

# Week 1

---

## Table of Contents

Week 1 .....
Introduction .....
What's Machine Learning? .....
Supervised Learning .....
Unsupervised Learning .....
Model and Cost Function .....
Model Representation .....
Cost Function .....
Cost Function Intuition I .....
Cost function - Intuition II .....
Parameter Learning .....
Gradient Descent .....
Gradient Descent Intuition .....
Gradient Descent For Linear Regression .....
Linear Algebra Review .....
Matrices and Vectors .....
Addition and Scalar Multiplication .....
Matrix Vector Multiplication .....
Matrix Matrix Multiplication .....
Matrix Multiplication Properties .....
Inverse and Transpose .....
FAQ for Week1: .....

## Introduction

---

Introduce the core idea of teaching a computer to learn concepts using data - without being explicitly programmed.

**linear regression with one variable.** Discuss the application of linear regression to housing price prediction, present the notion of a cost function, also introduce the gradient descent method for learning.

The course require the knowledge of **linear algebra** concepts.

# What's Machine Learning?

Andrew NG: Machine Learning is the science of getting computers to learn, without being explicitly programmed.

- 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).

## What's Machine Learning? [Definition]

Two definitions of Machine Learning are offered.

- Arthur Samuel described it as: "**the field of study that gives computers the ability to learn without being explicitly programmed.**" This is an older, informal definition.
- Tom Mitchell provides a more modern definition: "**A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E.**"

Examples:

playing checkers.

E = the experience of playing many games of checkers

T = the task of playing checkers.

P = the probability that the program will win the next game.

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. Explain the E T P according to Tom Mitchell's definition.

E = Watching you label emails as spam or not spam.

T = Classify emails as spam or not spam.

P = The number(or fraction) of emails correctly classified as spam/not spam.

In general, any machine learning problem can be assigned to one of two broad classifications:

### **Supervised learning and Unsupervised learning.**

Machine Learning algorithms

- Supervised Learning
- Unsupervised Learning

Others: Reinforcement Learning, recommender systems.

## **Supervised Learning**

---

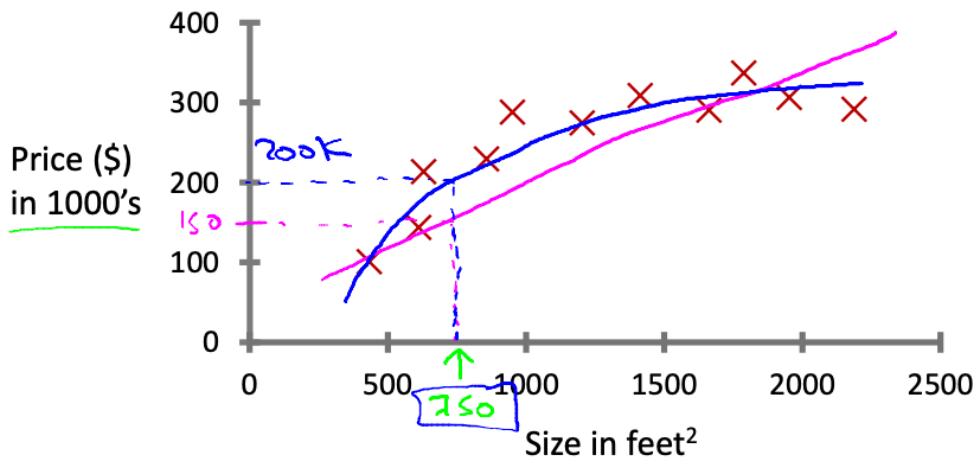
In supervised learning, we are given a data set and already know what our correct output should look like, **having the idea that there is a relationship between the input and the output.**

Supervised learning problems are categorized into "regression" and "classification" problems.

In a regression problem, we are trying to predict results within a **continuous output**, meaning that we are trying to map input variables to some continuous function.

In a classification problem, we are instead trying to predict results in a **discrete output**[离散输出]. In other words, we are trying to map input variables into discrete categories.

## Housing price prediction.

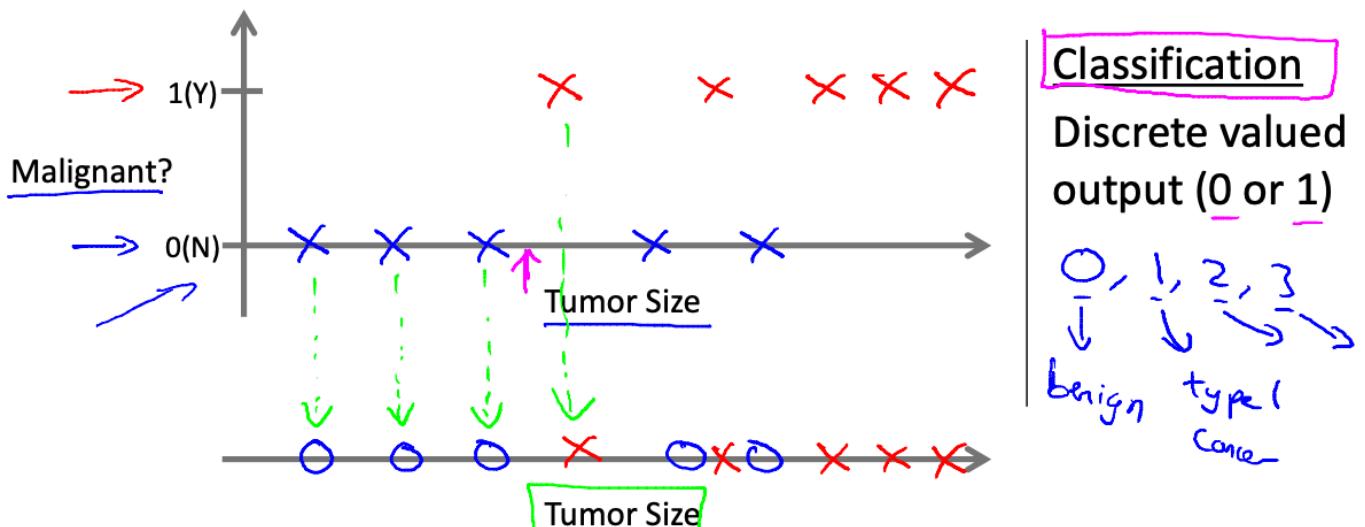


Supervised Learning  
"right answers" given

Regression: Predict continuous valued output (price)

Regression exp

## Breast cancer (malignant, benign)



Classification

### Example 1:

Given data about the size of houses on the real estate market, try to predict their price. Price as a function of size is a continuous output, so this is a regression problem.

We could turn this example into a classification problem by instead making our output about whether the house "sells for more or less than the asking price." Here we are classifying the

houses based on price into two discrete categories.

### **Example 2:**

(a) Regression - Given a picture of a person, we have to predict their age on the basis of the given picture

(b) Classification - Given a patient with a tumor, we have to predict whether the tumor is malignant or benign.

Practice :

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

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. Should you treat these as classification or as regression problems?

Answer:

Treat problem 1 as a regression problem, problem 2 as a classification problem.

Explain:

First one : we have thousand of items [Continuous]

Second one: 0 - not hacked, 1 - hacked.[Hacked or not, Discrete]

## **Unsupervised Learning**

---

Unsupervised learning is a type of machine learning that looks for previously undetected patterns in a data set with no pre-existing labels and with a minimum of human supervision. —  
Wikipedia

Unsupervised learning allows us to approach problems with little or no idea what our results should look like. We can derive structure from data where we don't necessarily know the effect of the variables.

We can derive this structure by clustering the data based on relationships among the variables in the data.

With unsupervised learning there is no feedback based on the prediction results.

### **Example:**

Clustering: Take a collection of 1,000,000 different genes, and find a way to automatically group these genes into groups that are somehow similar or related by different variables, such as lifespan, location, roles, and so on.

