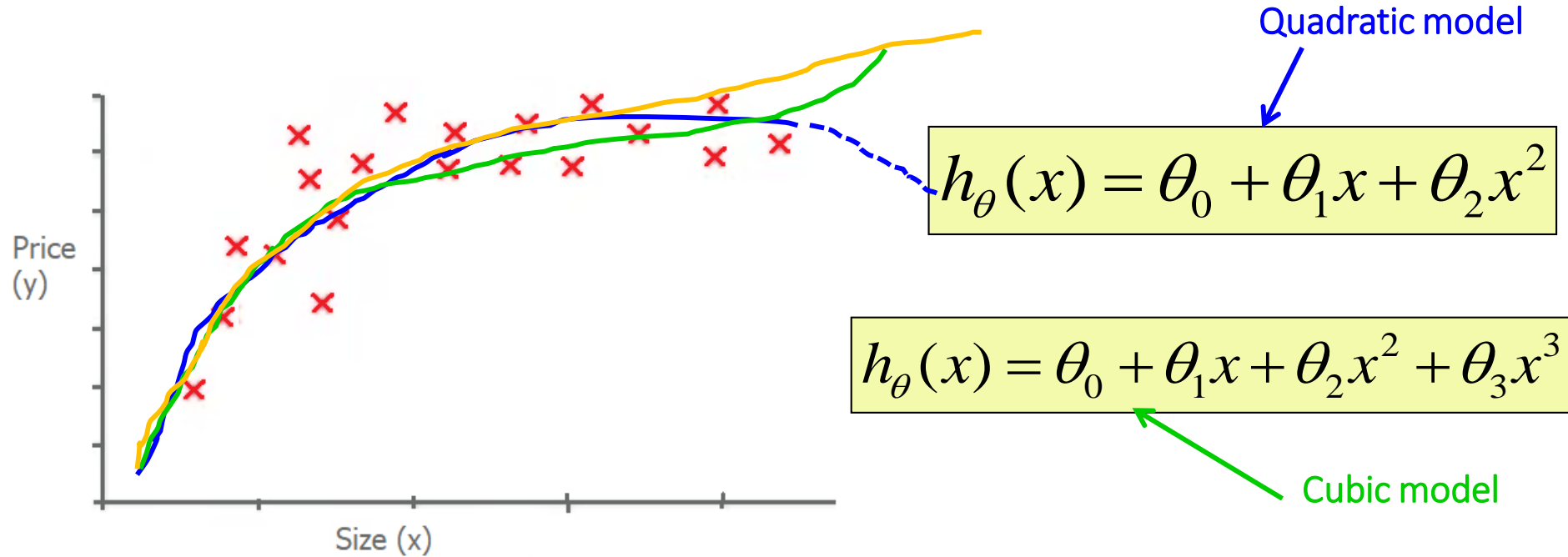
Area = frontage * depth

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



Depending on the insight of the problem one may create new features from existing ones

Polynomial Regression



$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$
$$= \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2 + \theta_3(\text{size})^3$$

Feature Scaling becomes important

Another choice

$$h_{\theta}(x) = \theta_0 + \theta_1(\text{size}) + \theta_2(\sqrt{\text{size}})$$

Normal Equation

- Method to solve for θ analytically ([Alternative for Gradient Descent](#))
- Example: Number of features, $n = 4$

Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
x_1	x_2	x_3	x_4	y
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 4 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$$

$m \times (n+1)$ matrix

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

m dimensional vector

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

Normal Equation

- Generalization: m examples are given as $(X^{(1)}, y^{(1)}), (X^{(2)}, y^{(2)}), \dots, (X^{(m)}, y^{(m)})$ ----
- n features per vector

$$\mathbf{X}^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ \vdots \\ x_n^{(i)} \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} \text{-----} (\mathbf{X}^{(1)})^T \text{-----} \\ \text{-----} (\mathbf{X}^{(2)})^T \text{-----} \\ \vdots \\ \vdots \\ \text{-----} (\mathbf{X}^{(m)})^T \text{-----} \end{bmatrix} \quad \begin{matrix} m \times (n+1) \\ \text{matrix} \end{matrix}$$

Design Matrix

Example

$$\mathbf{X}^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & x_1^{(1)} \\ 1 & x_1^{(2)} \\ \vdots & \vdots \\ \vdots & \vdots \\ 1 & x_1^{(m)} \end{bmatrix} \quad \begin{matrix} m \times 2 \\ \text{matrix} \end{matrix} \quad \mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ \vdots \\ y^{(m)} \end{bmatrix}$$

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

m training examples, n features per example

Gradient Descent	Normal Equation
Need to choose α	No need to choose α
Needs many iterations	Do not need to iterate
Feature scaling required	No feature scaling
Works well even when n is large (Typically, if $n > 10^5$ use gradient descent)	Need to compute $(X^T X)^{-1}$ If X is of size $(n \times n)$, Computation time approx. $O(n^3)$ → Slow if n is large

Unit 2

Logistic Regression

Machine Learning – Regression and Decision Trees

Classification

- Problem of Classification
 - Email: Spam / Not Spam?
 - Online Transactions: Fraudulent (Yes / No)?
 - Tumor: Malignant / Benign ?

