# Machine Learning Study Week 1

# Introduction

구 본 규 (protocolstack9@gmail.com)

Introduction

Welcome

Machine Learning

# Machine Learning

- Grew out of work in AI
- New capability for computers

# Examples:

- Database mining

  Large datasets from growth of automation/web.

  E.g., Web click data, medical records, biology, engineering
- Applications can't program by hand.

  E.g., Autonomous helicopter, handwriting recognition, most of Natural Language Processing (NLP), Computer Vision.
- Self-customizing programs

  E.g., Amazon, Netflix product recommendations
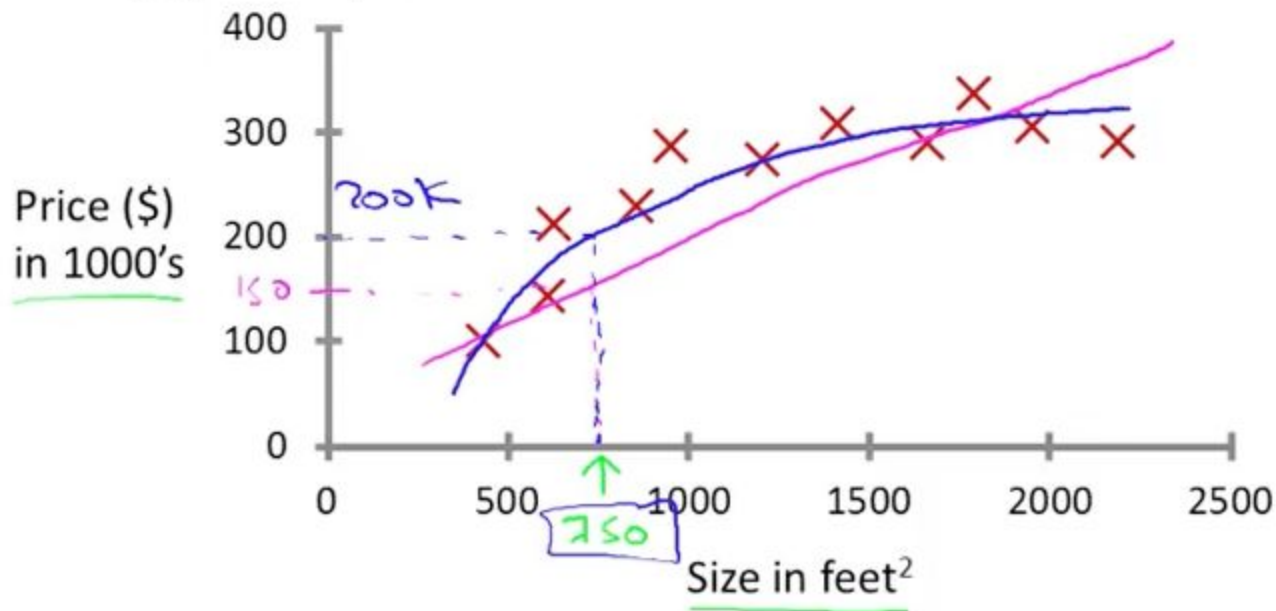- Understanding human learning (brain, real AI).

