

Notes

Machine Learning by Andrew Ng on Coursera

Sparsh Jain

September 23, 2020

Contents

1	Introduction	2
1.1	Supervised Learning	2
1.1.1	Regression Problem	2
1.1.2	Classification Problems	2
1.2	Unsupervised Learning	2
2	Linear Regression with One Variable	3
2.1	Notations	3
2.2	Supervised Learning	3
2.3	Gradient Descent	4
2.4	Gradient Descent for Linear Regression	4
	Appendices	5
	A Lecture 1	6
	B Lecture 2	46

Chapter 1

Introduction

Machine learning (task, experience, performance) can be classified into *Supervised* and *Unsupervised* learning.

1.1 Supervised Learning

Supervised learning can be basically classified into *Regression* and *Classification* problems.

1.1.1 Regression Problem

Regression problems work loosely on continuous range of outputs.

1.1.2 Classification Problems

Classification problems work loosely on discrete range of outputs.

1.2 Unsupervised Learning

An example is *Clustering Problem*.

Check Appendix A for more details.

Chapter 2

Linear Regression with One Variable

2.1 Notations

m = number of training examples

x 's = 'input' variables / features

y 's = 'output' variables / 'target' variables

(x, y) = single training example

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

2.2 Supervised Learning

We have a data set (*Training Set*).

Training Set \rightarrow Learning Algorithm $\rightarrow h$ (*hypothesis*, a function $X \rightarrow Y$)

To Represent h

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

Cost

$$\underset{\theta_0, \theta_1}{\text{minimize}} \frac{1}{2m} \sum_1^m (h_{\theta}(x) - y)^2$$

Cost Function

Squared Error Cost Function

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

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

2.3 Gradient Descent

Finds local optimum:

1. Start with some value
2. Get closer to optimum

Algorithm

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad \forall j$$

where α = learning rate

Important!

Simultaneous Update!

$$\begin{aligned} temp_j &:= \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad \forall j \\ \theta_j &:= temp_j \quad \forall j \end{aligned}$$

2.4 Gradient Descent for Linear Regression

Cost function for linear regression is convex!

Batch Gradient Descent: Each step of gradient descent uses all training examples.

Check Appendix B for more details.

Appendices

Appendix A

Lecture 1

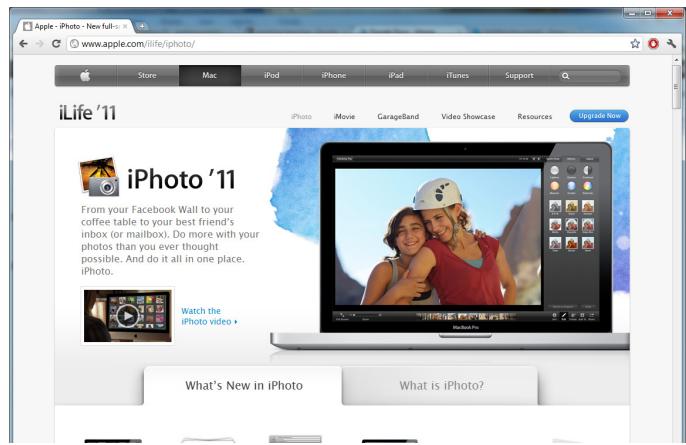


Machine Learning

Introduction

Welcome

Andrew Ng



Andrew Ng



Andrew Ng

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.

Machine Learning

- Grew out of work in AI

Exam

- I
- I
- I



ig
host of

Andrew Ng

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.

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

Andrew Ng

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).

Andrew Ng

Andrew Ng



Machine Learning

Introduction

What is machine learning

Machine Learning definition

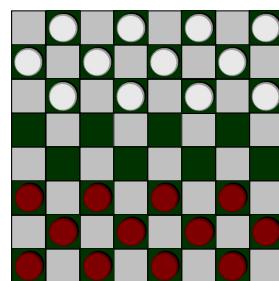
Andrew Ng

Machine Learning definition

- Arthur Samuel (1959). Machine Learning: Field of study that gives computers the ability to learn without being explicitly programmed.

Machine Learning definition

- Arthur Samuel (1959). Machine Learning: Field of study that gives computers the ability to learn without being explicitly programmed.



Andrew Ng

Machine Learning definition

- Arthur Samuel (1959). Machine Learning: Field of study that gives computers the ability to learn without being explicitly programmed.
- Tom Mitchell (1998) Well-posed Learning Problem: A computer program is said to *learn* from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E.

Andrew Ng

"A computer program is said to *learn* from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E."

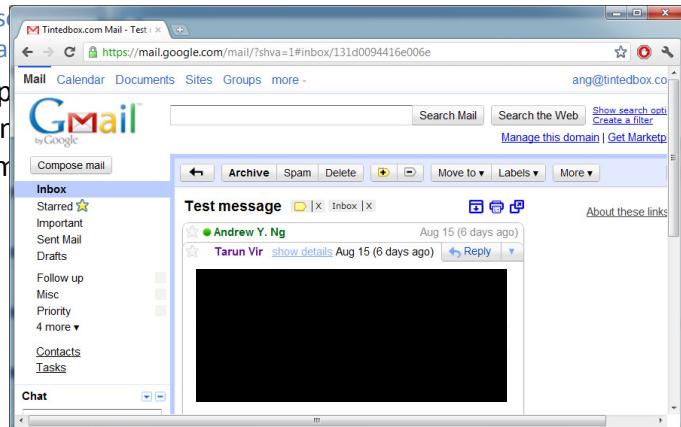
Suppose your email program watches which emails you do or do not mark as spam, and based on that learns how to better filter spam. What is the task T in this setting?

- Classifying emails as spam or not spam. $T \leftarrow$
- Watching you label emails as spam or not spam. $E \leftarrow$
- The number (or fraction) of emails correctly classified as spam/not spam.
- None of the above—this is not a machine learning problem. $P \leftarrow$

"A computer program is said to *learn* from experience E with respect to

on T,

Support
not re-
spam



or do
filter

spam.

"A computer program is said to *learn* from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E."

Suppose your email program watches which emails you do or do not mark as spam, and based on that learns how to better filter spam. What is the task T in this setting?

- Classifying emails as spam or not spam. $T \leftarrow$
- Watching you label emails as spam or not spam. $E \leftarrow$
- The number (or fraction) of emails correctly classified as spam/not spam.
- None of the above—this is not a machine learning problem. $P \leftarrow$

Machine learning algorithms:

- Supervised learning
 - Unsupervised learning
- 

Others: Reinforcement learning, recommender systems.

Also talk about: Practical advice for applying learning algorithms.



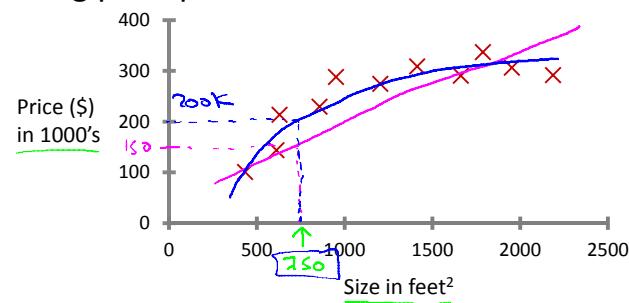
Andrew Ng



Machine Learning

Introduction Supervised Learning

Housing price prediction.