$$y \in \{0, 1\}$$

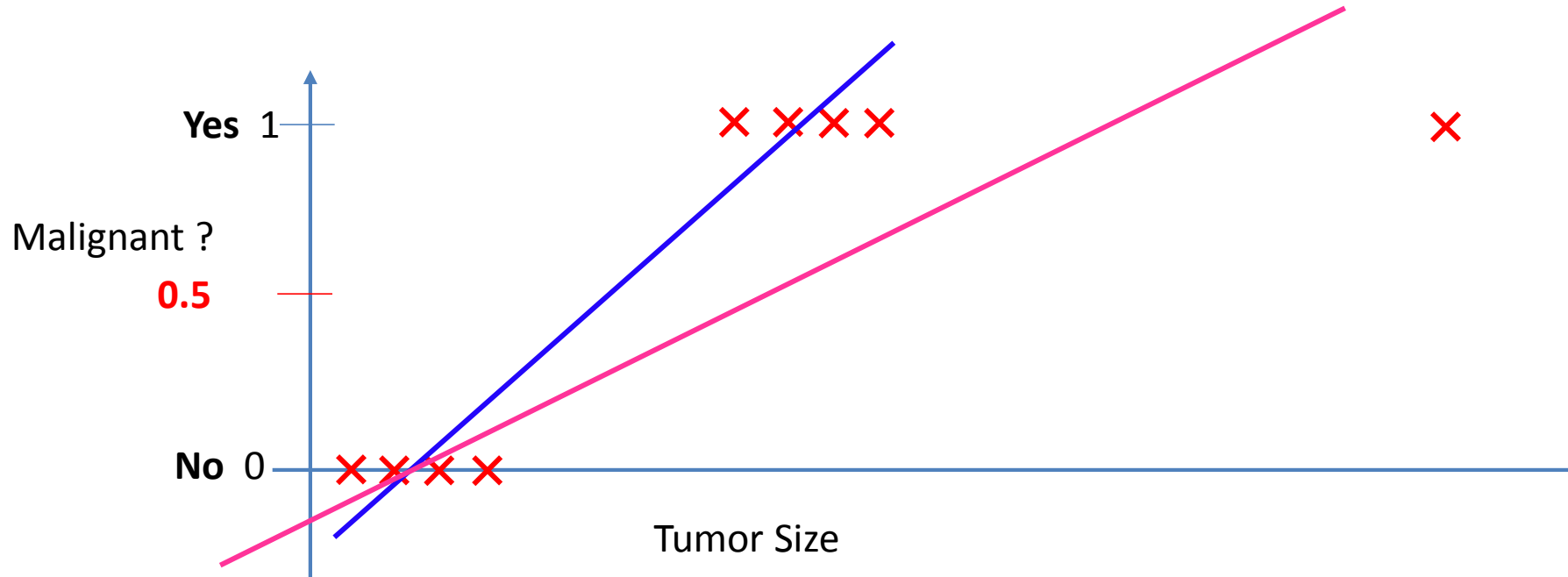
0: Negative Class (Benign)
1: Positive Class (Malignant)

Binary classification

$$y \in \{0, 1, 2, 3\}$$

Multiclass classification

How to develop a classification algorithm?



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

Threshold classifier output $h_{\theta}(x)$ at 0.5 :

If $h_{\theta}(x) \geq 0.5$, predict “ $y = 1$ ”

If $h_{\theta}(x) < 0.5$, predict “ $y = 0$ ”

Applying Linear Regression to Classification Problem is not a good idea !

Logistic Regression

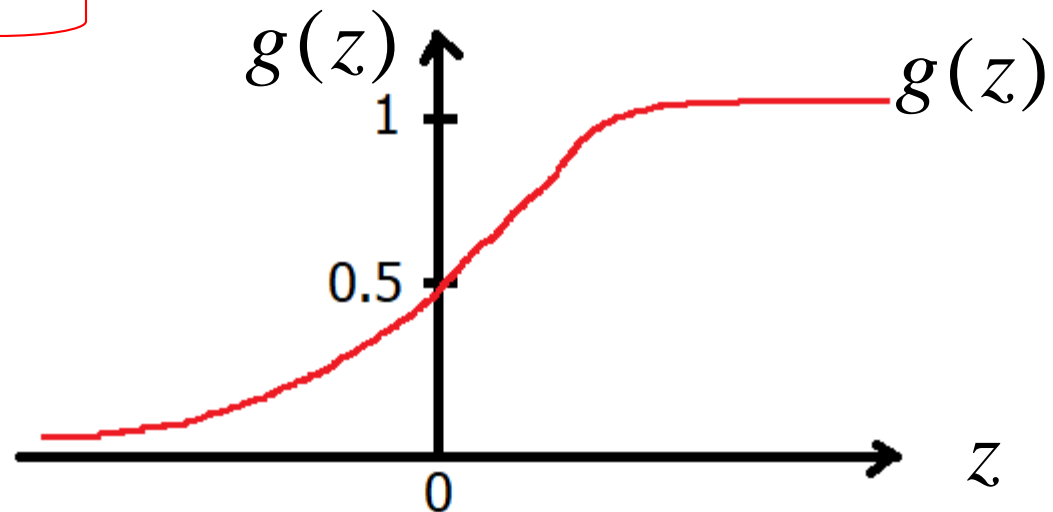
- Find hypothesis such as $0 \leq h_{\theta}(x) \leq 1$

$$h_{\theta}(\mathbf{X}) = g(\theta^T \mathbf{X})$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

Sigmoid function /
Logistic function

$$h_{\theta}(\mathbf{X}) = \frac{1}{1 + e^{-\theta^T \mathbf{X}}}$$



Logistic Regression

- How to interpret hypothesis ?

- $h_{\theta}(x)$ = estimated probability that $y = 1$ on input x

- Example $\mathbf{X} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ tumorSize \end{bmatrix}$