Non-clustering: The "Cocktail Party Algorithm", allows you to find structure in a chaotic environment. (i.e. identifying individual voices and music from a mesh of sounds at a [cocktail party](#)).

The Following examples using an unsupervised learning algorithm:

Given a set of news articles found on web, group them into sets of articles about the same stories.

Given a database of customer data, automatically discover market segments and group customers into different market segments.

## **Model and Cost Function**

---

### **Model Representation**

---

#### **Notes:**

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

Size in feet <sup>2</sup> (x)	Price (\$ in 1000's)(y)
→ 2104	460
1416	232
→ 1534	315
852	178
...	...

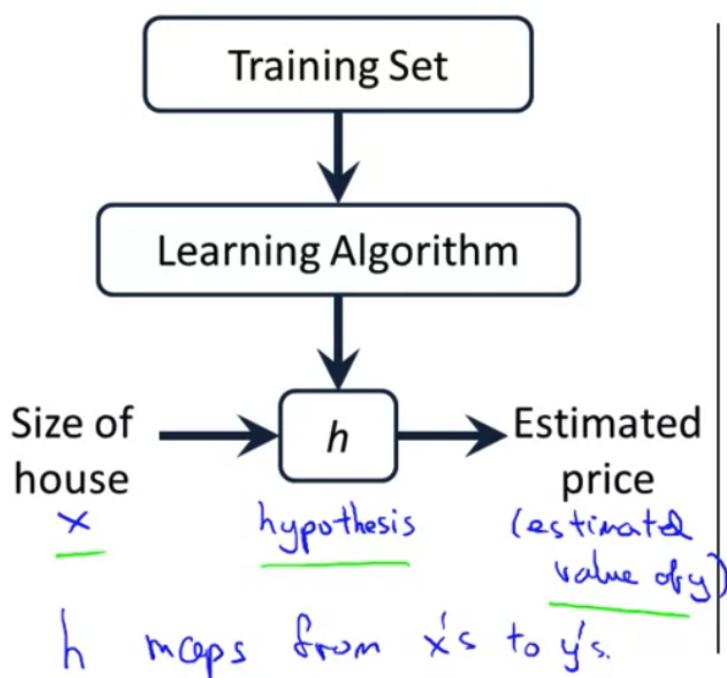
Notation:

- $m$  = Number of training examples
  - $x$ 's = "input" variable / features
  - $y$ 's = "output" variable / "target" variable
- $(x, y)$  - one training example  
 $(x^{(i)}, y^{(i)})$  -  $i^{\text{th}}$  training example

$$\left\{ \begin{array}{l} x^{(1)} = 2104 \\ x^{(2)} = 1416 \\ y^{(1)} = 460 \end{array} \right.$$

Andrew Ng

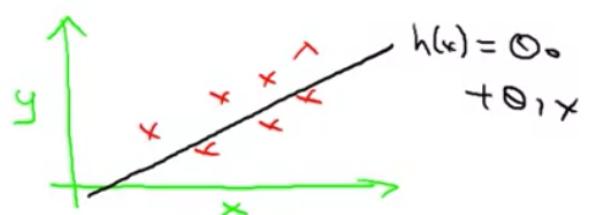
Linear regression model with only one variable [Also called univariate linear regression]



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

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

Shorthand:  $h(x)$



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

**Material:**

To describe the supervised learning problem slightly more formally, our goal is, given a training set, to learn a function  $h : X \rightarrow Y$  so that  $h(x)$  is a "good" predictor for the corresponding value of

y. For historical reasons, this function  $h$  is called a **hypothesis**. Seen pictorially, the process is therefore like this:

When the target variable that we're trying to predict is continuous, such as in our housing example, we call the learning problem a regression problem. When  $y$  can take on only a small number of discrete values (such as if, given the living area, we wanted to predict if a dwelling is a house or an apartment, say), we call it a classification problem.

## Cost Function

Notes:

How to choose the parameters?

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

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

$\theta_i$ 's: Parameters

How to choose  $\theta_i$ 's ?

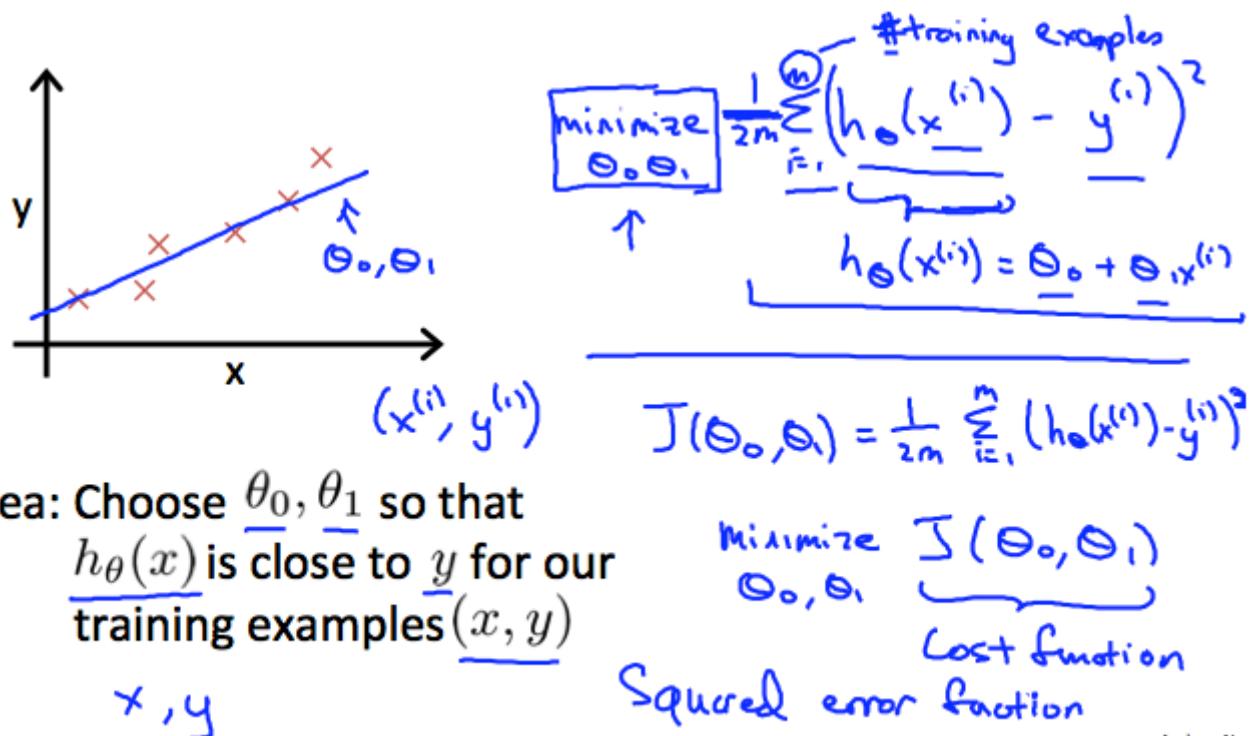
Material:

We can measure the accuracy of our hypothesis function by using a **cost function**. This takes an average difference (actually a fancier version of an average) of all the results of the hypothesis with inputs from  $x$ 's and the actual output  $y$ 's.

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

To break it apart, it is  $\frac{1}{2} \bar{x}$  where  $\bar{x}$  is the mean of the squares of  $h_\theta(x_i) - y_i$ , or the difference between the predicted value and the actual value.

This function is otherwise called the "Squared error function", or "Mean squared error". The mean is halved ( $\frac{1}{2}$ ) as a convenience for the computation of the gradient descent, as the derivative term of the square function will cancel out the  $\frac{1}{2}$  term. The following image summarizes what the cost function does:



MSE: the most common one used for regression problems.