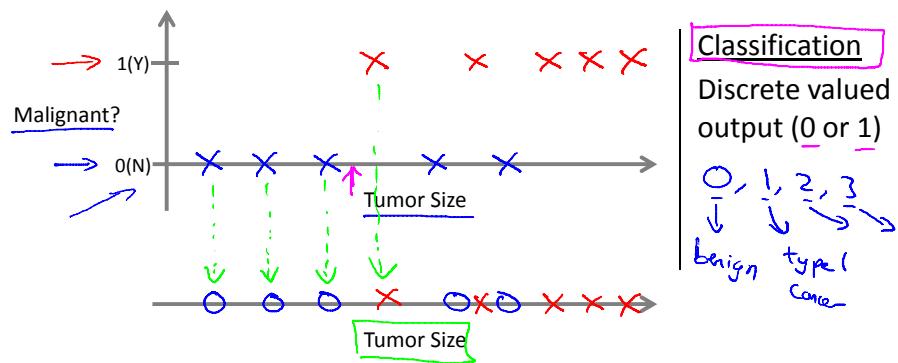


Supervised Learning
"right answers" given

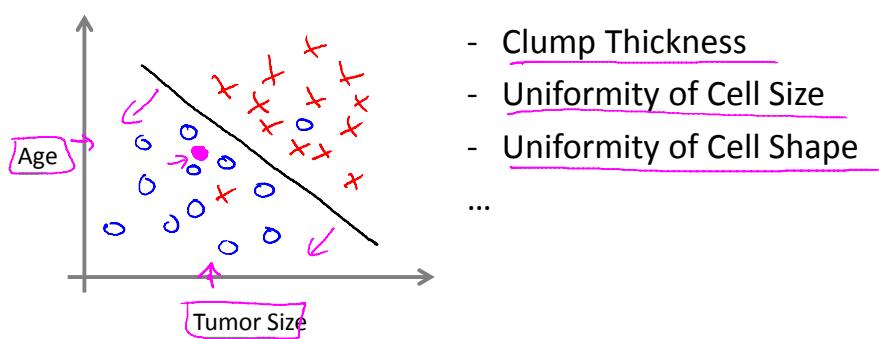
Regression: Predict continuous valued output (price)

Andrew Ng

Breast cancer (malignant, benign)



Andrew Ng



Andrew Ng

You're running a company, and you want to develop learning algorithms to address each of two problems.

1000's

→ Problem 1: You have a large inventory of identical items. You want to predict how many of these items will sell over the next 3 months.

→ Problem 2: You'd like software to examine individual customer accounts, and for each account decide if it has been hacked/compromised.

→ 0 - not hacked
→ 1 - hacked

Should you treat these as classification or as regression problems?

- Treat both as classification problems.
- Treat problem 1 as a classification problem, problem 2 as a regression problem.
- Treat problem 1 as a regression problem, problem 2 as a classification problem.
- Treat both as regression problems.

Andrew Ng

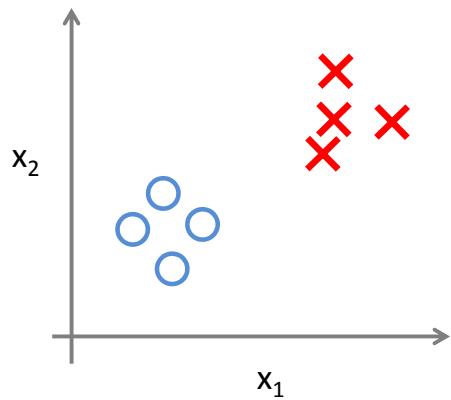


Machine Learning

Introduction

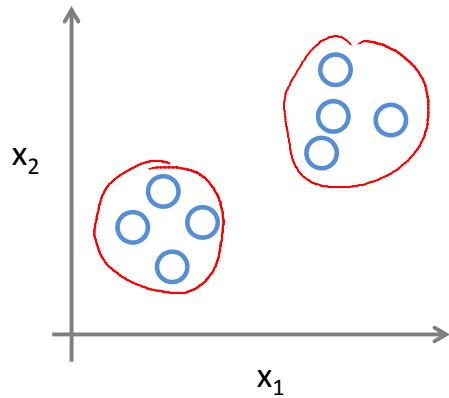
Unsupervised Learning

Supervised Learning



Andrew Ng

Unsupervised Learning



Andrew Ng

The screenshot shows the Google News homepage. A red arrow points from the address bar at the top left to the URL 'news.google.com'. Another red arrow points from the left sidebar to the 'BP Oil Well' headline. The main content area displays several news stories, with the 'BP Oil Well' story highlighted by a red box. The story title is 'BP Oil Well, Site of National Catastrophe, Dies at One'. Below the title is a small image of a burning oil rig.

Top Stories

- Deepwater Horizon
- Fed meeting
- Foreign exchange market
- Lindsay Lohan
- IBM
- Tim Brady
- Toronto International Film Festival
- Paris Hilton
- Iran
- Hurricane Igor
- Starred
- San Francisco Bay Area
- World
- U.S.
- Business
- Sci/Tech
- More Top Stories
- Sport
- Health
- Sports
- Entertainment

All news

- Headlines
- Images

Top Stories

- Christine O'Donnell » **White House official denies Tea Party-focused ad campaign** CNN.com - 1 hour ago
- Democratic sources say the White House is not considering an ad campaign tying Republicans to the Tea Party. Washington (CNN) -- A top White House official says the White House is not considering an ad campaign that claims President Obama's political advisers are weighing a national Tea Party is misplaced the blame, former President Bill Clinton claims. (USA Today)
- GOP tea party backer defends Christine O'Donnell The Associated Press Atlantic Journal Constitution - Politico Daily - MyFox Washington DC - Salon (Salon)
- US Stocks Climb After Recession Called Over, Homebuilders Gain MarketWatch - Kristina Peterson - 16 minutes ago
- REC
- HELIOS Energy Corp.'s US stocks climbed Monday, gaining speed after a key nonprofit organization officially called the recession over, giving investors a boost of confidence in the gradual economic recovery. Los Angeles Times - 1 hour ago
- Downturn Was Longest in Decades, Panel Confirms New York Times (New York Times) - AP - CNN - USA Today
- Deepwater Horizon

BP Oil Well, Site of National Catastrophe, Dies at One

The BP oil well, site of the Deepwater Horizon explosion that led to the worst oil spill in US history, died today at one year old. By Matt Weisert and Michael S. Weiss Doubts BP Would End Operations in Gulf of Mexico: Video Bloomberg (CNN International) - Wall Street Journal (blog) - The Guardian - New York Times - AP - USA Today - all 2,292 news articles

Recent

- Recession officially ended in June 2009
- CNNMoney - Chris Isidore - 39 minutes ago
- Hurricane Igor lashes Bermuda
- USA Today - Gerry Broome - 5 minutes ago
- Explain what you want from us - reads front-page editorial
- msnbc.com - Olivia Torres - 10 minutes ago

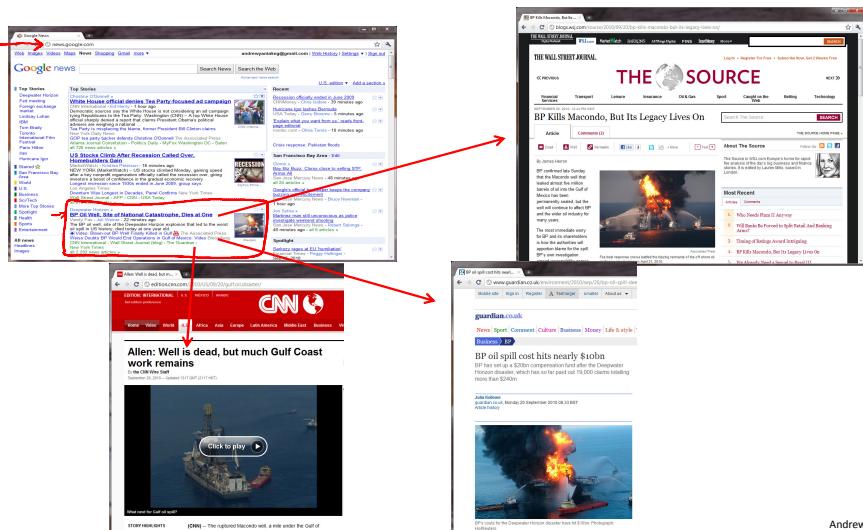
San Francisco Bay Area - Edit

- Chrono »
- Rev. Big Buzz - Clooney close to selling STP.
- Armor All
- San Jose Mercury News - 48 minutes ago - all 24 articles
- Google's official beekeeper keeps the company buzzing with excitement
- San Jose Mercury News - Bruce Newman - 1 hour ago
- Jon Sylvia »
- Martinez man still unconscious as police investigate weekend shooting
- San Jose Mercury News - Robert Salonga - 48 minutes ago - all 6 articles