- $h_{\theta}(x) = 0.8$ \rightarrow probability that $y = 1$ is 0.8
 \rightarrow 80 % chance to have tumor as malignant

$$h_{\theta}(\mathbf{X}) = P(y = 1 | \mathbf{X}; \theta)$$

“probability that $y = 1$, given X ,
parameterized by θ ”

$$P(y = 0 | \mathbf{X}; \theta) + P(y = 1 | \mathbf{X}; \theta) = 1$$

$$P(y = 0 | \mathbf{X}; \theta) = 1 - P(y = 1 | \mathbf{X}; \theta)$$

Unit 2

Logistic Regression – Decision Boundary, Cost Function

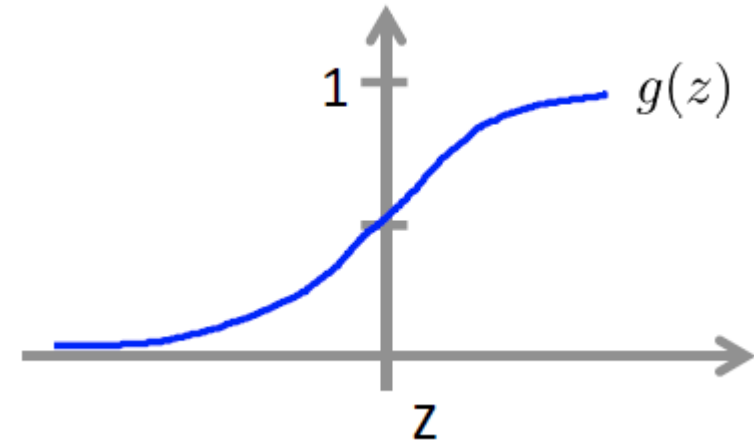
Machine Learning – Regression and Decision Trees

Decision Boundary

- Logistic Regression

$$h_{\theta}(\mathbf{X}) = g(\theta^T \mathbf{X}) = P(y = 1 | \mathbf{X}; \theta)$$

$$g(z) = \frac{1}{1 + e^{-z}}$$



Predict $y = \text{"1"}$ if $h_{\theta}(x) \geq 0.5$

whenever $\theta^T X \geq 0$

If $z \geq 0$, $g(z) \geq 0.5$

$\Rightarrow h_{\theta}(x) = g(\theta^T X) \geq 0.5$

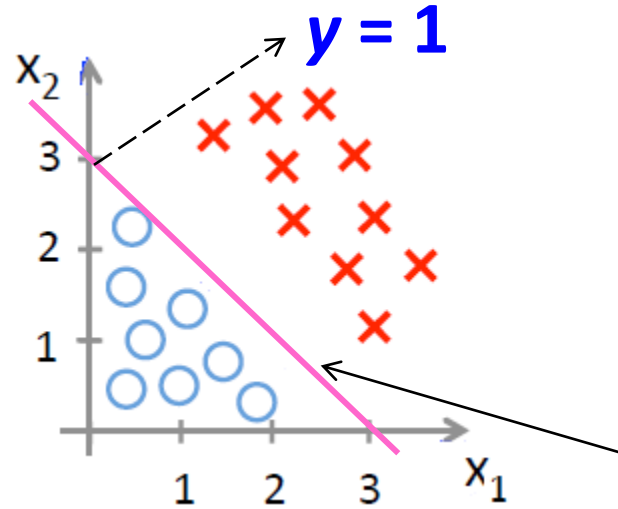
whenever $\theta^T X \geq 0$

Predict $y = \text{"0"}$ if $h_{\theta}(x) < 0.5$

whenever $\theta^T X < 0$

If $z < 0$, $g(z) < 0.5$

Decision Boundary



Decision Boundary
 $x_1 + x_2 = 3$

$$h_{\theta}(\mathbf{X}) = g(\underbrace{\theta_0}_{=-3} + \underbrace{\theta_1}_{=1}x_1 + \underbrace{\theta_2}_{=1}x_2)$$

$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

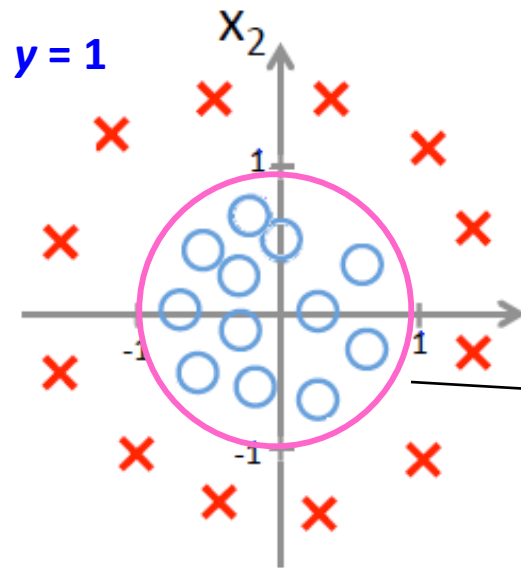
$$h_{\theta}(\mathbf{X}) = g(-3 + x_1 + x_2)$$

Predict $y = "1"$ if $h_{\theta}(x) \geq 0.5$
whenever $\theta^T X \geq 0$

Predict $y = "1"$ if $h_{\theta}(x) \geq 0.5 \Rightarrow$ if $-3 + x_1 + x_2 \geq 0$
 $\Rightarrow x_1 + x_2 \geq 3$

Predict $y = "0"$ if $h_{\theta}(x) < 0.5 \Rightarrow$ if $-3 + x_1 + x_2 < 0$
 $\Rightarrow x_1 + x_2 < 3$

Non-linear decision boundaries

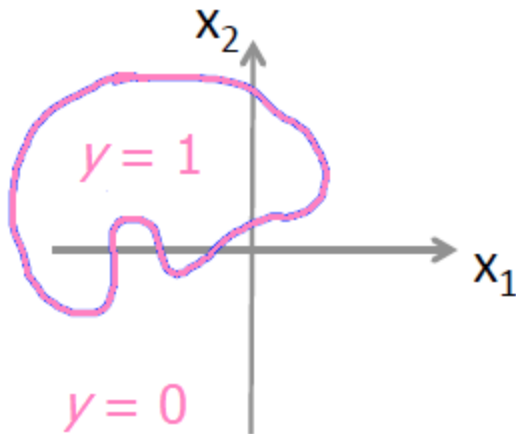


$$h_{\theta}(\mathbf{X}) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

Decision Boundary
 $x_1^2 + x_2^2 = 1$

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

Predict $y = "1"$ if $h_{\theta}(x) \geq 0.5 \quad \Rightarrow \text{if } -1 + x_1^2 + x_2^2 \geq 0$
 $\Rightarrow x_1^2 + x_2^2 \geq 1$



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^2 x_2^2 + \dots)$$

Logistic Regression – Cost Function

Machine Learning – Regression and Decision Trees

Logistic Regression – Cost Function

- Training Set $\{(X^{(1)}, y^{(1)}), (X^{(2)}, y^{(2)}), \dots, (X^{(m)}, y^{(m)})\}$ ---- m examples, n features each

$$\mathbf{X} = \begin{bmatrix} x_0 \\ x_1 \\ \cdot \\ \cdot \\ x_n \end{bmatrix}$$

$$x_0 = 1, y \in \{0, 1\}$$

How to choose
parameters θ ?