Mean squared error

## Cost Function Intuition I

1. Set the theta(0) to zero.

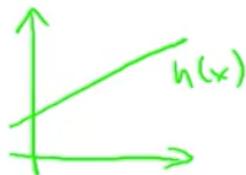
## Simplified

Hypothesis:

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

Parameters:

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



Cost Function:

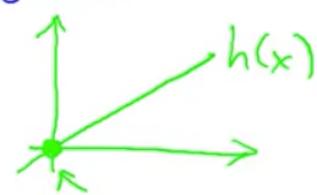
$$\rightarrow J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

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

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

$$\underline{\theta_0 = 0}$$

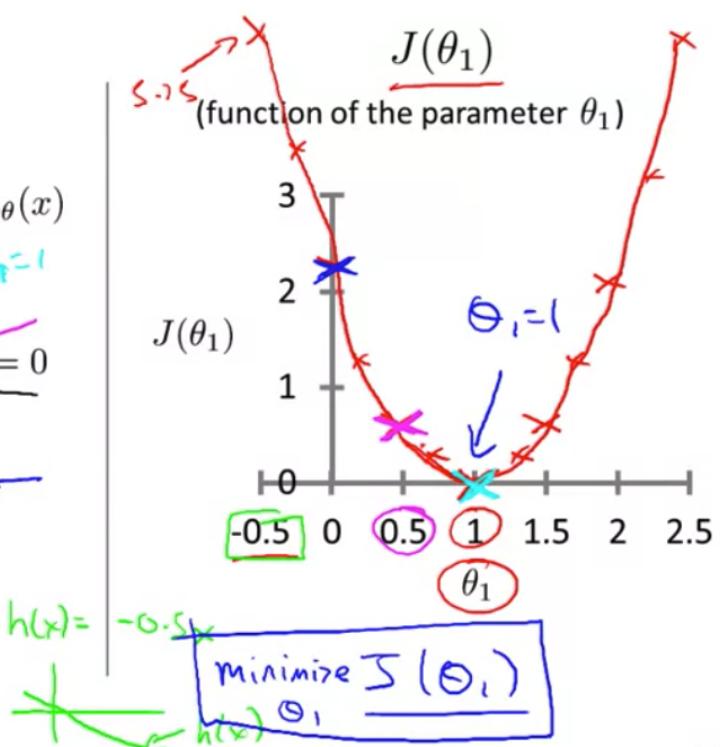
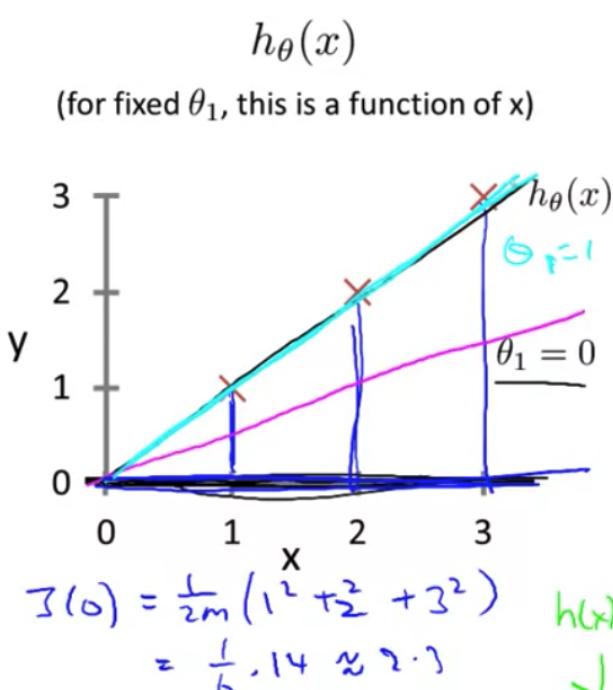
$$\underline{\theta_1}$$



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

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

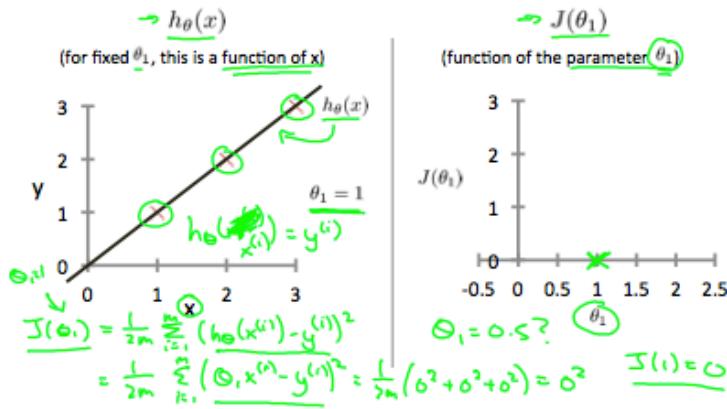
2. The optimization objective for our learning algorithm is we want to choose the value of theta (1). The Minimizes  $J(\theta_1)$ .



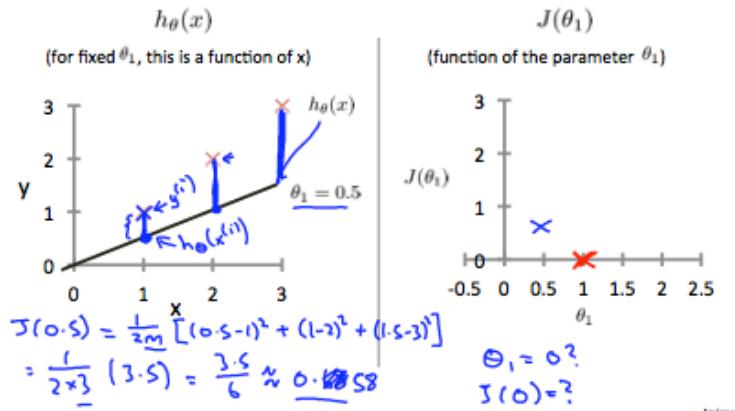
Material:

If we try to think of it in visual terms, our training data set is scattered on the x-y plane. We are trying to make a straight line (defined by  $h_{\theta}(x)$ ) which passes through these scattered data points.

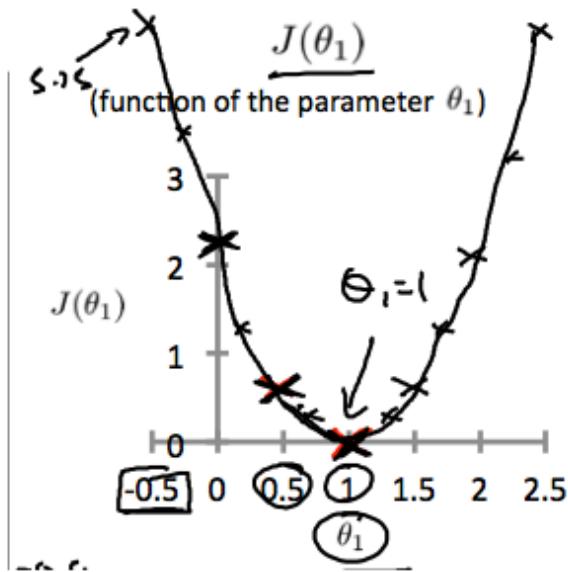
Our objective is to get the best possible line. The best possible line will be such so that the average squared vertical distances of the scattered points from the line will be the least. Ideally, the line should pass through all the points of our training data set. In such a case, the value of  $J(\theta_0, \theta_1)$  will be 0. The following example shows the ideal situation where we have a cost function of 0.



When  $\theta_1 = 1$ , we get a slope of 1 which goes through every single data point in our model. Conversely, when  $\theta_1 = 0.5$ , we see the vertical distance from our fit to the data points increase.