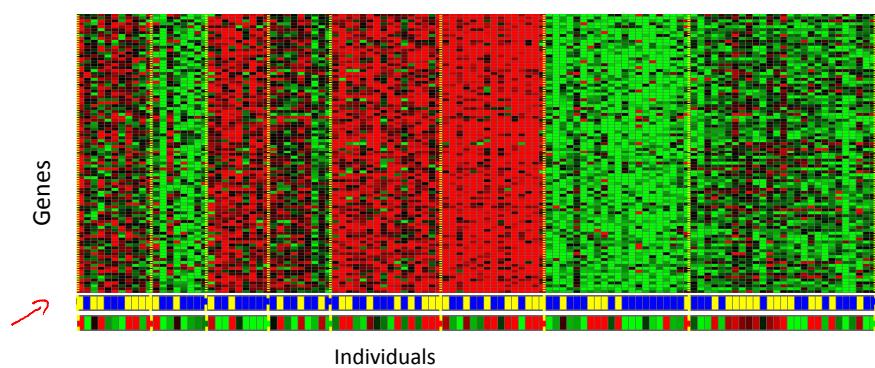
Spotlight

- Sarkozy: Pages at EU 'humiliation'
- Financial Times - Peggy Hollinger - Sep 16, 2010

Andrew Ng

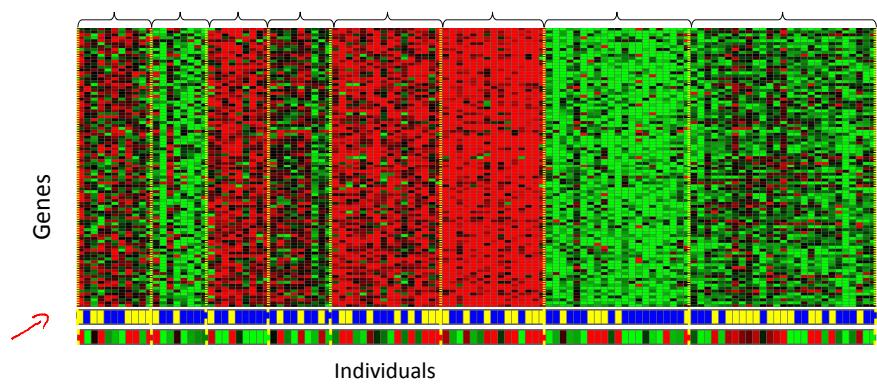


Andrew Ng



[Source: Daphne Koller]

Andrew Ng



[Source: Daphne Koller]

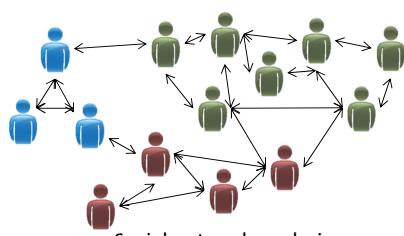
Andrew Ng



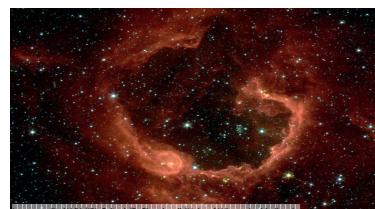
Organize computing clusters



Market segmentation



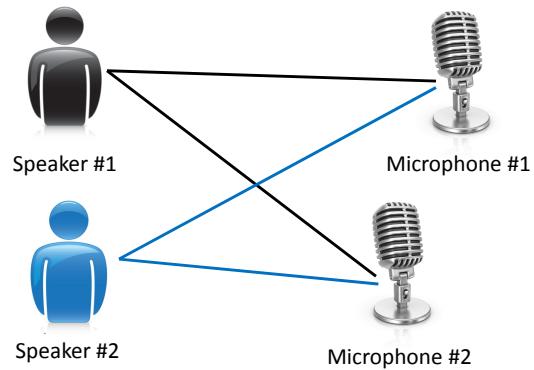
Social network analysis



Astronomical data analysis

Andrew Ng

Cocktail party problem



Andrew Ng

Microphone #1:  Output #1: 

Microphone #2:  Output #2: 

Microphone #1:  Output #1: 

Microphone #2:  Output #2: 

[Audio clips courtesy of Te-Won Lee.]

Andrew Ng

Cocktail party problem algorithm

```
[W,s,v] = svd((repmat(sum(x.*x,1),size(x,1),1).*x)*x');
```

[Source: Sam Roweis, Yair Weiss & Eero Simoncelli]

Andrew Ng

Of the following examples, which would you address using an unsupervised learning algorithm? (Check all that apply.)

- Given email labeled as spam/not spam, learn a spam filter.
- Given a set of news articles found on the web, group them into set of articles about the same story.
- Given a database of customer data, automatically discover market segments and group customers into different market segments.
- Given a dataset of patients diagnosed as either having diabetes or not, learn to classify new patients as having diabetes or not.

Andrew Ng

Appendix B

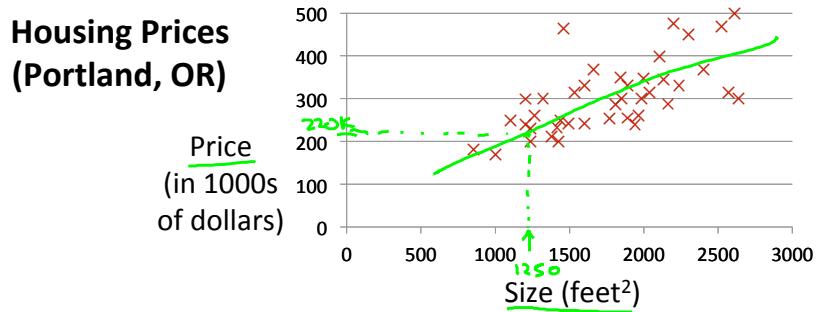
Lecture 2



Machine Learning

Linear regression
with one variable

Model
representation



Supervised Learning

Given the “right answer” for each example in the data.

Regression Problem

Predict real-valued output

Classification: Discrete-valued output

Andrew Ng

<u>Training set of housing prices (Portland, OR)</u>	<u>Size in feet²(x)</u>	<u>Price (\$)² in 1000's(y)</u>
	→ 2104	460
	→ 1416	232
	→ 1534	315
	852	178

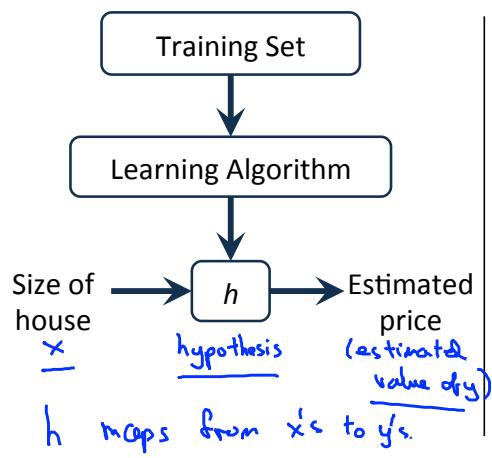
Notation:

- 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

$$\begin{cases} x^{(1)} = 2104 \\ x^{(2)} = 1416 \\ y^{(1)} = 460 \end{cases}$$

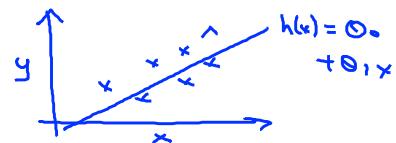
Andrew Ng



How do we represent h ?