- Hypothesis

$$h_{\theta}(\mathbf{X}) = \frac{1}{1 + e^{-\theta^T \mathbf{X}}}$$

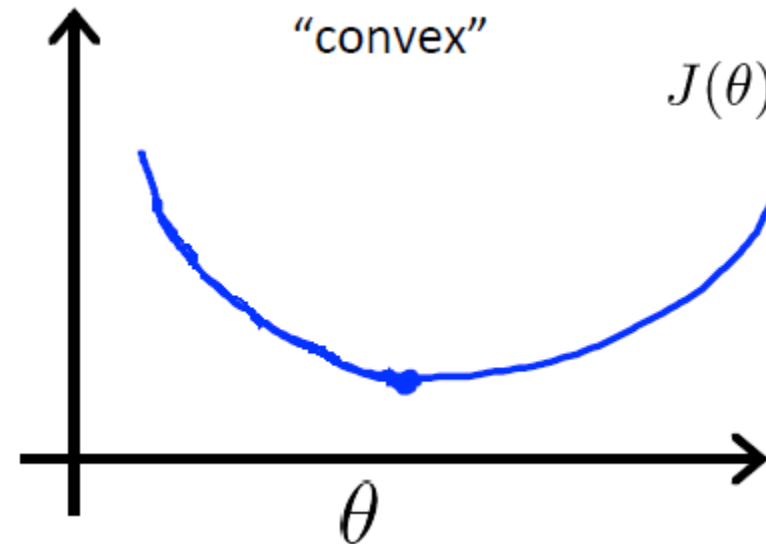
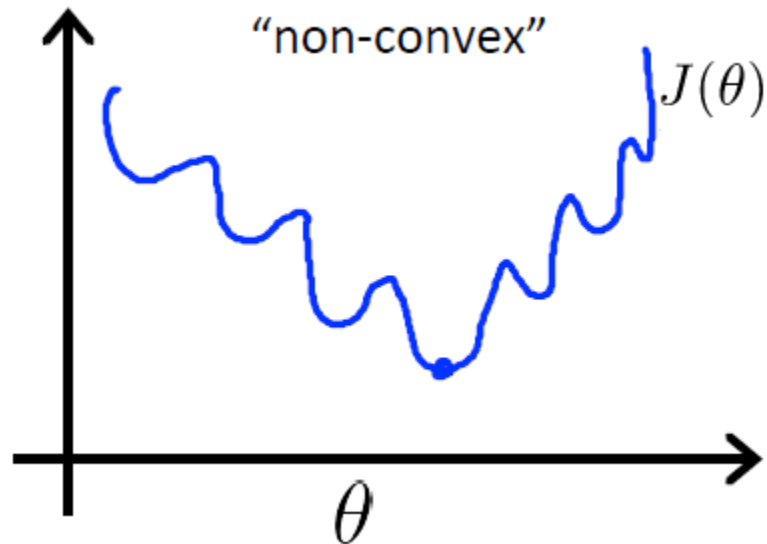
Cost Function

- Linear Regression

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

$$\text{Cost}(h_{\boldsymbol{\theta}}(\mathbf{X}^{(i)}), y^{(i)})$$

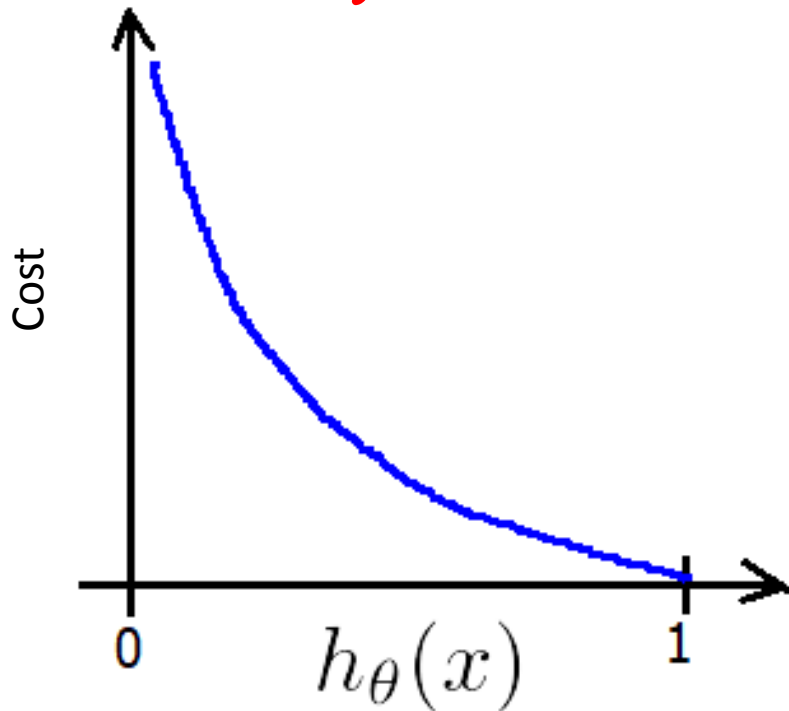
$$h_{\boldsymbol{\theta}}(\mathbf{X}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{X}}}$$



Logistic Regression – Cost Function

$$Cost(h_{\theta}(\mathbf{X}), y) = \begin{cases} -\log(h_{\theta}(\mathbf{X})) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(\mathbf{X})) & \text{if } y = 0 \end{cases}$$