Machine Learning

Introduction

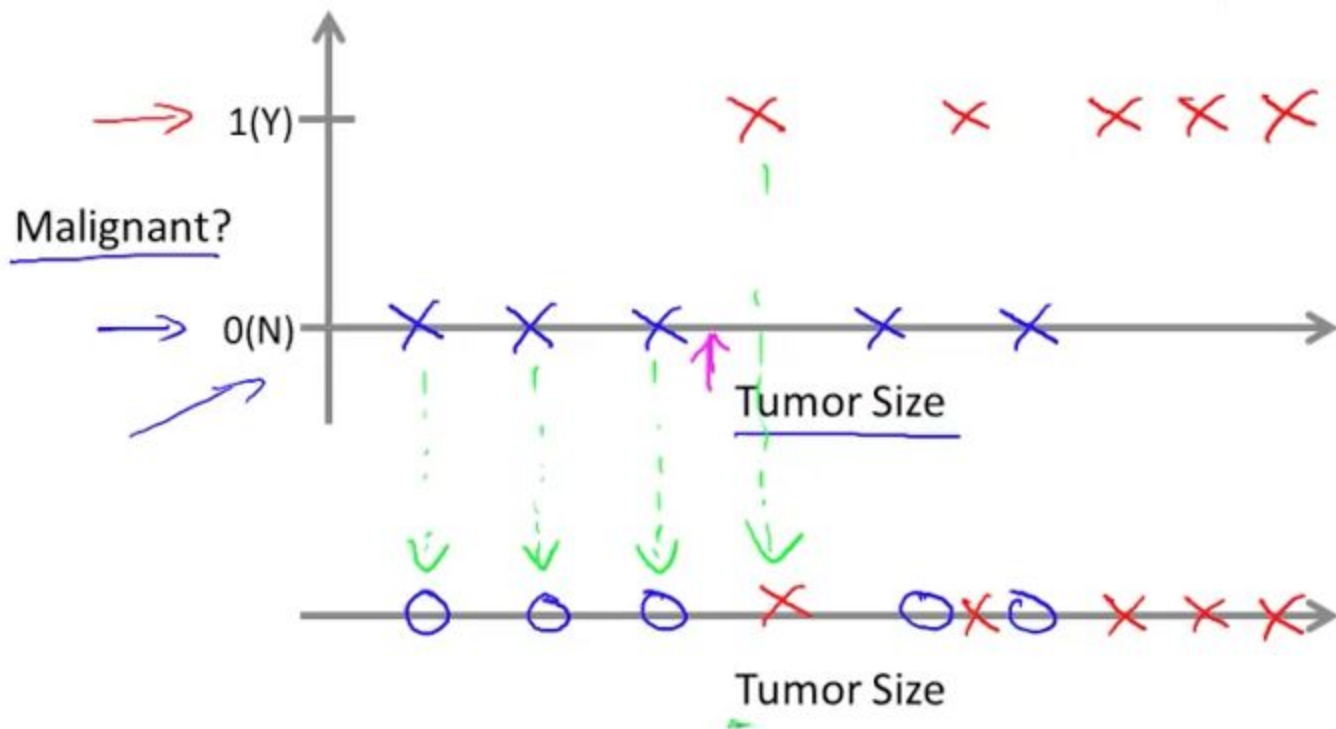Supervised
Learning

# Housing price prediction.



Price ($) in 1000's

Size in feet²

200K

150

750

**Supervised Learning**
"right answers" given

**Regression:** Predict continuous valued output (price)

# Breast cancer (malignant, benign)



Malignant?

1(Y)

0(N)

Tumor Size

Tumor Size

Classification

Discrete valued output (0 or 1)

$0, 1, 2, 3$

benign    type1 Cancer
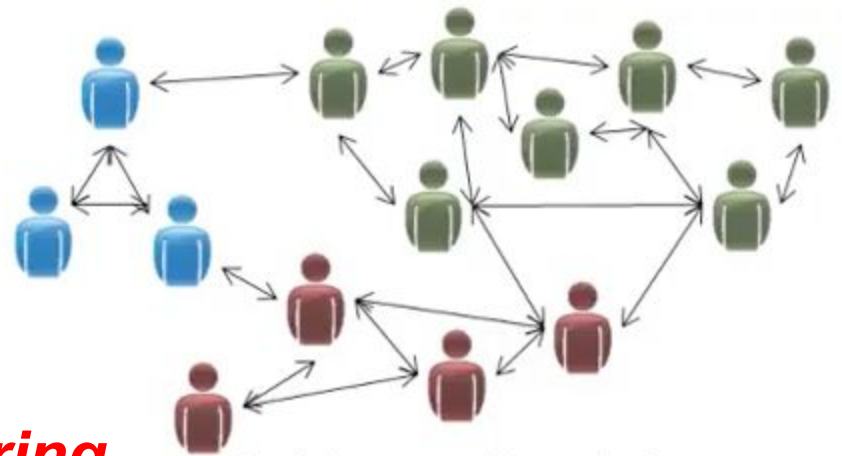
Machine Learning

Introduction

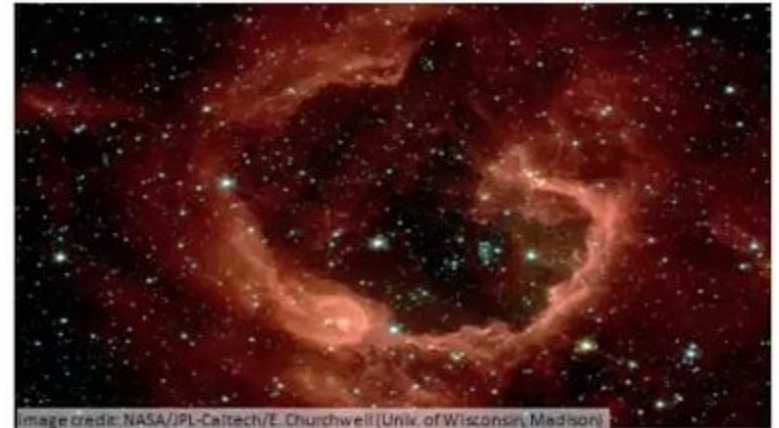Unsupervised Learning

*Clustering*

Organize computing clusters
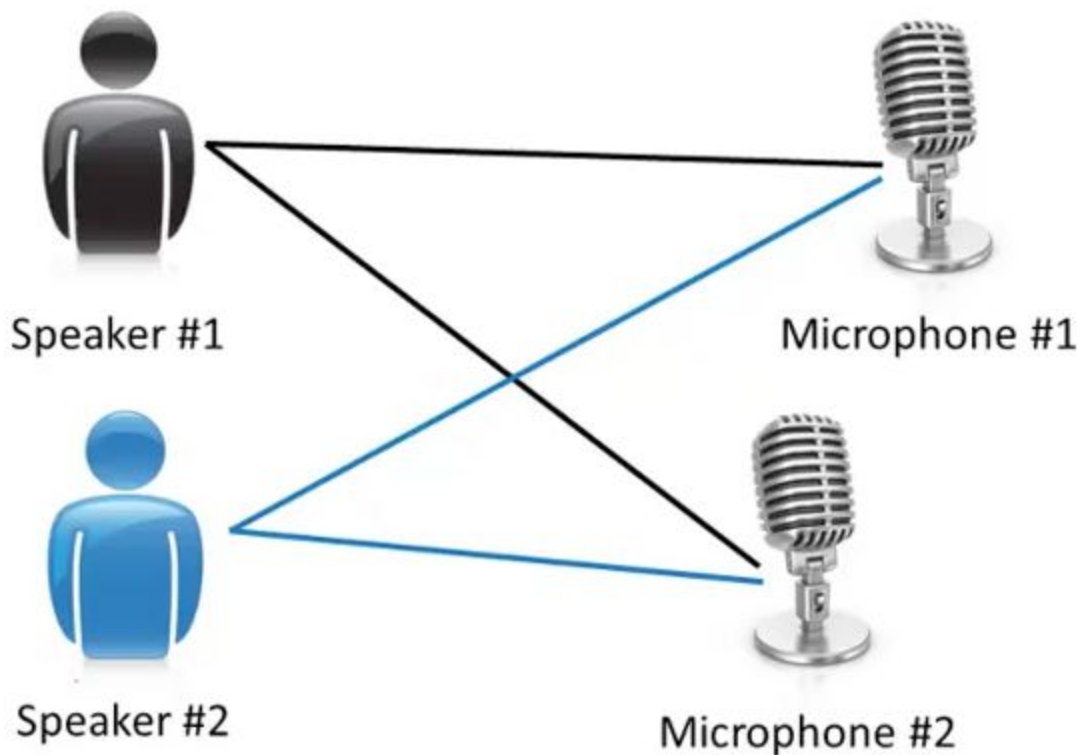
*Clustering*