This increases our cost function to 0.58. Plotting several other points yields to the following graph:



Thus as a goal, we should try to minimize the cost function. In this case,  $\theta_1 = 1$  is our global minimum.

## Cost function - Intuition II

1. Keep both two parameters of the cost function.

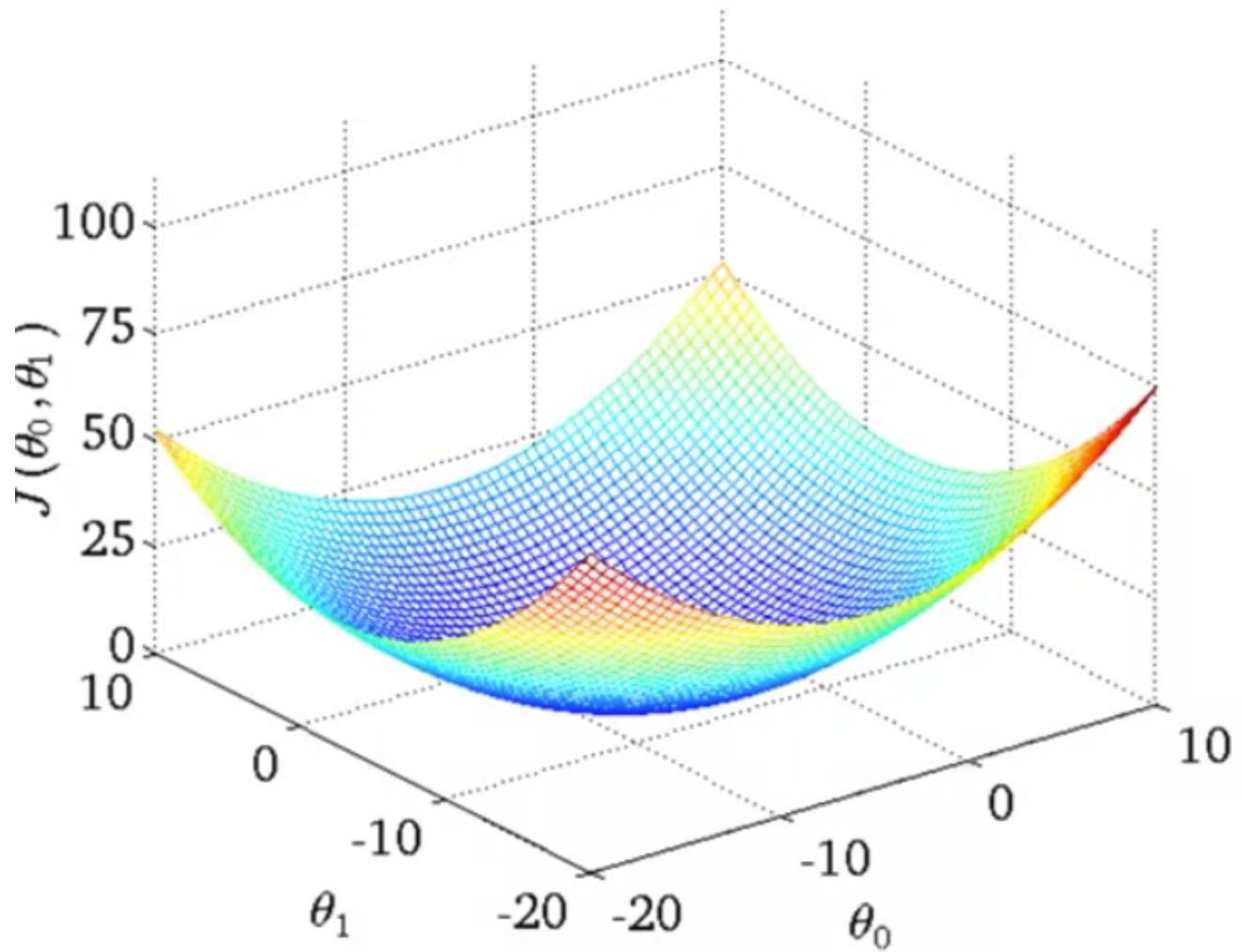
**Hypothesis:** 
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

**Parameters:**  $\theta_0, \theta_1$

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

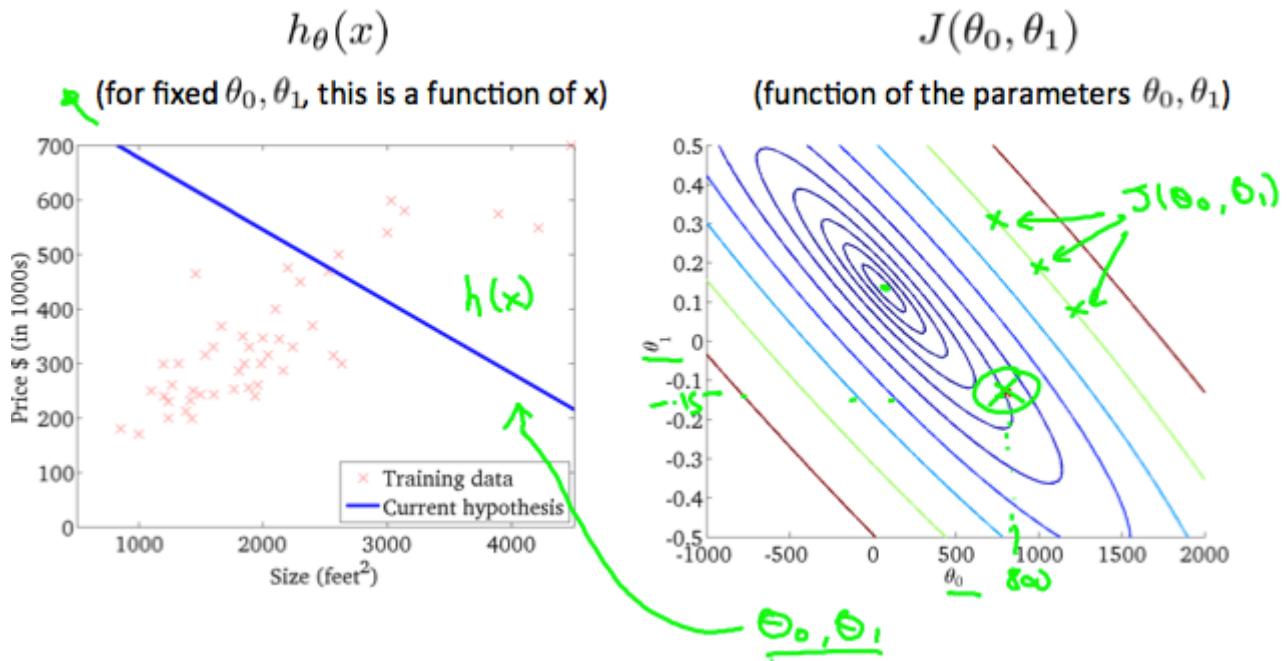
**Goal:** 
$$\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_0, \theta_1)$$

The cost function may look like below:



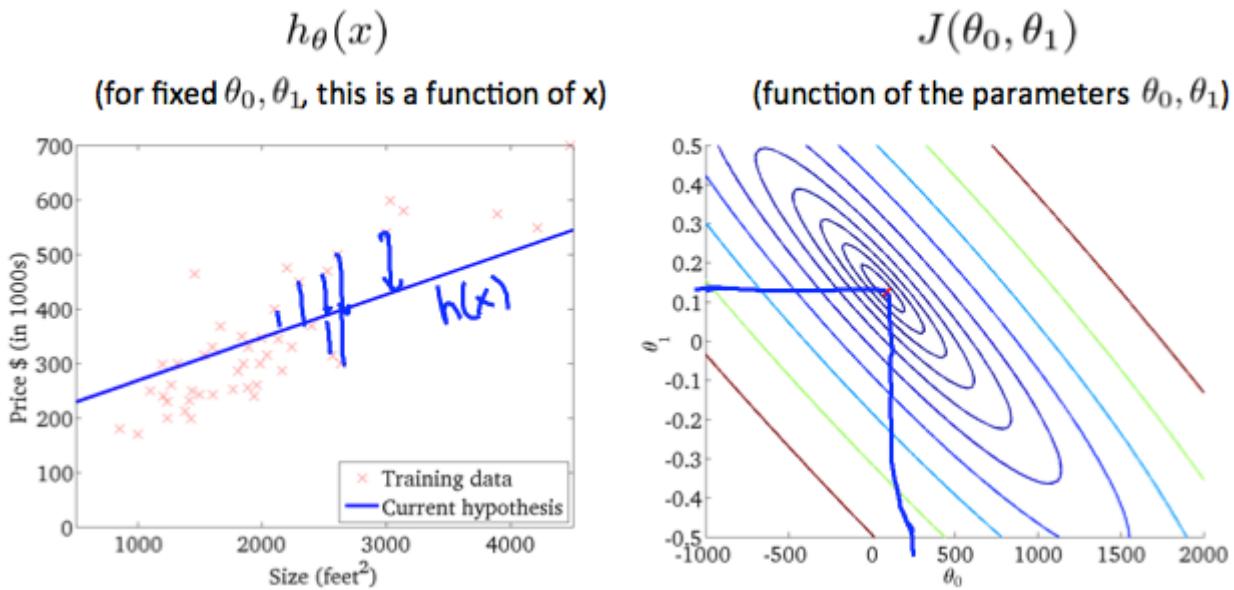
Contour plots / figure

A contour plot is a graph that contains many contour lines. A contour line of a two variable function has constant value at all points of same line. An example of such a graph is the one to the right below:



3 green x have the same cost.

When  $\theta_0=360$  and  $\theta_1 = 0$ , the value of  $J(\theta_0, \theta_1)$  in the contour plot gets closer to the center thus reducing the cost function error. Now Giving our hypothesis function a slightly positive slope result in a better fit of the data.



The graph above minimizes the cost function as much as possible and consequently, the result of  $\theta_1$  and  $\theta_0$  tend to be around 0.12 and 250 respectively. Plotting those values on our graph to the right seems to put our point in the center of the inner most 'circle'.

# Parameter Learning

## Gradient Descent

The algorithm called gradient descent for **minimizing the cost function  $j()$** .

Gradient is a general algorithm.

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

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

### 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

The outline of Gradient descent algorithm:

Understand the concept between Assignment and Truth assertion.

Simultaneous update both parameters.

## 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  $j = 1$ )

Simultaneously update  $\theta_0$  and  $\theta_1$

learning rate

Assignment  $a := b$  ↗  
 $\underline{a := a + 1}$  ↘

Truth assertion  $a = b$  ↗  
 $a = a + 1$  ↘ X

### 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:

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

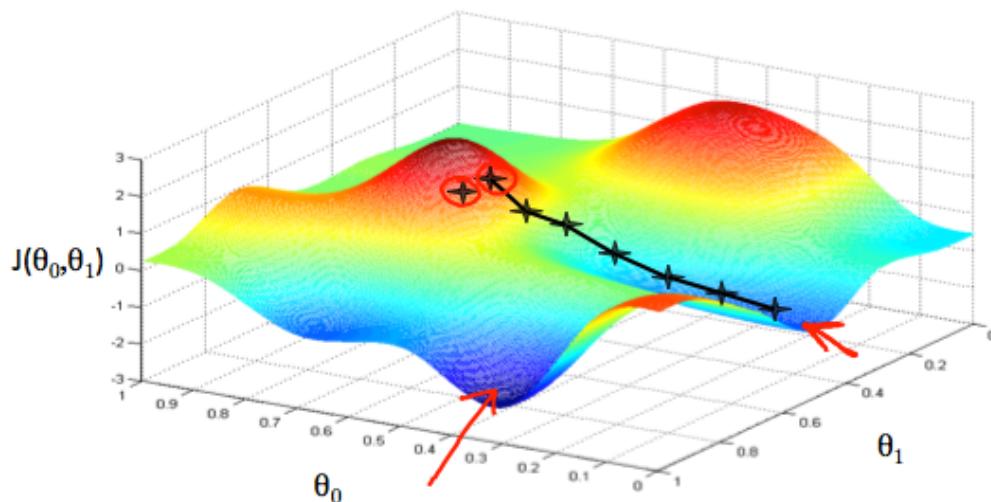
Andrew Ng

## Material:

So we have our hypothesis function and we have a way of measuring how well it fits into the data. Now we need to estimate the parameters in the hypothesis function. That's where gradient descent comes in.

Imagine that we graph our hypothesis function based on its fields  $\theta_0$  and  $\theta_1$  (actually we are graphing the cost function as a function of the parameter estimates). We are not graphing  $x$  and  $y$  itself, but the parameter range of our hypothesis function and the cost resulting from selecting a particular set of parameters.

We put  $\theta_0$  on the x axis and  $\theta_1$  on the y axis, with the cost function on the vertical z axis. The points on our graph will be the result of the cost function using our hypothesis with those specific theta parameters. The graph below depicts such a setup.



We will know that we have succeed when our cost function is at the very bottom of the pits in our graph. when its values is the minimum. The red arrows show the minimum points in the graph.

The way we do is by taking the derivative of our cost function. The slope of the tangent is the derivative at the point and it will give us a direction move towards. We make steps down the cost function in the direction with the steepest descent. The size of each step is determined by the parameter  $\alpha$ , which is called the learning rate.

We will know that we have succeeded when our cost function is at the very bottom of the pits in our graph, when its value is the minimum. The red arrows show the minimum points in the graph.

The way we do this is by taking the derivative (the tangential line to a function) of our cost function. The slope of the tangent is the derivative at that point and it will give us a direction to move towards. We make steps down the cost function in the direction with the steepest descent. The size of each step is determined by the parameter  $\alpha$ , which is called the learning rate.

For example, the distance between each 'star' in the graph above represents a step determined by our parameter  $\alpha$ . A smaller  $\alpha$  would result in a smaller step and a larger  $\alpha$  results in a larger step. The direction which the step is taken is determined by the partial derivative of  $J(\theta_0, \theta_1)$ . Depending on where one starts the graph, one could end up at different points. The image above shows us two different starting points that end up in two different places.

The gradient descent algorithm is:

repeat until convergence:

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

where

$j=0,1$  represents the feature index number.

At each iteration  $j$ , one should simultaneously update the parameters  $\theta_1, \theta_2, \dots, \theta_n$ . Updating a specific parameter prior to calculating another one on the  $j^{(th)}$  iteration would yield to a wrong implementation.