If $y = 1$



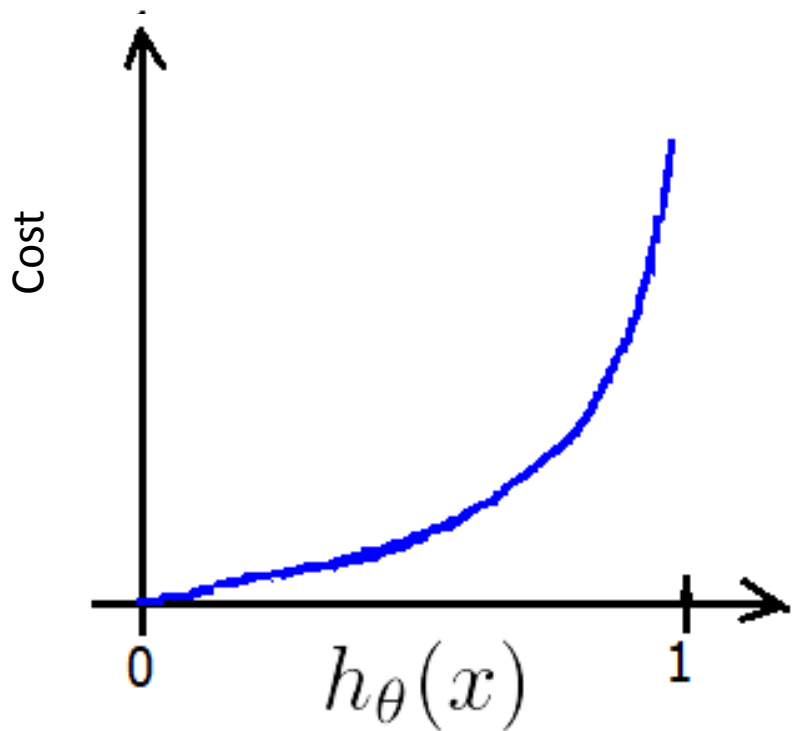
Cost = 0 if $y = 1$ and $h_{\theta}(x) = 1$

But as $h_{\theta}(x) \rightarrow 0$, Cost $\rightarrow \infty$

Logistic Regression – Cost Function

$$Cost(h_{\theta}(\mathbf{X}) - y) = \begin{cases} -\log(h_{\theta}(\mathbf{X})) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(\mathbf{X})) & \text{if } y = 0 \end{cases}$$

If $y = 0$



Cost = 0 if $y = 0$ and $h_{\theta}(x) = 0$

But as $h_{\theta}(x) \rightarrow 1$, Cost $\rightarrow \infty$

Logistic Regression – Simplified Cost Function

$$J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\boldsymbol{\theta}}(\mathbf{X}^{(i)}), y^{(i)})$$

$$\text{Cost}(h_{\boldsymbol{\theta}}(\mathbf{X}), y) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\mathbf{X})) & \text{if } y = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\mathbf{X})) & \text{if } y = 0 \end{cases}$$

Simplified Cost Function

$$\text{Cost}(h_{\boldsymbol{\theta}}(\mathbf{X}), y) = -y \log(h_{\boldsymbol{\theta}}(\mathbf{X})) - (1 - y) \log(1 - h_{\boldsymbol{\theta}}(\mathbf{X}))$$

Logistic Regression – Simplified Cost Function

$$Cost(h_{\theta}(\mathbf{X}), y) = -y \log(h_{\theta}(\mathbf{X})) - (1 - y) \log(1 - h_{\theta}(\mathbf{X}))$$



This is for single example / single sample

For m examples / m samples

$$\begin{aligned} J(\boldsymbol{\theta}) &= \frac{1}{m} \sum_{i=1}^m Cost(h_{\theta}(\mathbf{X}^{(i)}), y^{(i)}) \\ &= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h_{\theta}(\mathbf{X}^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(\mathbf{X}^{(i)})) \right] \end{aligned}$$

Find parameters θ such that,

$$\min_{\theta} J(\theta)$$

Predict $h_{\theta}(x)$ for new x

$$h_{\theta}(\mathbf{X}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{X}}}$$

$$\nearrow P(y = 1 | \mathbf{x}; \boldsymbol{\theta})$$

Logistic Regression and Gradient Descent

$$J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\boldsymbol{\theta}}(\mathbf{X}^{(i)}), y^{(i)})$$
$$= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h_{\boldsymbol{\theta}}(\mathbf{X}^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\boldsymbol{\theta}}(\mathbf{X}^{(i)})) \right]$$

Find parameters $\boldsymbol{\theta}$ such that,

$$\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

Gradient Descent

$$\frac{1}{m} \sum_{i=1}^m (h_{\boldsymbol{\theta}}(\mathbf{X}^{(i)}) - y^{(i)}) x_j^{(i)}$$

Repeat {

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

Simultaneous update all θ_j

Logistic Regression and Gradient Descent

Observation:

For Linear Regression,

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

For Logistic Regression,

$$h_{\theta}(\mathbf{X}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{X}}}$$

$$- y^{(i)} \log(1 - h_{\theta}(\mathbf{X}^{(i)})) \Big]$$

$$\min_{\theta} J(\theta)$$

Gradient Descent

Whether it is identical to linear regression ?