Social network analysis

Market segmentation

Astronomical data analysis

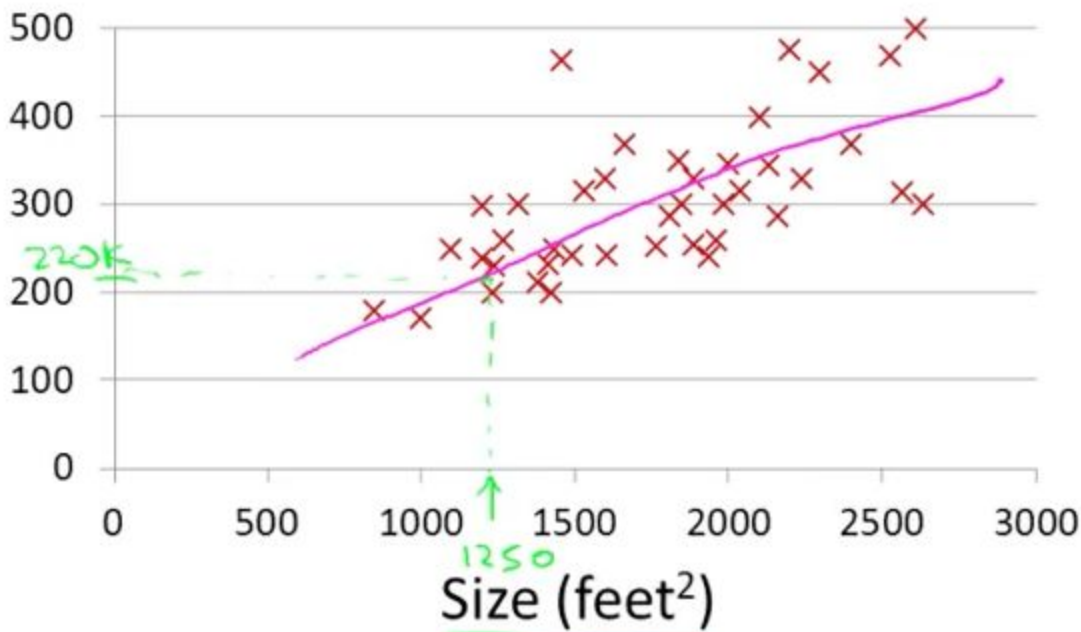image credit: NASA/JPL-Caltech/E. Churchwell (Univ. of Wisconsin, Madison)

Cocktail party problem

# Housing Prices (Portland, OR)



Price (in 1000s of dollars) vs Size (feet²)

220k, 1250

## Supervised Learning

Given the "right answer" for each example in the data.

## Regression Problem

Predict real-valued output

Classification: Discrete-valued output

**Training set of housing prices (Portland, OR)**

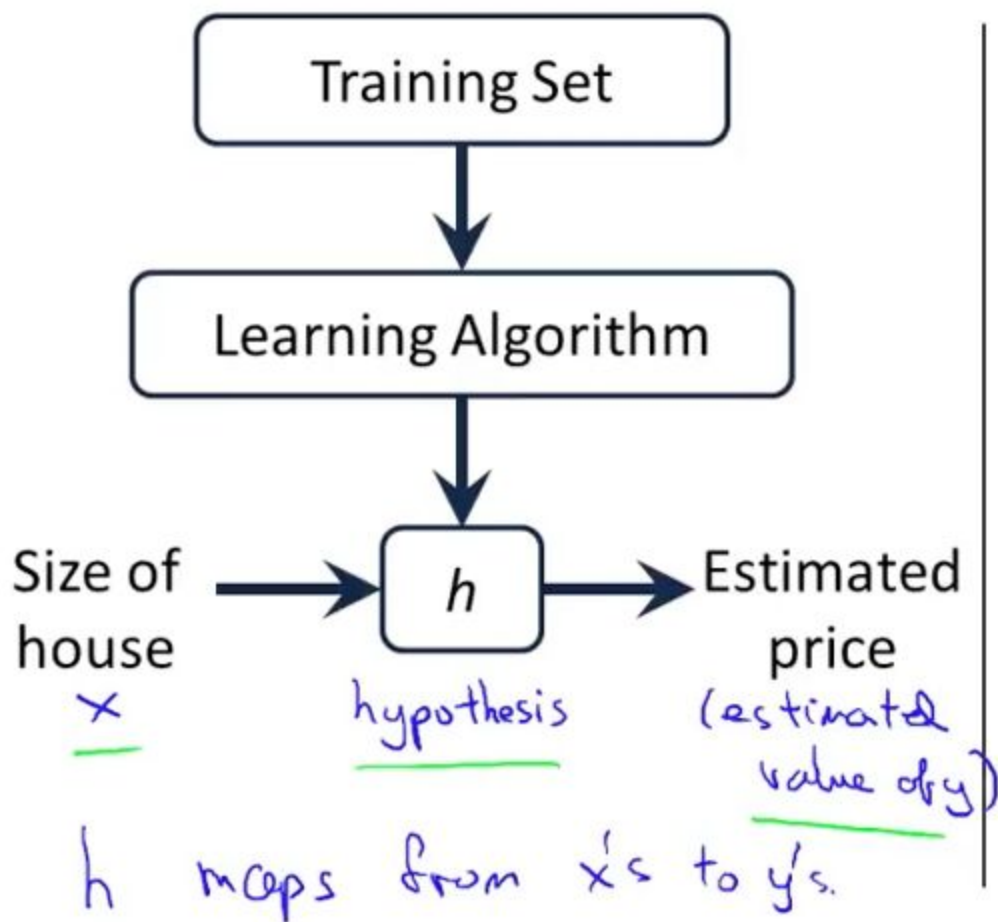| Size in feet² (x) | Price ($) in 1000's (y) |
|---|---|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| ... | ... |