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

```

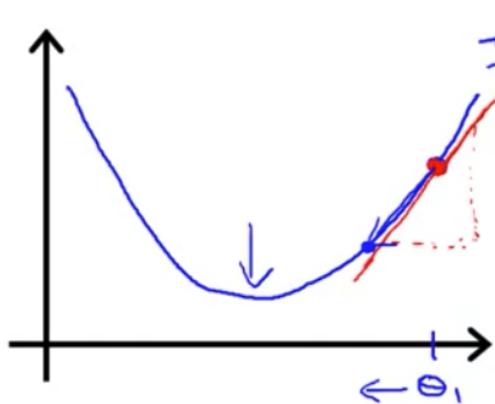
## Gradient Descent Intuition

Apply gradient descent to minimize our squared cost function.

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

don't need to change the a.

Explain the Gradient descent with one variable:

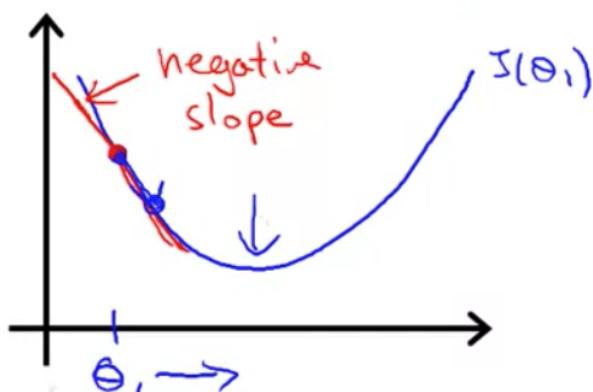


$$J(\theta_1) \quad (\theta_1 \in \mathbb{R})$$

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

$$\frac{\partial}{\partial \theta_1} J(\theta_1) \geq 0$$

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



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

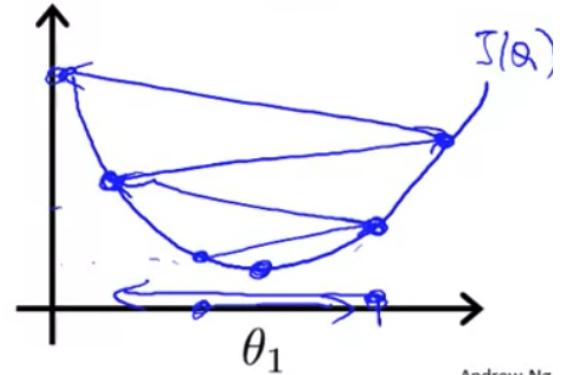
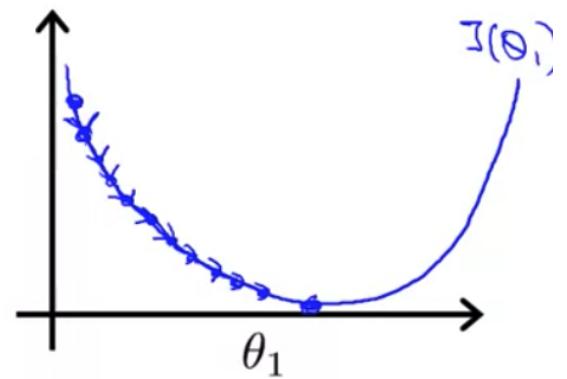
$$\leq 0$$

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

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

If  $\alpha$  is too small, gradient descent can be slow.

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



Andrew Ng

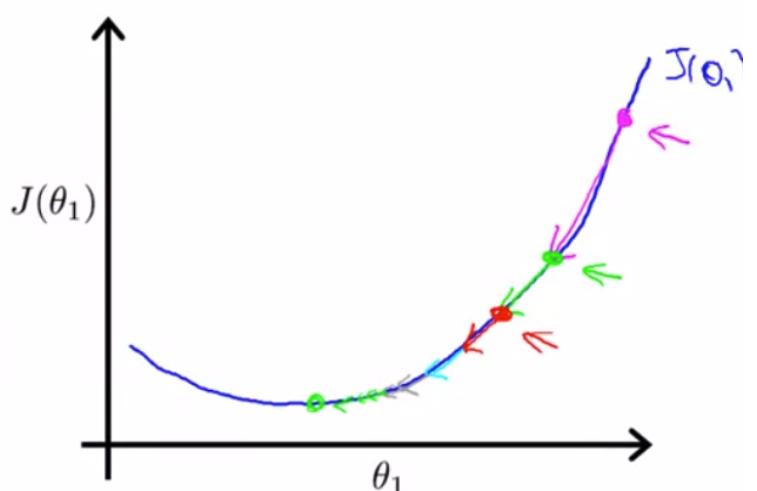
Features of Gradient descent:

When reach the bottom of our convex function, the derivative will always be 0, thus we get the settled  $\theta_1$ .

Gradient descent can converge to a local minimum, even with the learning rate  $\alpha$  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  $\alpha$  over time.



Andrew Ng

## Gradient Descent For Linear Regression

Simplify the function:

When we apply the Gradient descent algorithm with Linear Regression Model, we will derived a new form of the gradient descent equation.

### 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

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

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

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

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) &= \frac{2}{2m} \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\ &= \frac{2}{2m} \cdot \frac{1}{m} \sum_{i=1}^m (\underline{\theta_0 + \theta_1 x^{(i)}} - y^{(i)})^2 \end{aligned}$$

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

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

## Gradient descent algorithm

repeat until convergence {

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

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

}

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

update  
 $\theta_0$  and  $\theta_1$   
simultaneously

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

Convex function equals to Bowl shaped

"Batch" Gradient Descent.

Which of the following are true statements? Select all that apply.

- To make gradient descent converge, we must slowly decrease  $\alpha$  over time.
- Gradient descent is guaranteed to find the global minimum for any function  $J(\theta_0, \theta_1)$ .
- Gradient descent can converge even if  $\alpha$  is kept fixed. (But  $\alpha$  cannot be too large, or else it may fail to converge.)

✓ 正确

- For the specific choice of cost function  $J(\theta_0, \theta_1)$  used in linear regression, there are no local optima (other than the global optimum).

✓ 正确

Material:

A new form of the gradient descent equation:

repeat until convergence: {

$$\begin{aligned}\theta_0 &:= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x_i) - y_i) \\ \theta_1 &:= \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m ((h_\theta(x_i) - y_i)x_i) \\ &\quad }\end{aligned}$$

where  $m$  is the size of the training set,  $\theta_0$  a constant that will be changing simultaneously with  $\theta_1$  and  $x_i, y_i$  are values of the given training set (data).

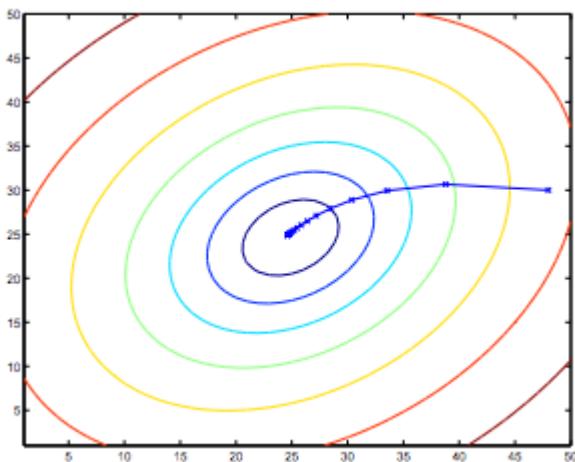
A single example of derivative.

Note that we have separated out the two cases for  $\theta_j$  into separate equations for  $\theta_0$  and  $\theta_1$ ; and that for  $\theta_1$  we are multiplying  $x_i$  at the end due to the derivative. The following is a derivation of  $\frac{\partial}{\partial \theta_j} J(\theta)$  for a single example :

$$\begin{aligned}\frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_\theta(x) - y)^2 \\ &= 2 \cdot \frac{1}{2} (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_\theta(x) - y) \\ &= (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left( \sum_{i=0}^n \theta_i x_i - y \right) \\ &= (h_\theta(x) - y) x_j\end{aligned}$$

So this is simply gradient descent on the original cost function  $j$ . This method looks at every example in the entire training set on every step, which is called **batch gradient descent**. 【批量梯度下降】 Note that, while gradient descent can be susceptible to **local minima** in general. the optimization problem we have posed here for linear regression has only one global, and no other local, optima; thus gradient descent always converges to the global minimum.

Indeed,  $j$  is a **convex quadratic function**. Here is an example of gradient descent as it runs to minimize a quadratic function.



The ellipses shown above are the contours of a quadratic function. Also shown is the trajectory taken by gradient descent, which was initialized at (48,30). The x's in the figure (joined by straight lines) mark the successive values of  $\theta$  that gradient descent went through as it converged to its minimum.

### Rating Test:

1. Consider the problem of predicting how well a student does in her second year of college/university, given how well she did in her first year. 1 分

Specifically, let  $x$  be equal to the number of "A" grades (including A-. A and A+ grades) that a student receives in their first year of college (freshmen year). We would like to predict the value of  $y$ , which we define as the number of "A" grades they get in their second year (sophomore year).