Repeat {

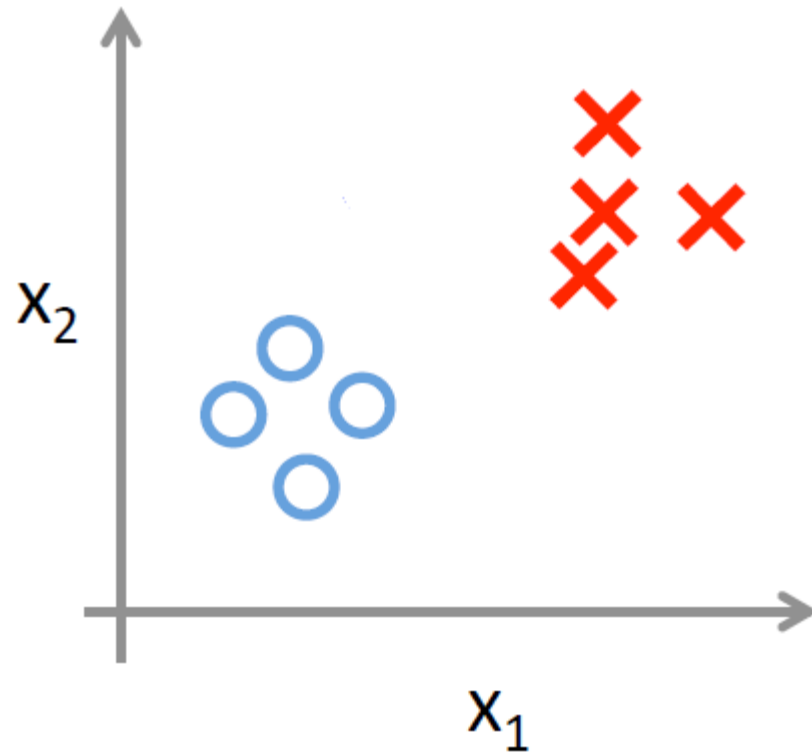
$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(\mathbf{X}^{(i)}) - y^{(i)}) x_j^{(i)}$$

Simultaneous update all θ_j

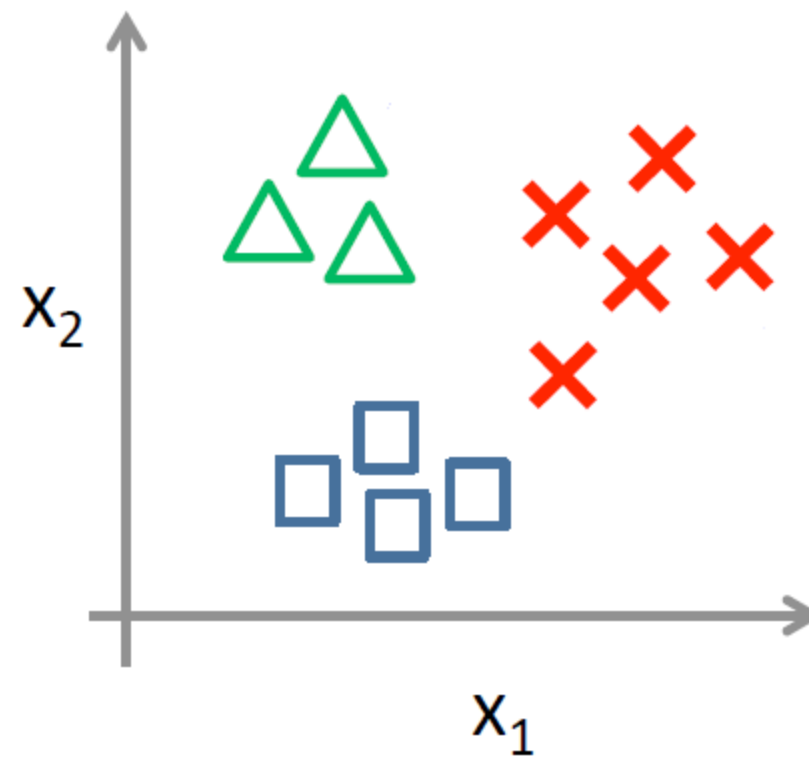
Logistic Regression for Multi-class Classification

- Categorize emails as:
Work, Friends, Family
- Categorize Weather as:
Sunny, Cloudy, Rainy, Snow

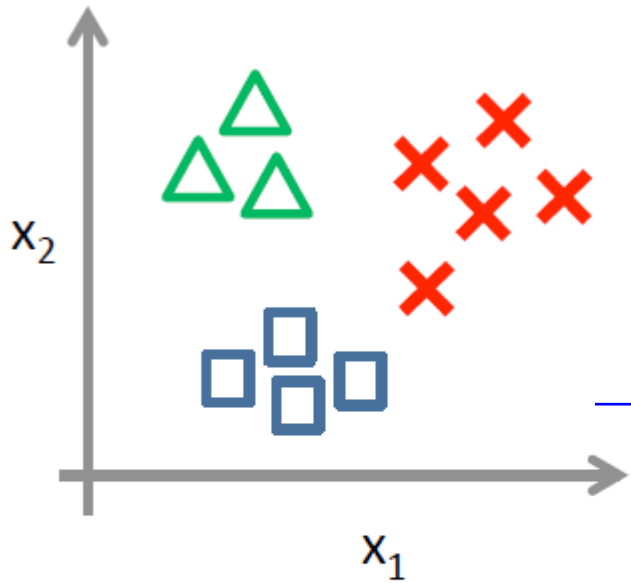
Binary Classification





Multi-class Classification




One-vs-all (one-vs-rest)

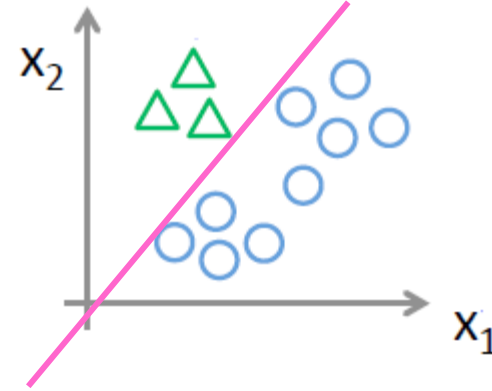


Class 1: 

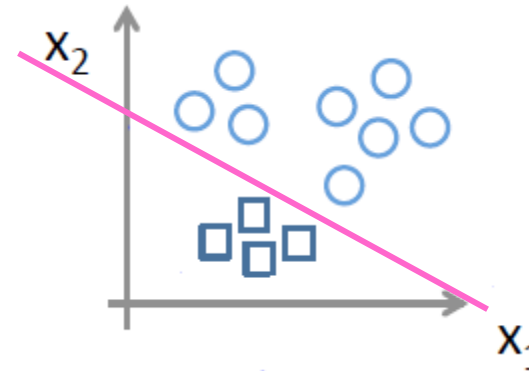
Class 2: 