$m = 47$

Notation:

$m$ = Number of training examples

$x$'s = "input" variable / features

$y$'s = "output" variable / "target" variable

$(x, y)$ — one training example

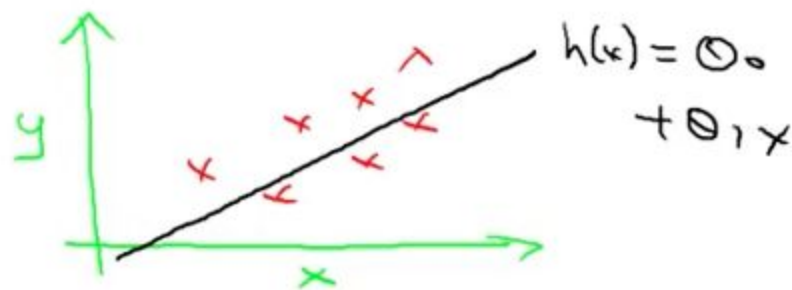$(x^{(i)}, y^{(i)})$ — $i^{th}$ training example

$x^{(1)} = 2104$

$x^{(2)} = 1416$

$y^{(1)} = 460$

Andrew N.

Training Set

↓

Learning Algorithm

↓

Size of house → h → Estimated price

$x$

hypothesis    (estimated value of $y$)

h maps from $x$'s to $y$'s.

**How do we represent $h$ ?**

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

Shorthand: $h(x)$



$h(x) = \theta_0 + \theta_1 x$

Linear regression with one variable.
Univariate linear regression.

Training Set

| Size in feet² (x) | Price ($) in 1000's (y) |
|:---:|:---:|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| ... | ... |

$m = 47$

Hypothesis: $h_\theta(x) = \theta_0 + \theta_1 x$

$\theta_i$'s:  Parameters

Training Set

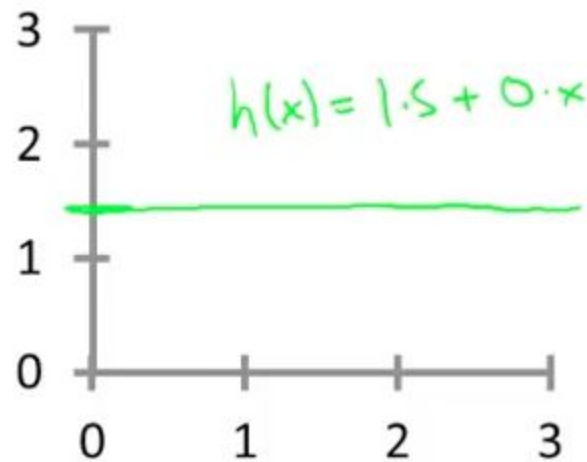| Size in feet² (x) | Price ($) in 1000's (y) |
|---|---|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| ... | ... |

$M = 47$

Hypothesis: $h_\theta(x) = \theta_0 + \theta_1 x$

$\theta_i$'s:     Parameters

How to choose $\theta_i$'s ?
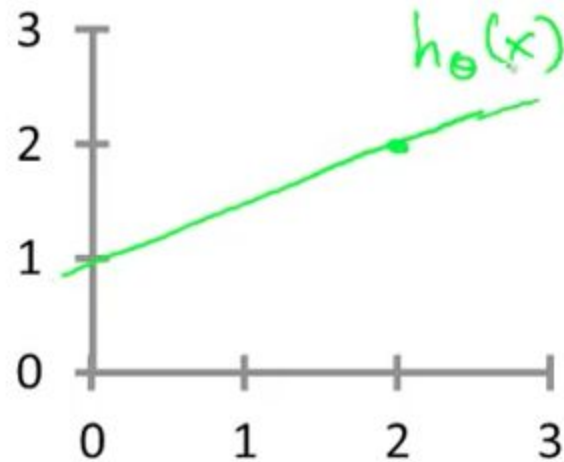
$$h_\theta(x) = \theta_0 + \theta_1 x$$



$h(x) = 1.5 + 0 \cdot x$

$\rightarrow \theta_0 = 1.5$
$\rightarrow \theta_1 = 0$

$h(x) = 0.5x$

$\rightarrow \theta_0 = 0$
$\rightarrow \theta_1 = 0.5$

$h_\theta(x)$

$\rightarrow \theta_0 = 1$
$\rightarrow \theta_1 = 0.5$

$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

$m$ - #training examples

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

$$(x^{(i)}, y^{(i)})$$

**Idea:** Choose $\theta_0, \theta_1$ so that $h_\theta(x)$ is close to $y$ for our training examples $(x, y)$

$x, y$

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

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

Cost function

Squared error function

**Hypothesis:**

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

**Parameters:**

$$\theta_0, \theta_1$$

**Cost Function:**

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

**Goal:** $\underset{\theta_0, \theta_1}{\text{minimize}} \; J(\theta_0, \theta_1)$

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

$$\theta_0 = 0$$

$$\theta_1$$

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

$$\underset{\theta_1}{\text{minimize}} \; J(\theta_1) \qquad \theta_1 x^{(i)}$$

$\rightarrow h_\theta(x)$

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

$\rightarrow J(\theta_1)$

(function of the parameter $\theta_1$)