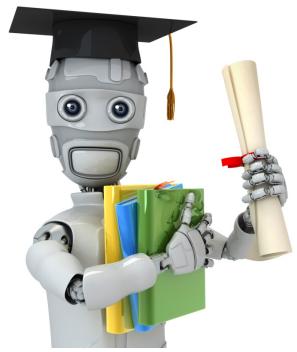
$$h_{\theta}(x) = \underline{\theta_0 + \theta_1 x}$$

Shorthand: $h(x)$



Linear regression with one variable. (x)
Univariate linear regression.
one variable

Andrew Ng



Machine Learning

Linear regression
with one variable

Cost function

Training Set	Size in feet ² (x)	Price (\$) in 1000's (y)
	2104	460
	1416	232
	1534	315
	852	178

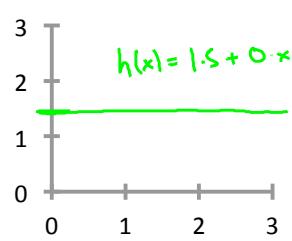
$\left. \begin{matrix} \\ \\ \\ \\ \end{matrix} \right\} m=47$

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

θ_i 's: Parameters

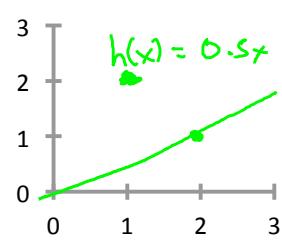
How to choose θ_i 's ?

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



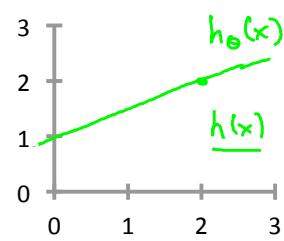
$$\begin{aligned} \rightarrow \theta_0 &= 1.5 \\ \rightarrow \theta_1 &= 0 \end{aligned}$$

$$h(x) = 0.5x$$



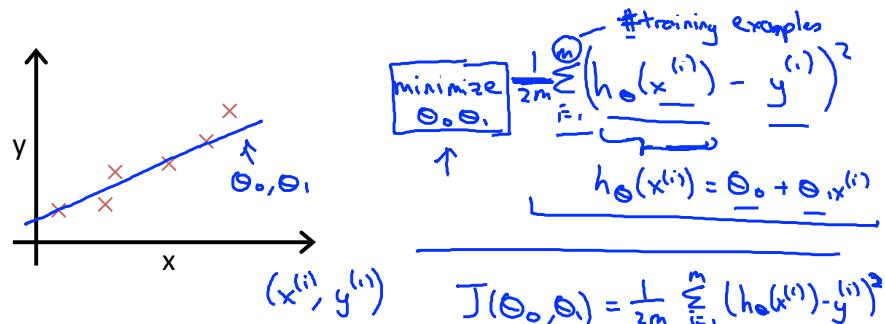
$$\begin{aligned} \rightarrow \theta_0 &= 0 \\ \rightarrow \theta_1 &= 0.5 \end{aligned}$$

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



$$\begin{aligned} \rightarrow \theta_0 &= 1 \\ \rightarrow \theta_1 &= 0.5 \end{aligned}$$

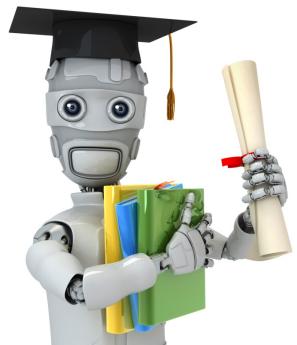
Andrew Ng



Idea: Choose $\underline{\theta}_0, \underline{\theta}_1$ so that
 $\underline{h}_\theta(\underline{x})$ is close to \underline{y} for our
 training examples $(\underline{x}, \underline{y})$

$\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$
 Cost function
 Squared error function

Andrew Ng



Machine Learning

Linear regression
with one variable

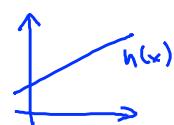
Cost function
intuition I

Hypothesis:

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

Parameters:

$$\underline{\theta_0, \theta_1}$$



Cost Function:

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

Goal: minimize $J(\theta_0, \theta_1)$

$$\nearrow \theta_0, \theta_1$$

Simplified

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

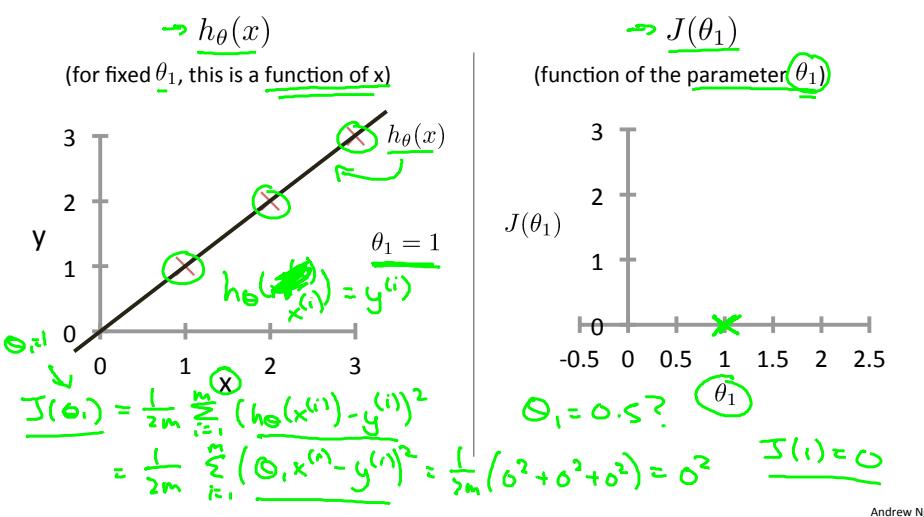
$$\theta_0 = 0$$

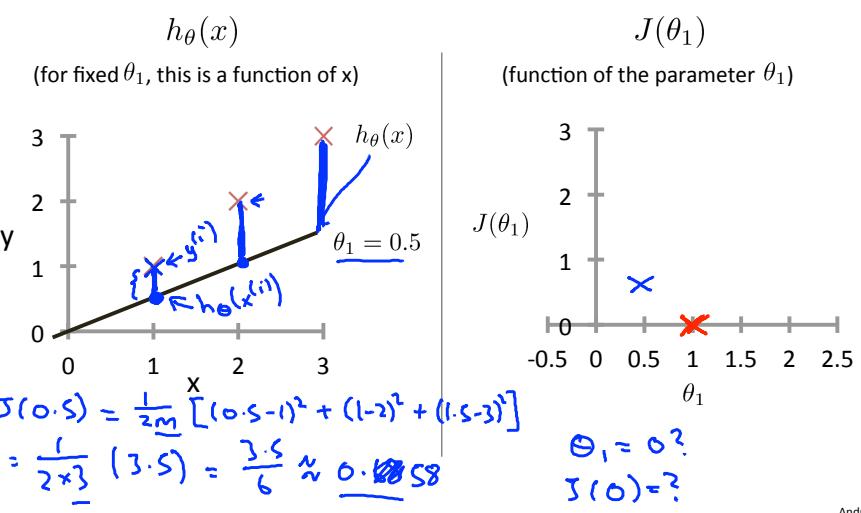


$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

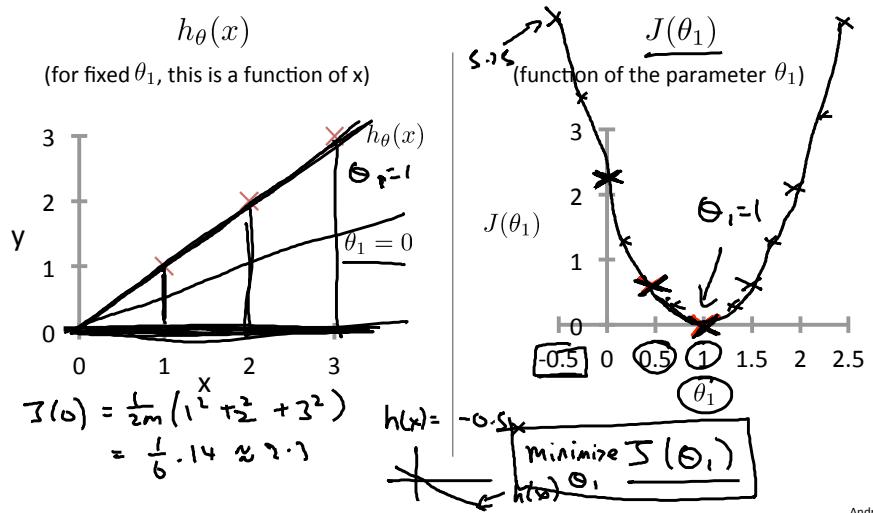
$$\text{minimize } \underline{\theta_1} \quad \theta_1, x^{(i)}$$

Andrew Ng





Andrew Ng





Machine Learning

Linear regression
with one variable

Cost function
intuition II

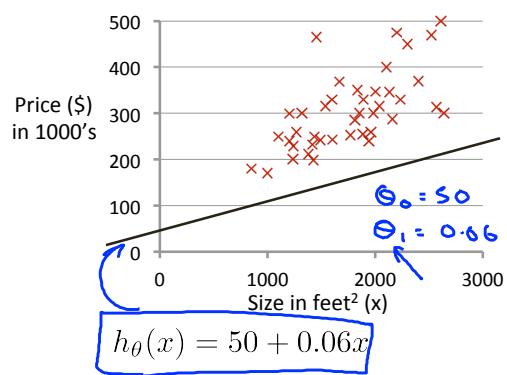
Hypothesis: $h_{\theta}(x) = \theta_0 + \theta_1 x$

Parameters: θ_0, θ_1

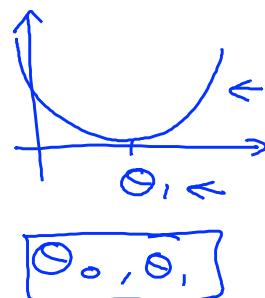
Cost Function: $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Goal: minimize $J(\theta_0, \theta_1)$

$\underline{h_\theta(x)}$
(for fixed θ_0, θ_1 , this is a function of x)

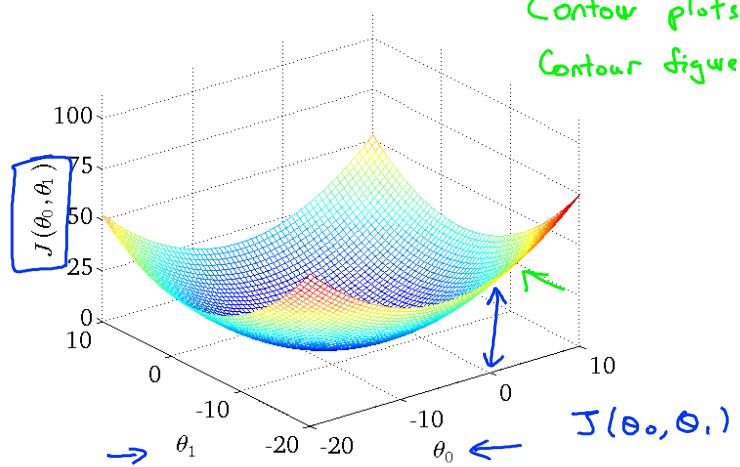


$\underline{J(\theta_0, \theta_1)}$
(function of the parameters θ_0, θ_1)