Class 3: 

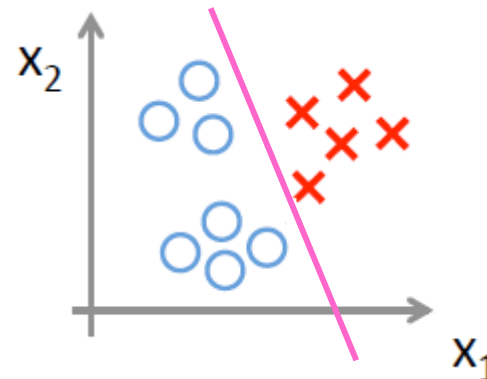
Convert multi-class classification
problem into three separate binary
class classification problems



$$h_{\theta}^{(1)}(\mathbf{X})$$



$$h_{\theta}^{(2)}(\mathbf{X})$$



$$h_{\theta}^{(3)}(\mathbf{X})$$

$$h_{\theta}^{(i)}(\mathbf{X}) = P(y = i | \mathbf{X}; \theta) \quad i = 1, 2, 3$$

One-vs-all (one-vs-rest)

- Train a logistic regression classifier $h_{\theta}^{(i)}(\mathbf{X})$ for each class i to predict the probability that $y = i$
- On a new input x , to make a prediction, pick the class i that maximizes

$$\max_i h_{\theta}^{(i)}(\mathbf{X})$$

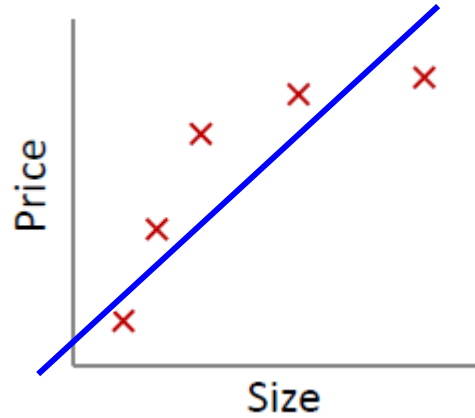
Unit 2

The Problem of Overfitting

Machine Learning – Regression and Decision Trees

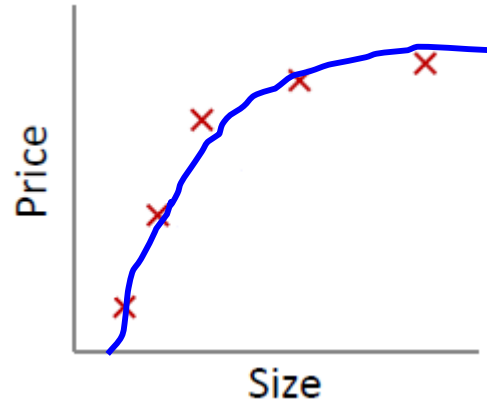
The Problem of Overfitting

Example: Linear Regression (Housing Price)



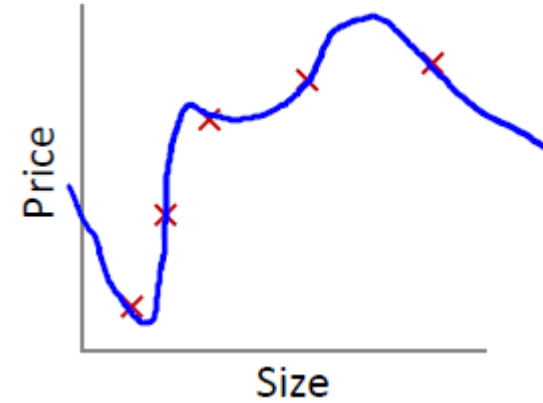
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Under fitting



$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

Right fitting



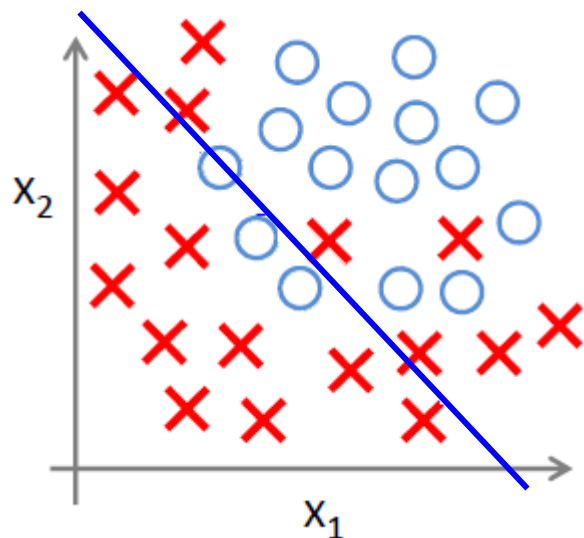
$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Overfitting

Overfitting: If we have too many features, the learned hypothesis may fit the training set very well. i.e. $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(\mathbf{X}^{(i)}) - y^{(i)})^2 \approx 0$, but fail to generalize to new examples