$h_\theta(x)$

$\theta_1 = 1$

$h_\theta(x^{(i)}) = y^{(i)}$

$\theta_1 = 1$

$J(\theta_1)$

$\theta_1$

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

$= \frac{1}{2m} \sum_{i=1}^{m} (\theta_1 x^{(i)} - y^{(i)})^2 = \frac{1}{2m}(0^2 + 0^2 + 0^2) = 0^2$

$\theta_1 = 0.5?$

$\theta_1$

$J(1) = 0$

$h_\theta(x)$

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

$J(\theta_1)$

(function of the parameter $\theta_1$)

$h_\theta(x)$

$\theta_1 = 0.5$

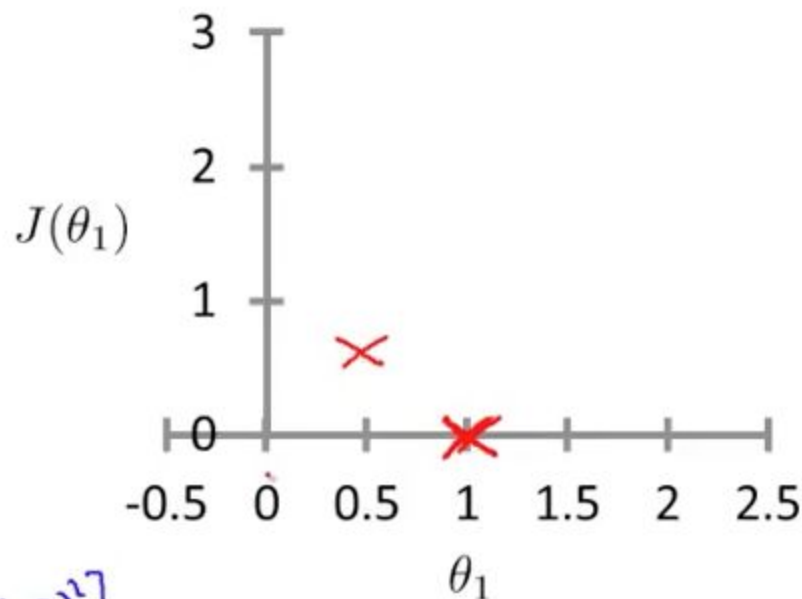$\leftarrow y^{(i)}$

$\leftarrow h_\theta(x^{(i)})$

$J(\theta_1)$

$J(0.5) = \frac{1}{2m}\left[(0.5-1)^2 + (1-2)^2 + (1.5-3)^2\right]$

$= \frac{1}{2\times3}(3.5) = \frac{3.5}{6} \approx 0.58$

# $h_\theta(x)$

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



$\theta_1 = 1$

$\theta_1 = 0$

$J(0) = \frac{1}{2m}(1^2 + 2^2 + 3^2)$

$= \frac{1}{6} \cdot 14 \approx 2 \cdot 3$

# $J(\theta_1)$

(function of the parameter $\theta_1$)

5.25



$J(\theta_1)$

$\theta_1 = 1$

-0.5   0   0.5   1   1.5   2   2.5

$\theta_1$

$h(x) = -0.5x$

$h(x)$

minimize $J(\theta_1)$
$\theta_1$

Andrew N

Hypothesis: $h_\theta(x) = \theta_0 + \theta_1 x$

Parameters: $\theta_0, \theta_1$

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

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

<span style="color:red">*Cost Function을 최소화 하는<br>theta_0, theta_1을 찾는다.*</span>

# $h_\theta(x)$

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

# $J(\theta_0, \theta_1)$

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



$\theta_0 = 50$

$\theta_1 = 0.06$

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

$\theta_1$

$\theta_0, \theta_1$

$h_\theta(x)$

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

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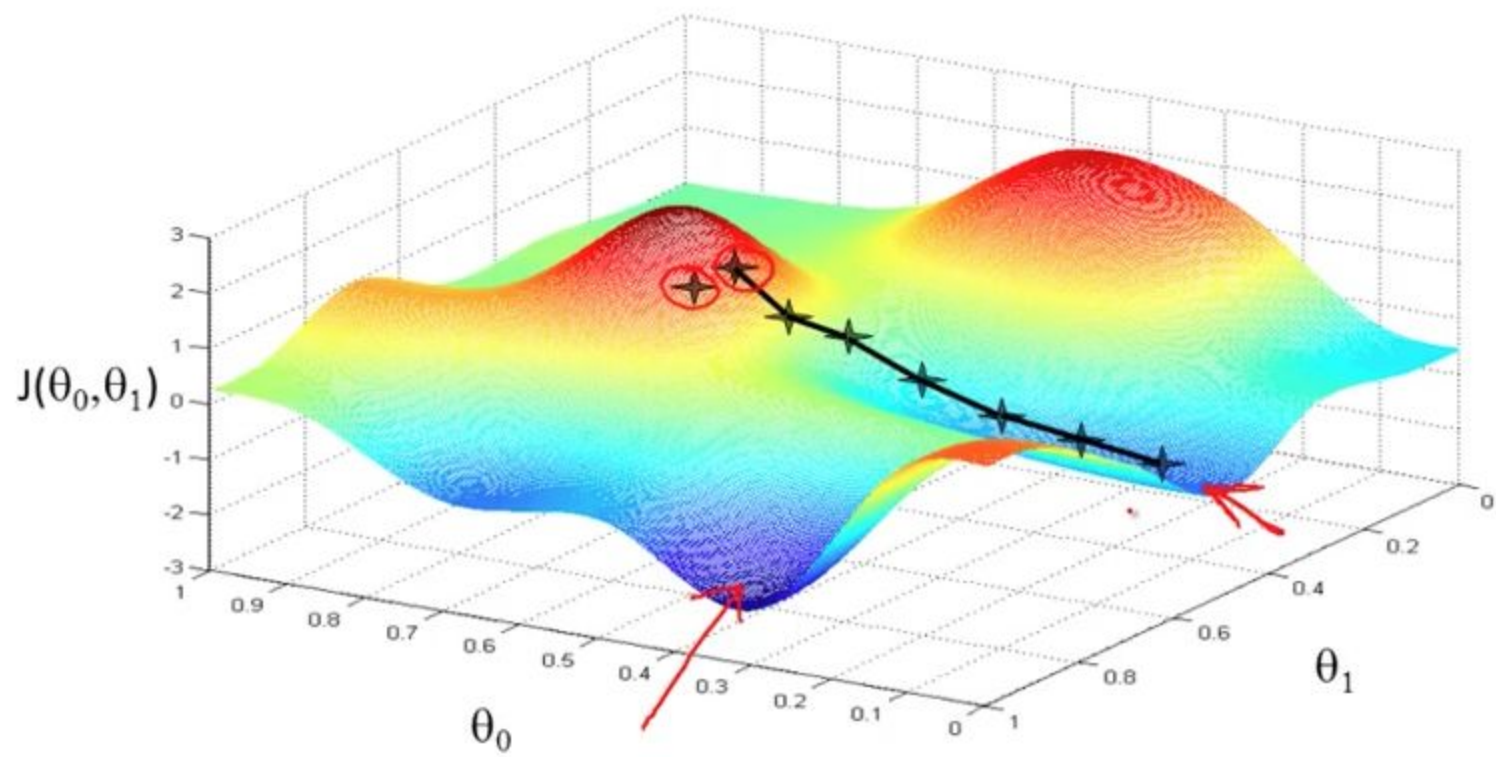
(function of the parameters $\theta_0, \theta_1$)