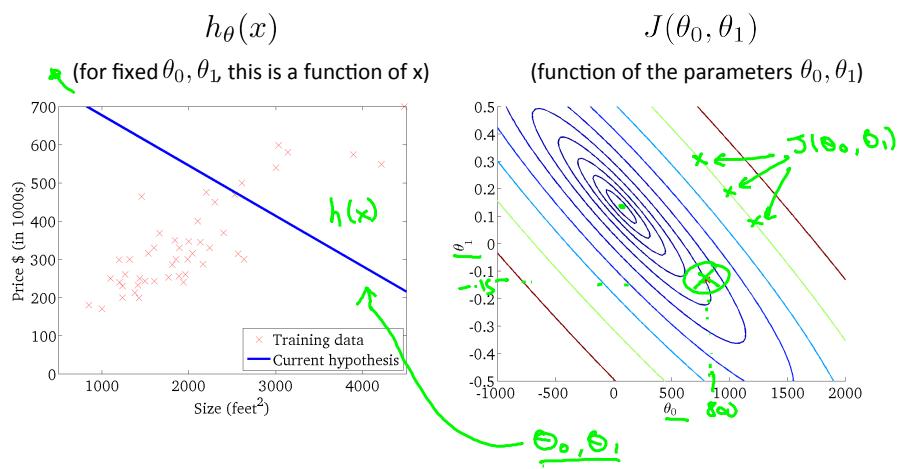


Andrew Ng

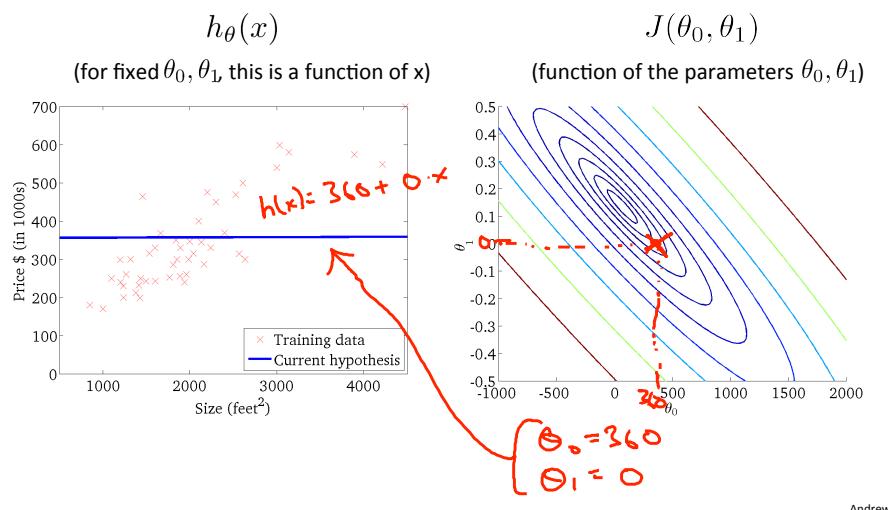
Contour plots
Contour figures -



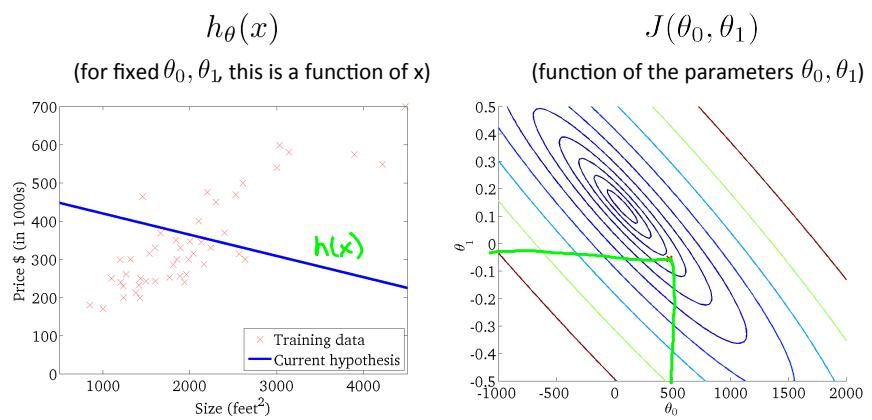
Andrew Ng



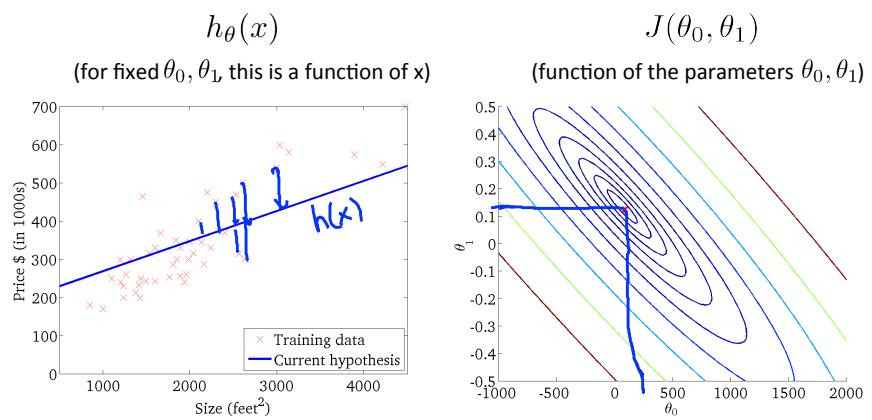
Andrew Ng



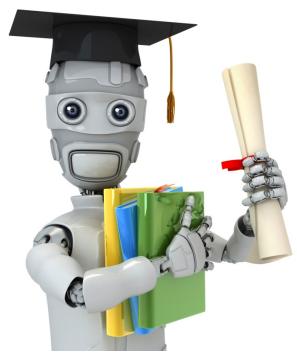
Andrew Ng



Andrew Ng



Andrew Ng



Machine Learning

Linear regression
with one variable

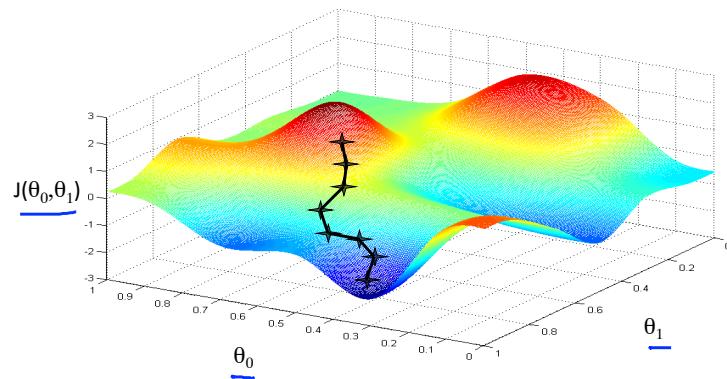
Gradient
descent

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

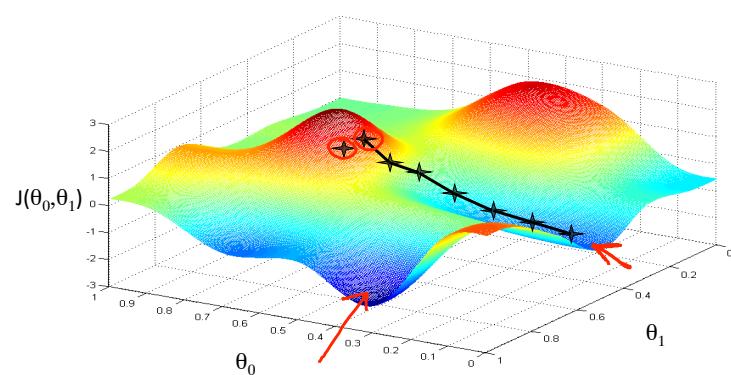
Want $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$ $\min_{\theta_0, \dots, \theta_n} J(\theta_0, \dots, \theta_n)$

Outline:

- Start with some θ_0, θ_1 (say $\theta_0 = 0, \theta_1 = 0$)
- Keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$
until we hopefully end up at a minimum



Andrew Ng



Andrew Ng

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 $j = 1$)
learning rate

Assignment
 $a := b$ Truth assertion
 $a := a + 1$ $a = b \Leftarrow$
 $a = a + 1 \times$

Simultaneously update
 θ_0 and θ_1

Correct: Simultaneous update

```

    → temp0 :=  $\theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$ 
    → temp1 :=  $\theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$ 
    →  $\theta_0 := \text{temp0}$ 
    →  $\theta_1 := \text{temp1}$ 
    
```

Incorrect:

```

    → temp0 :=  $\theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$ 
    →  $\theta_0 := \text{temp0}$ 
    → temp1 :=  $\theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$   $\leftarrow$ 
    →  $\theta_1 := \text{temp1}$ 
    
```

Andrew Ng



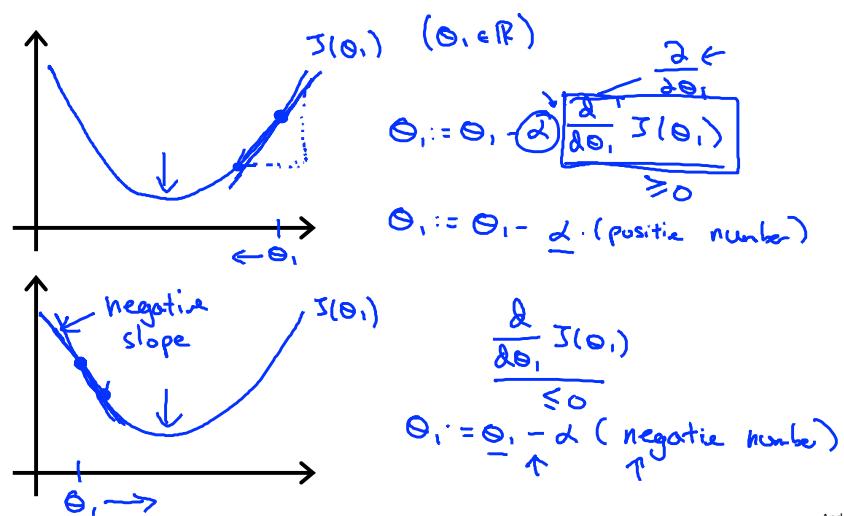
Machine Learning

Linear regression
with one variable

Gradient descent
intuition

Gradient descent algorithm

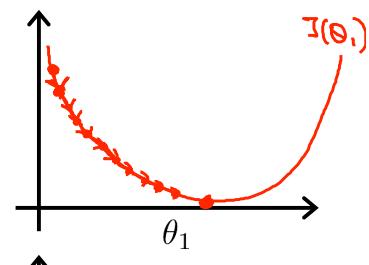
repeat until convergence {
 $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$ (simultaneously update
 }
 ↑ learning rate ↑ derivative
 $j = 0$ and $j = 1$)
 $\min_{\theta_1} J(\theta_1) \quad \theta_1 \in \mathbb{R}$.



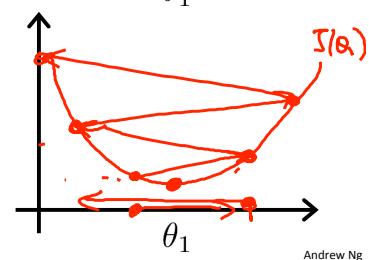
Andrew Ng

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

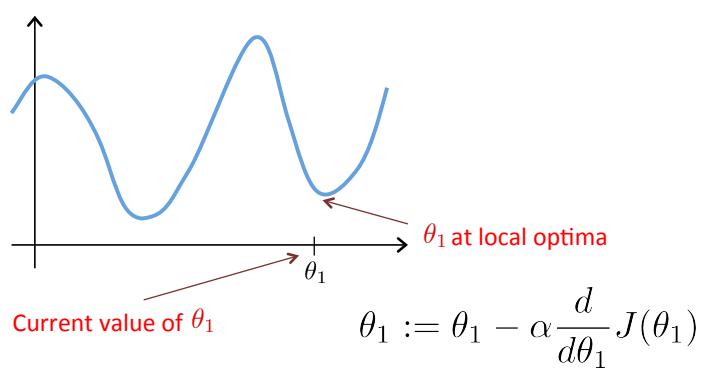
If α is too small, gradient descent can be slow.