The problem of Overfitting

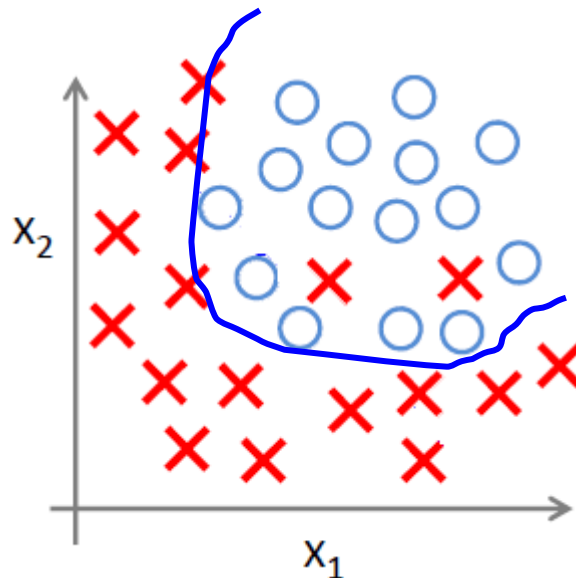
Example: Logistic Regression



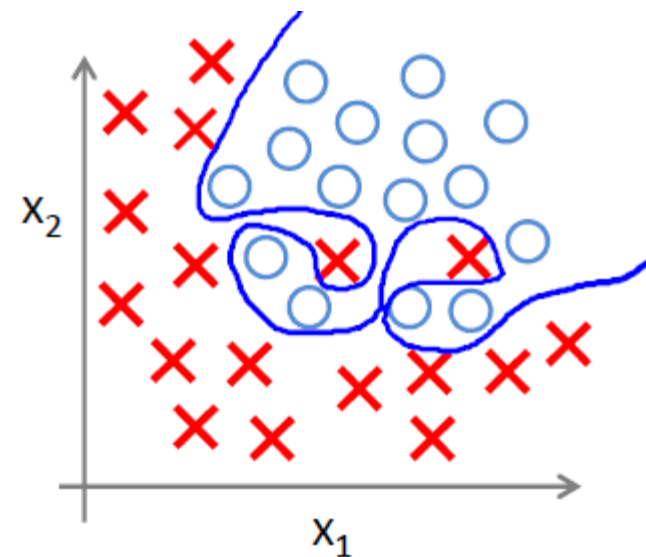
$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

Sigmoid function

Under fitting



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \dots)$$

Overfitting

How to address overfitting ?

x_1 = size of house

x_2 = no. of bedrooms

x_3 = no. of floors

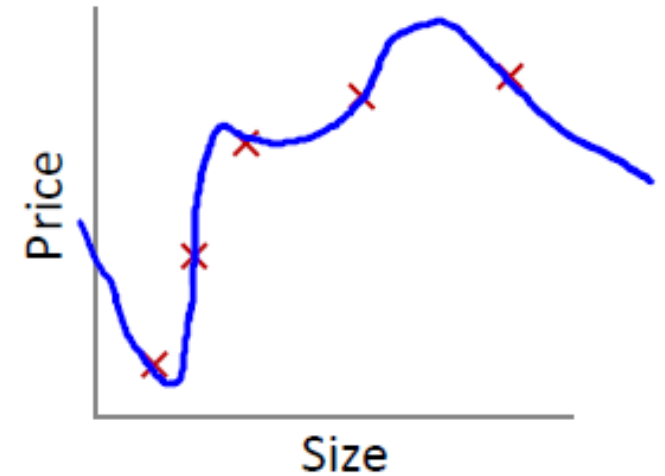
x_4 = age of house

x_5 = average income in neighborhood

x_6 = kitchen size

⋮

x_{100}



If we have a lot of features, and, very little training data, then, overfitting can become a problem

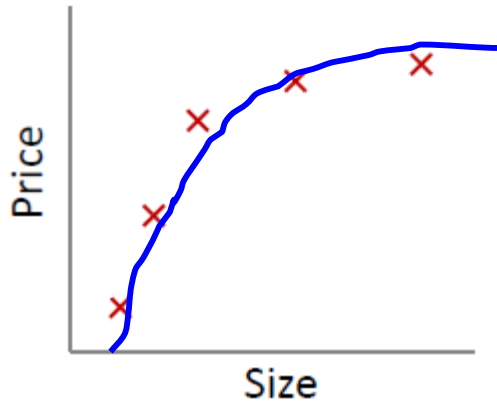
How to address overfitting ?

1. Reduce number of features

- Manual selection of features
- Feature selection algorithm

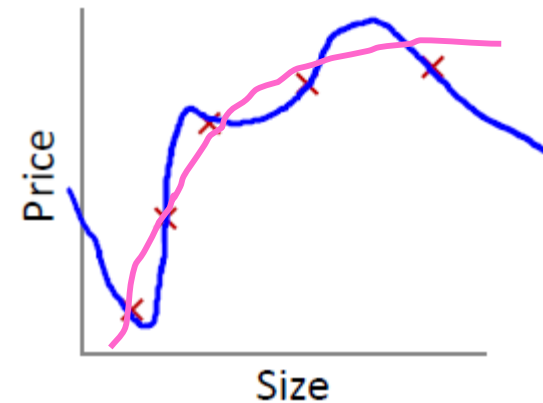
2. Regularization

- Regularisation is a technique used to reduce the errors by fitting the function appropriately on the given training set and avoid overfitting.
- Keep all the features but reduce magnitude / values of parameter θ_j
- Works well when we have a lot of features, each of which contributes a bit to predicting



$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(\mathbf{X}_i) - y_i)^2 \quad \min_{\theta} J(\theta)$$



$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Make θ_3, θ_4 very small

Make $\theta_3 \approx 0, \theta_4 \approx 0$

Regularization

- Small values for parameters $\theta_1, \theta_2, \dots, \theta_m$
 - Simpler hypothesis
 - Less prone to overfitting
- Example : $(x_0, x_1, x_2, x_3, \dots, x_{100})$ and

$$(\theta_1, \theta_2, \dots, \theta_{100})$$

$$\min_{\theta} J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(\mathbf{X}^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

Regularization parameter

Regularization term

Regularized Linear Regression

$$J(\boldsymbol{\theta}) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\boldsymbol{\theta}}(\mathbf{X}^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right] \quad \min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

- Gradient Descent

Repeat {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\boldsymbol{\theta}}(\mathbf{X}^{(i)}) - y^{(i)})$$

$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\boldsymbol{\theta}}(\mathbf{X}^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right]$$

}

$j = 1, 2, 3, \dots, n$

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

< 1