Have some function $J(\theta_0, \theta_1)$    $J(\theta_0, \theta_1, \theta_2, \ldots, \theta_n)$

Want    $\min\limits_{\theta_0, \theta_1} J(\theta_0, \theta_1)$     $\min\limits_{\theta_0 \cdots \theta_n} J(\theta_0, \ldots, \theta_n)$

**Outline:**

- Start with some $\theta_0, \theta_1$   $(Say \;\; \theta_0 = 0, \theta_1 = 0)$

- Keep changing $\theta_0, \theta_1$ to reduce $J(\theta_0, \theta_1)$

  until we hopefully end up at a minimum

# Gradient descent algorithm

Assignment | Truth assertion

$a := b$

$a := a+1$

$a = b$

$a = a+1$ ✗

$\theta_0, \theta_1$

repeat until convergence {

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

}

learning rate

Simultaneously update $\theta_0$ and $\theta_1$

---

## Correct: Simultaneous update

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

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

$\theta_0 := \text{temp0}$

$\theta_1 := \text{temp1}$

## Incorrect:

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

$\theta_0 := \text{temp0}$

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

$\theta_1 := \text{temp1}$

Andrew Ng

Parameter Learning

Linear regression with one variable

Gradient descent intuition

Machine Learning

# Gradient descent algorithm

repeat until convergence {

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

(simultaneously update
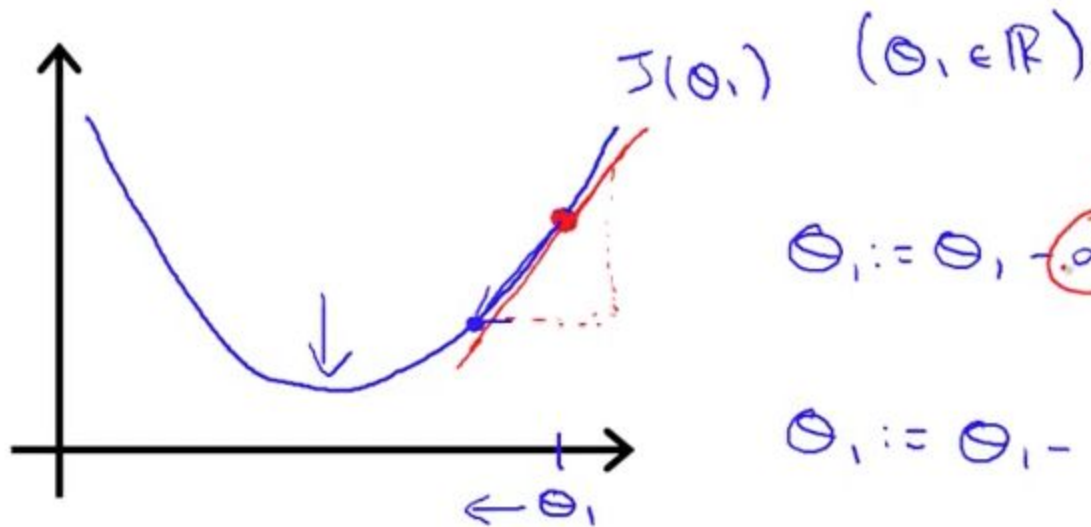$j = 0$ and $j = 1$)