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



Andrew Ng

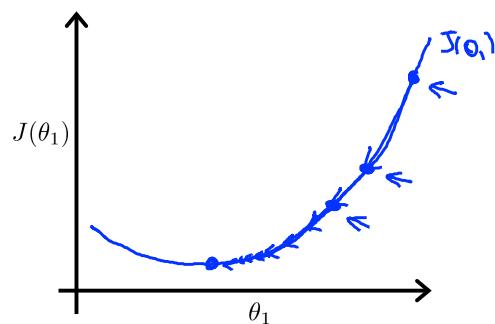


Andrew Ng

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.



Andrew Ng



Machine Learning

Linear regression
with one variable

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  $j = 1$  and  $j = 0$ )  
}
```

Linear Regression Model

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

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

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

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

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

Andrew Ng

Gradient descent algorithm

repeat until convergence {

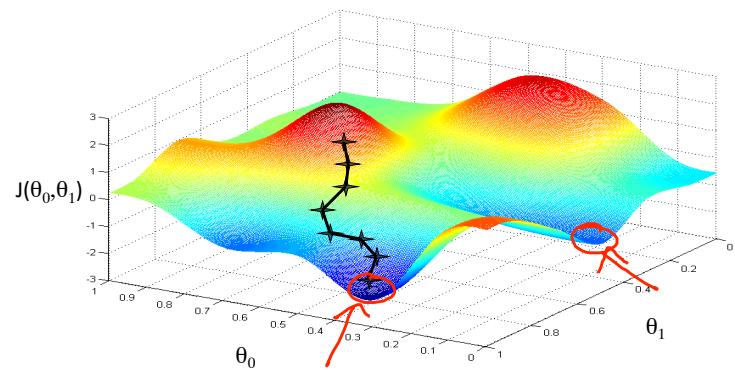
$$\theta_0 := \theta_0 - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \right]$$

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