}

learning rate

derivative
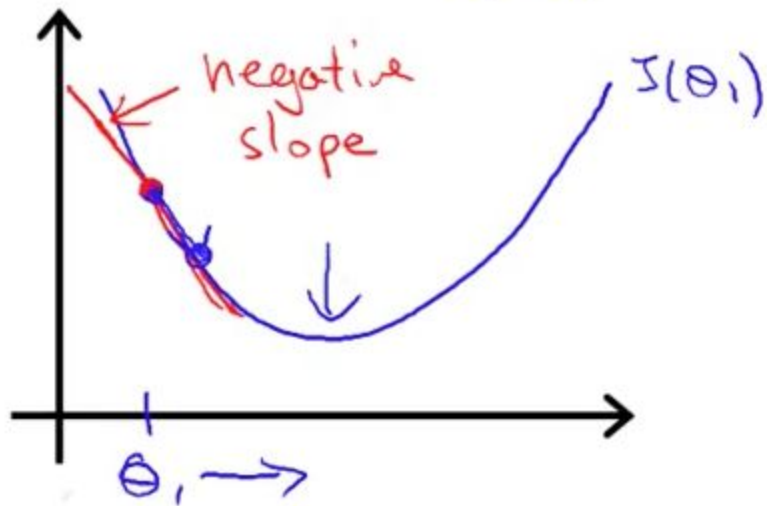
$\min\limits_{\theta_1} J(\theta_1)$

$\theta_1 \in \mathbb{R}.$

$J(\theta_1)$    $(\theta_1 \in \mathbb{R})$

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

$\frac{d}{d\theta_1} \geq 0$

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

negative slope

$$\frac{\frac{d}{d\theta_1} J(\theta_1)}{\leq 0}$$

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

Andrew N

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

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.

slope $= 0$

$\theta_1$ at local optima

Current value of $\theta_1$

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

$\Theta_1 := \Theta_1 - \alpha \cdot 0$

$= 0$

$\Theta_1 := \Theta_1$

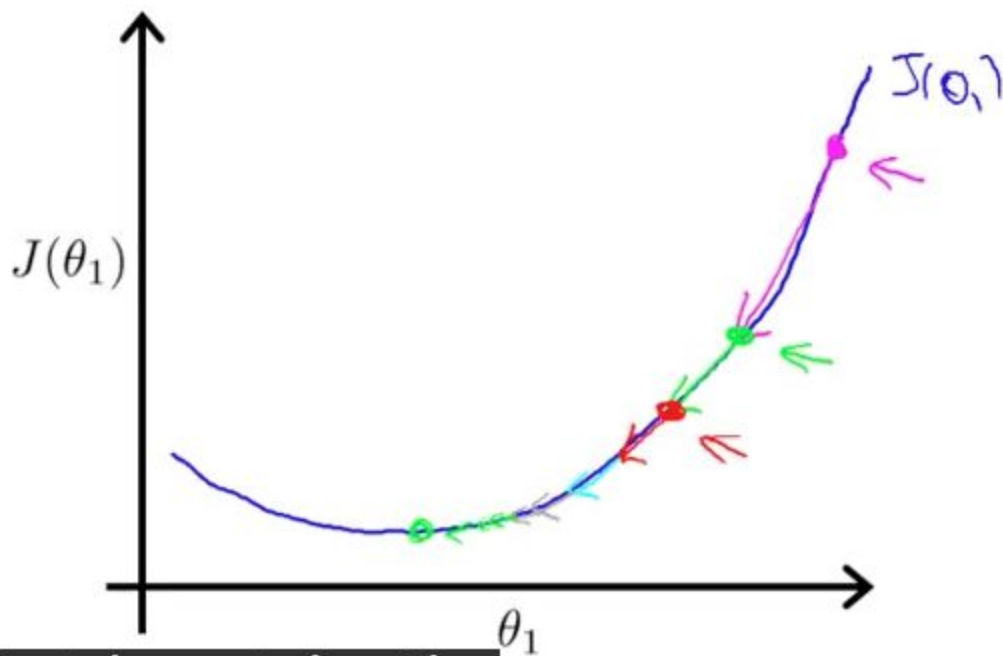Gradient descent can converge to a local minimum, even with the learning rate α fixed.

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

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

$J(\theta_1)$

$J(\theta_1)$

$\theta_1$

So that's the gradient descent algorithm and you can use it to try to minimize

**Parameter Learning**

Linear regression
with one variable

Gradient descent for
linear regression

Machine Learning

## Gradient descent algorithm

repeat until convergence {

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

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

}

## Linear Regression Model

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

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

$$\underset{\theta_0, \theta_1}{\text{Min}} \; J(\theta_0, \theta_1)$$

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

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

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

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

# Gradient descent algorithm

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

repeat until convergence {

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

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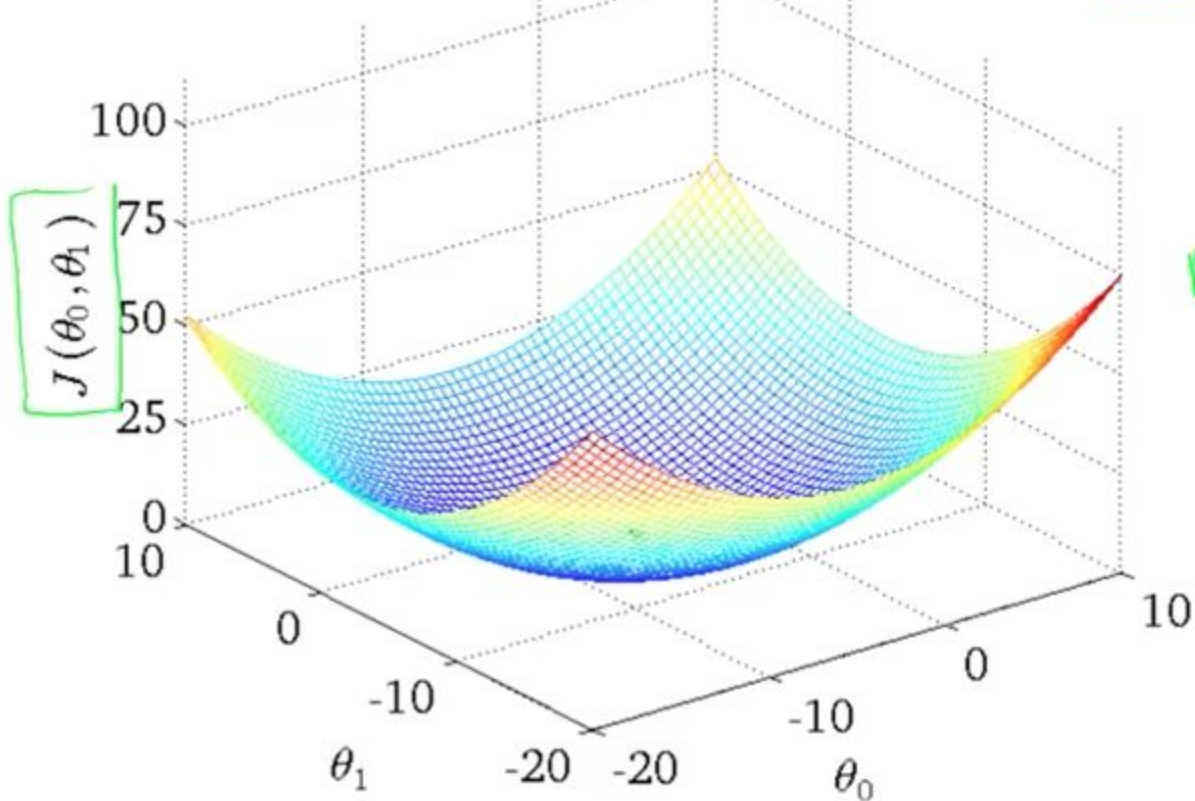
}

update $\theta_0$ and $\theta_1$ simultaneously

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

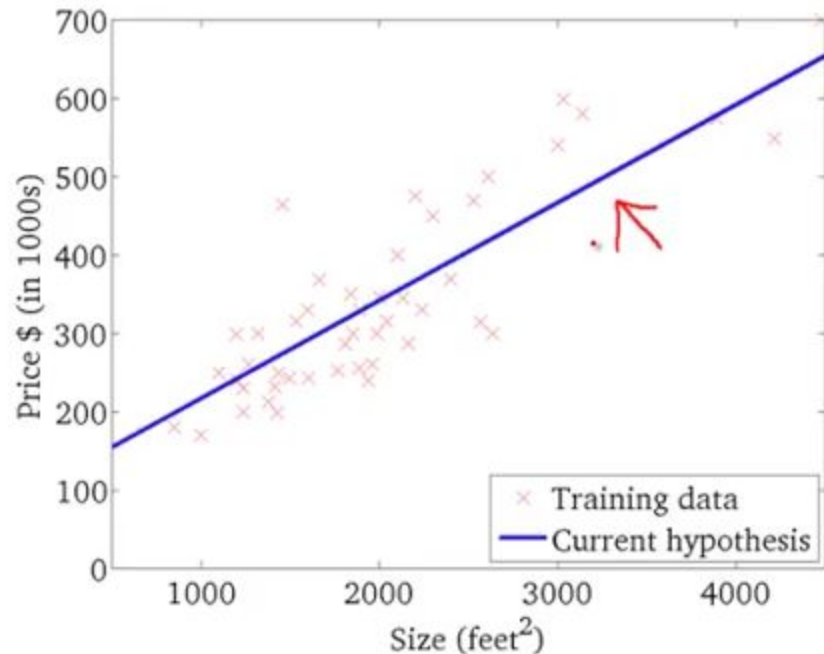*Linear regression*의 *cost function*은
*convex function*이므로 최저값을 갖는다.

$h_\theta(x)$

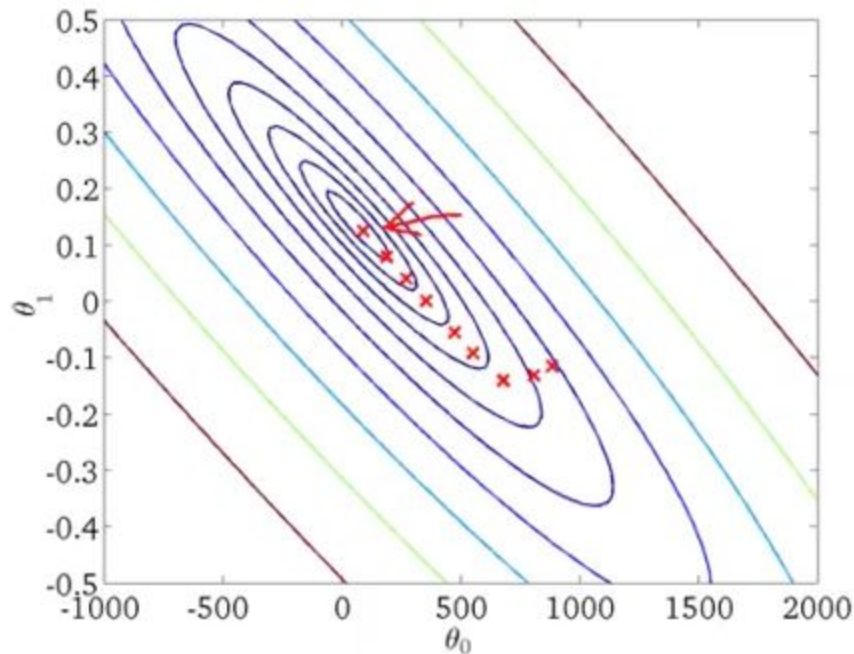(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

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

# "Batch" Gradient Descent

"Batch": Each step of gradient descent uses all the training examples.

$$\rightarrow \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)$$

우리가 다루는 *'Batch'* 외에 서브셋만 다루는 등의
다른 *Gradient Descent* 알고리즘이 존재한다.

Linear Algebra
review (optional)

Matrices and
vectors   *Skipping...*

Machine Learning