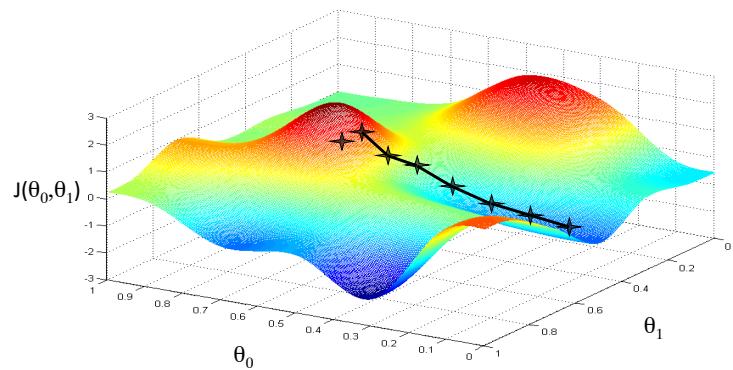
}

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

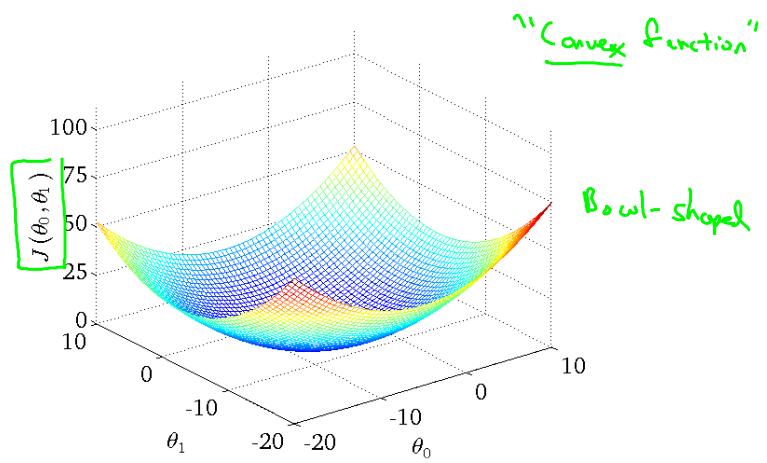
update θ_0 and θ_1 simultaneously



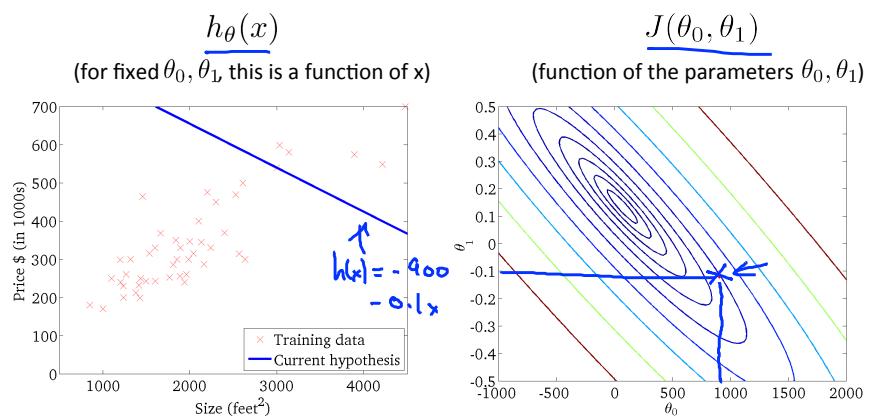
Andrew Ng



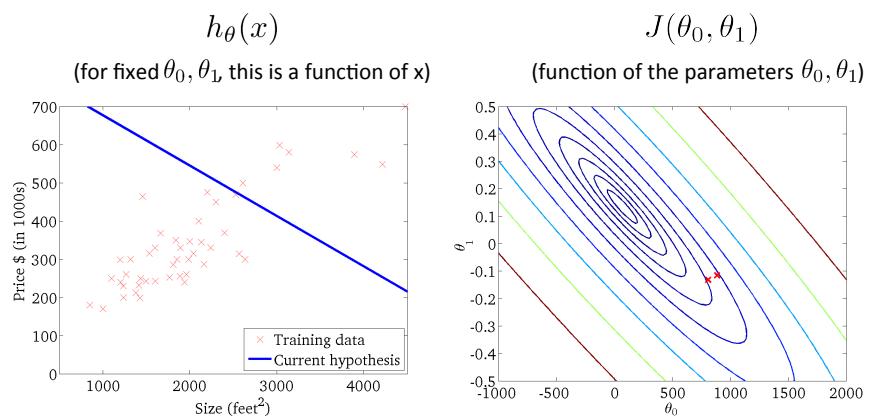
Andrew Ng



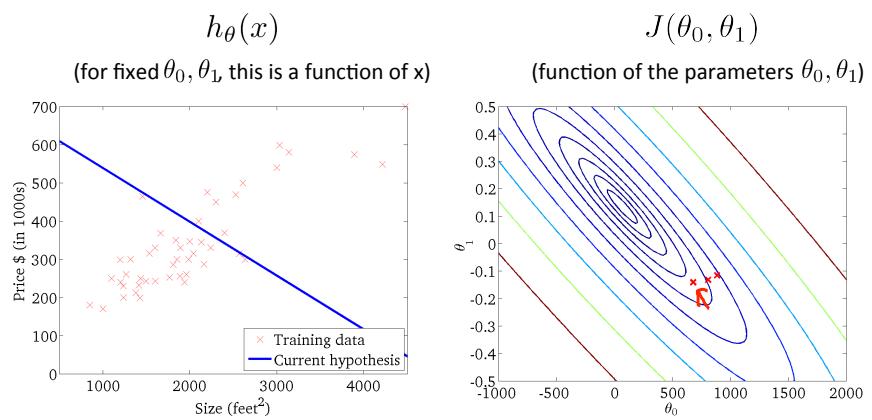
Andrew Ng



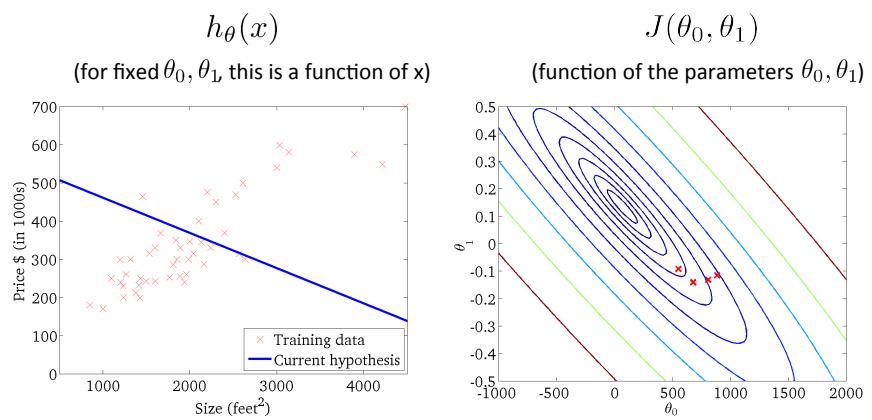
Andrew Ng



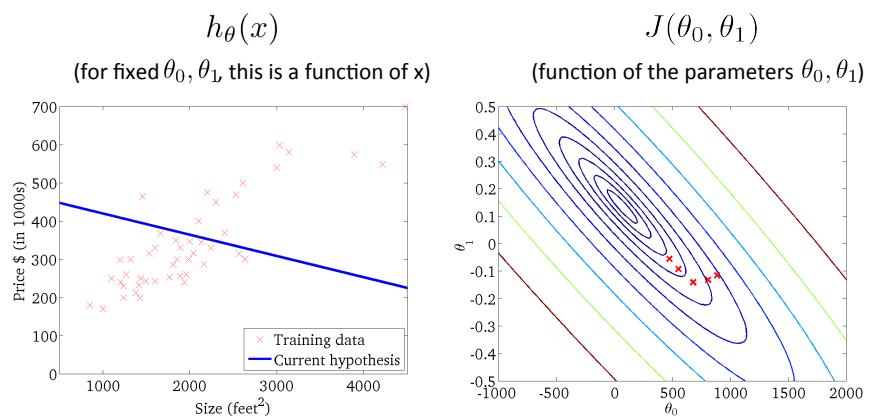
Andrew Ng



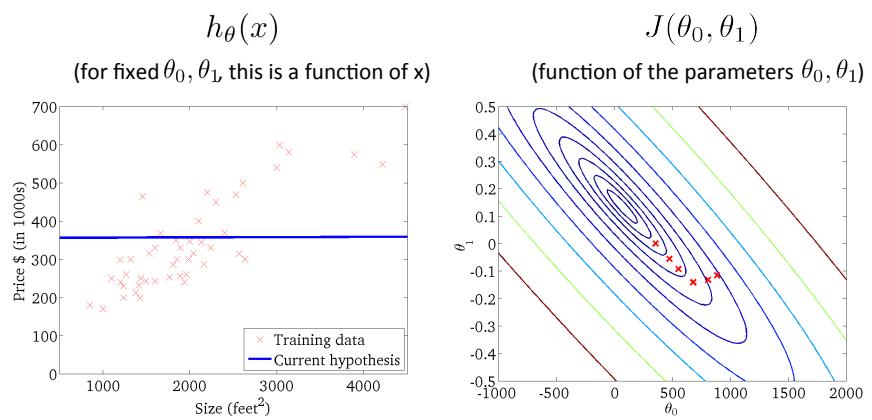
Andrew Ng



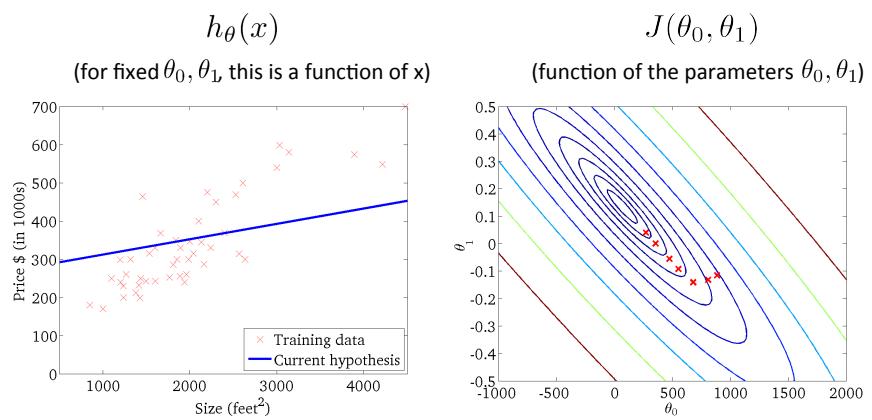
Andrew Ng



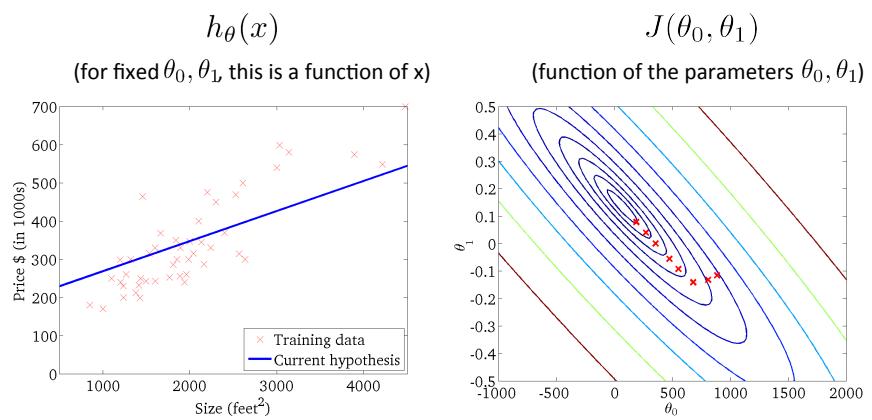
Andrew Ng



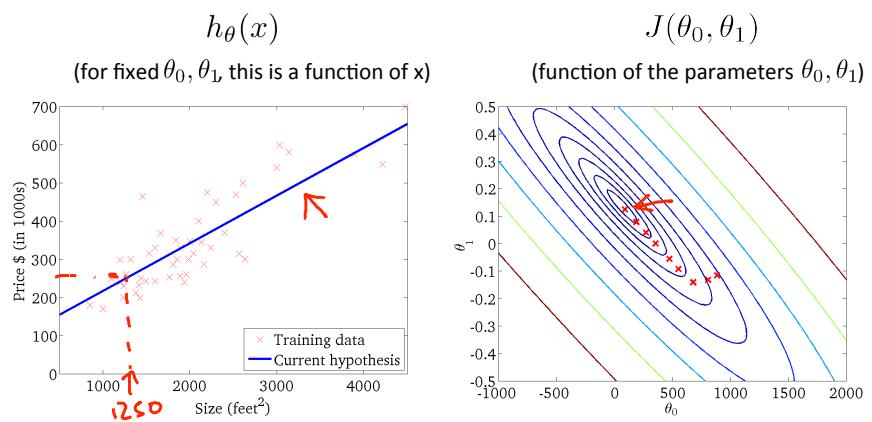
Andrew Ng



Andrew Ng



Andrew Ng



Andrew Ng

“Batch” Gradient Descent

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

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