Here each row is one training example. Recall that in linear regression, our hypothesis is  $h_{\theta}(x) = \theta_0 + \theta_1 x$ , and we use  $m$  to denote the number of training examples.

<b>x</b>	<b>y</b>
3	2
1	2
0	1
4	3

For the training set given above (note that this training set may also be referenced in other questions in this quiz), what is the value of  $m$ ? In the box below, please enter your answer (which should be a number between 0 and 10).

2. For this question, assume that we are

1 分

using the training set from Q1. Recall our definition of the

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

What is  $J(0, 1)$ ? In the box below,

please enter your answer (Simplify fractions to decimals when entering answer, and '.' as the decimal delimiter e.g., 1.5).

0.5

3. Suppose we set  $\theta_0 = 0, \theta_1 = 1.5$  in the linear regression hypothesis from Q1. What is  $h_\theta(2)$ ?

1 分

3|

4. Let  $f$  be some function so that

$f(\theta_0, \theta_1)$  outputs a number. For this problem,

$f$  is some arbitrary/unknown smooth function (not necessarily the

cost function of linear regression, so  $f$  may have local optima).

Suppose we use gradient descent to try to minimize  $f(\theta_0, \theta_1)$

as a function of  $\theta_0$  and  $\theta_1$ . Which of the

following statements are true? (Check all that apply.)

- Setting the learning rate  $\alpha$  to be very small is not harmful, and can only speed up the convergence of gradient descent.
- No matter how  $\theta_0$  and  $\theta_1$  are initialized, so long as  $\alpha$  is sufficiently small, we can safely expect gradient descent to converge to the same solution.
- If the first few iterations of gradient descent cause  $f(\theta_0, \theta_1)$  to **increase** rather than decrease, then the most likely cause is that we have set the learning rate  $\alpha$  to too large a value.
- If  $\theta_0$  and  $\theta_1$  are initialized at the global minimum, then one iteration will not change their values.

Choose C and D

5. Suppose that for some linear regression problem (say, predicting housing prices as in the lecture), we have some training set, and for our training set we managed to find some  $\theta_0, \theta_1$  such that  $J(\theta_0, \theta_1) = 0$ . 1分

Which of the statements below must then be true? (Check all that apply.)

- For these values of  $\theta_0$  and  $\theta_1$  that satisfy  $J(\theta_0, \theta_1) = 0$ ,

we have that  $h_\theta(x^{(i)}) = y^{(i)}$  for every training example  $(x^{(i)}, y^{(i)})$

- This is not possible: By the definition of  $J(\theta_0, \theta_1)$ , it is not possible for there to exist

$\theta_0$  and  $\theta_1$  so that  $J(\theta_0, \theta_1) = 0$

- For this to be true, we must have  $\theta_0 = 0$  and  $\theta_1 = 0$

so that  $h_\theta(x) = 0$

- We can perfectly predict the value of  $y$  even for new examples that we have not yet seen.

(e.g., we can perfectly predict prices of even new houses that we have not yet seen.)

5. Suppose that for some linear regression problem (say, predicting housing prices as in the lecture), we have some training set, and for our training set we managed to find some  $\theta_0, \theta_1$  such that  $J(\theta_0, \theta_1) = 0$ . 1分

Which of the statements below must then be true? (Check all that apply.)

- Gradient descent is likely to get stuck at a local minimum and fail to find the global minimum.

- For this to be true, we must have  $y^{(i)} = 0$  for every value of  $i = 1, 2, \dots, m$ .

- For this to be true, we must have  $\theta_0 = 0$  and  $\theta_1 = 0$

so that  $h_\theta(x) = 0$

- Our training set can be fit perfectly by a straight line,

i.e., all of our training examples lie perfectly on some straight line.

## Answers

# Linear Algebra Review

## Matrices and Vectors

Matrix: Rectangular array of numbers.

Matrices are 2-dimensional arrays.

## Matrix Elements (entries of matrix)

$$A = \begin{bmatrix} 1402 & 191 \\ 1371 & 821 \\ 949 & 1437 \\ 147 & 1448 \end{bmatrix}$$

$$A_{11} = 1402$$

$$A_{12} = 191$$

$$A_{32} = 1437$$

$$A_{41} = 147$$

~~$$A_{43} = \text{unc}$$~~

Vector: An  $n \times 1$  matrix.

A vector is a matrix with one column and many rows.

So vectors are a **subset** of matrices.

**Vector:** An  $n \times 1$  matrix.

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix} \quad \begin{array}{l} \uparrow \uparrow \\ n=4 \end{array}$$

$\leftarrow$  4-dimensional vector.

~~$\mathbb{R}^{3 \times 2}$~~

$\mathbb{R}^4$

$y_i = i^{\text{th}}$  element

$$y_1 = 460$$

$$y_2 = 232$$

$$y_3 = 315$$

A, B, C, X

a, b, x, y

1-indexed vs 0-indexed:

$$y[1] = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} \quad \begin{array}{l} \leftarrow \\ \leftarrow \\ \leftarrow \\ \leftarrow \end{array}$$

1-indexed

$$y = \begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ y_3 \end{bmatrix} \quad \begin{array}{l} \leftarrow \\ \leftarrow \\ \leftarrow \\ \leftarrow \end{array}$$

0-indexed

## Notation and terms [Essential]

### Notation and terms:

- $A_{ij}$  refers to the element in the  $i^{\text{th}}$  row and  $j^{\text{th}}$  column of matrix A.
- A vector with ' $n$ ' rows is referred to as an ' $n$ '-dimensional vector.
- $v_i$  refers to the element in the  $i^{\text{th}}$  row of the vector.
- In general, all our vectors and matrices will be 1-indexed. Note that for some programming languages, the arrays are 0-indexed.
- Matrices are usually denoted by uppercase names while vectors are lowercase.
- "Scalar" means that an object is a single value, not a vector or matrix.
- $\mathbb{R}$  refers to the set of scalar real numbers.
- $\mathbb{R}^n$  refers to the set of  $n$ -dimensional vectors of real numbers.

### Matlab commands below:

```

% The ; denotes we are going back to a new row.
A = [1, 2, 3; 4, 5, 6; 7, 8, 9; 10, 11, 12]

% Initialize a vector
v = [1;2;3]

% Get the dimension of the matrix A where m = rows and n = columns
[m,n] = size(A)

% You could also store it this way
dim_A = size(A)

% Get the dimension of the vector v
dim_v = size(v)

% Now let's index into the 2nd row 3rd column of matrix A
A_23 = A(2,3)

```

## Addition and Scalar Multiplication

---

### **Outline:**

Matrix addition and subtraction, and how to multiply a matrix by a number. SO called Scalar Multiplication.

### **Matrix Addition:**

# Matrix Addition

$$\begin{array}{c} \downarrow \quad \downarrow \\ \rightarrow \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} + \begin{bmatrix} 4 & 0.5 \\ 2 & 5 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 5 & 0.5 \\ 4 & 10 \\ 3 & 2 \end{bmatrix} \end{array}$$

$\underbrace{\hspace{10em}}_{3 \times 2 \text{ matrix}}$        $\underbrace{\hspace{10em}}_{3 \times 2}$        $\underbrace{\hspace{10em}}_{3 \times 2}$

$$\begin{array}{c} \rightarrow \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} + \begin{bmatrix} 4 & 0.5 \\ 2 & 5 \end{bmatrix} = \text{error} \\ \rightarrow \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} \quad \begin{bmatrix} 4 & 0.5 \\ 2 & 5 \end{bmatrix} \\ \underbrace{\hspace{10em}}_{3 \times 2} \quad \underbrace{\hspace{10em}}_{2 \times 2} \end{array}$$

Scalar Multiplication:

## Scalar Multiplication

real number

$$3 \times \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} = \boxed{\begin{bmatrix} 3 & 0 \\ 6 & 15 \\ 9 & 3 \end{bmatrix}} = \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} \times 3$$

$\underbrace{\hspace{10em}}_{3 \times 2} \quad \underbrace{\hspace{10em}}_{3 \times 2}$

$$\boxed{\begin{bmatrix} 4 & 0 \\ 6 & 3 \end{bmatrix}} / 4 = \frac{1}{4} \begin{bmatrix} 4 & 0 \\ 6 & 3 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ \frac{3}{2} & \frac{3}{4} \end{bmatrix}$$

Combination of Operands:

## Combination of Operands

$$\begin{aligned} & \quad \boxed{3 \times \begin{bmatrix} 1 \\ 4 \\ 2 \end{bmatrix}} + \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix} - \boxed{\begin{bmatrix} 3 \\ 0 \\ 2 \end{bmatrix} / 3} \\ & \text{Scalar multiplication} \qquad \qquad \qquad \text{Scalar division} \\ = & \begin{bmatrix} 3 \\ 12 \\ 6 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \\ \frac{2}{3} \end{bmatrix} \\ & \qquad \qquad \qquad \text{matrix subtraction /} \\ & \qquad \qquad \qquad \text{vector subtraction} \\ = & \begin{bmatrix} 2 \\ 12 \\ 10 \frac{1}{3} \end{bmatrix} \qquad \qquad \qquad \text{matrix addition /} \\ & \qquad \qquad \qquad \text{vector addition} \end{aligned}$$

Matlab command:

Anderson 31

```

% Initialize matrix A and B
A = [1, 2, 4; 5, 3, 2]
B = [1, 3, 4; 1, 1, 1]

% Initialize constant s
s = 2

% See how element-wise addition works
add_AB = A + B

% See how element-wise subtraction works
sub_AB = A - B

% See how scalar multiplication works
mult_As = A * s

% Divide A by s
div_As = A / s

% What happens if we have a Matrix + scalar?
add_As = A + s

A =
    1    2    4
    5    3    2

B =
    1    3    4
    1    1    1
s = 2

add_AB =
    2    5    8
    6    4    3

sub_AB =
    0   -1    0
    4    2    1

mult_As =
    2    4    8
   10    6    4

div_As =
    0.50000    1.00000    2.00000
    2.50000    1.50000    1.00000

```

```

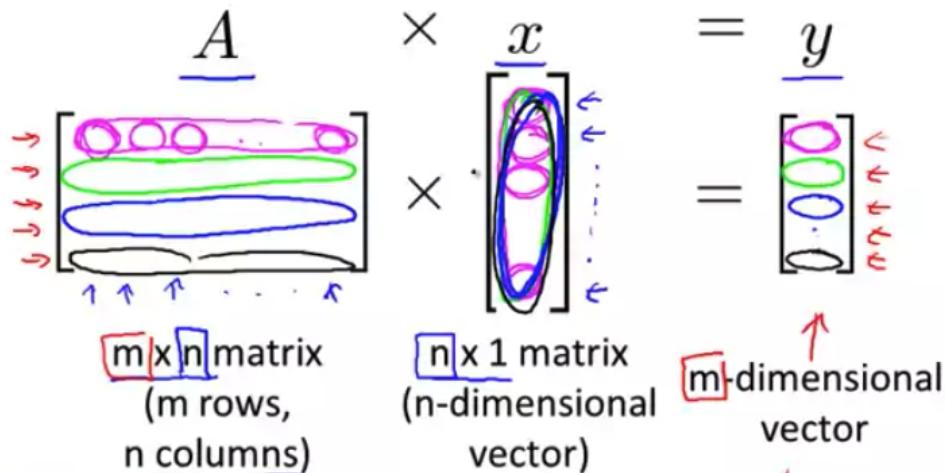
add_As =
3 4 6
7 5 4

```

## Matrix Vector Multiplication

The number of columns of matrix must match the number of rows of vector.

### Details:



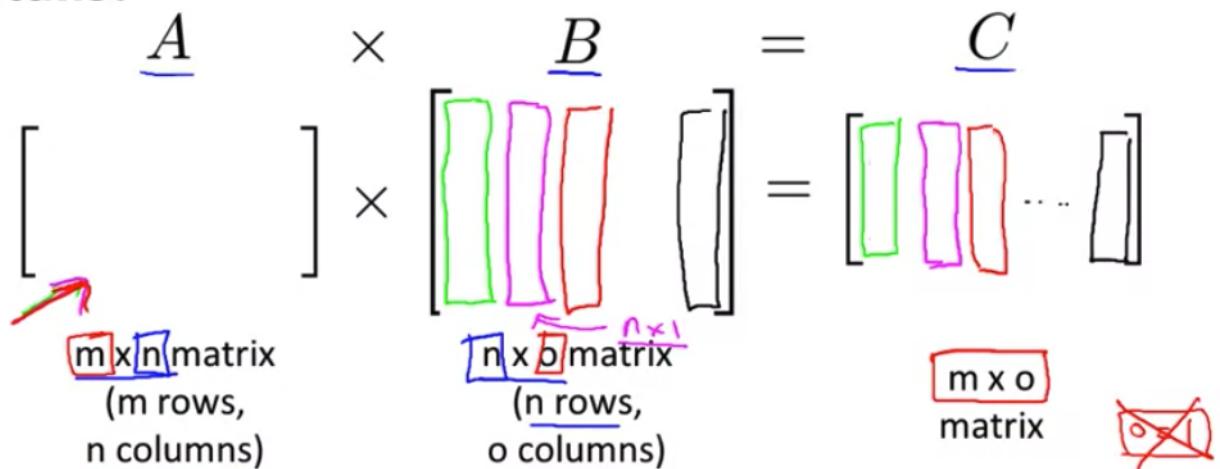
To get  $y_i$ , multiply  $A$ 's  $i^{th}$  row with elements of vector  $x$ , and add them up.

$m \times n$  matrix multiply  $n \times 1$  matrix [n-dimensional vector] =  $m$ -dimensional vector

## Matrix Matrix Multiplication

### Details:

## Details:



The  $i^{th}$  column of the matrix  $C$  is obtained by multiplying  $A$  with the  $i^{th}$  column of  $B$ . (for  $i = 1, 2, \dots, o$ )

Andrew Ng

Apply:

House sizes:

$$\begin{cases} 2104 \\ 1416 \\ 1534 \\ 852 \end{cases}$$

Have 3 competing hypotheses:

1.  $h_{\theta}(x) = -40 + 0.25x$
2.  $h_{\theta}(x) = 200 + 0.1x$
3.  $h_{\theta}(x) = -150 + 0.4x$

Matrix

$$\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix}$$

Matrix

$$\begin{bmatrix} -40 \\ 200 \\ -150 \\ 0.25 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$$\begin{bmatrix} 486 \\ 314 \\ 342 \\ 410 \\ 416 \\ 692 \\ 344 \\ 353 \\ 464 \\ 173 \\ 285 \\ 191 \end{bmatrix}$$

Prediction  
of 1st  
 $h_{\theta}$

Predictions  
of 2nd  
 $h_{\theta}$

Andrew Ng

Matlab command:

```

% Initialize a 3 by 2 matrix
A = [1, 2; 3, 4; 5, 6]

% Initialize a 2 by 1 matrix
B = [1; 2]

% We expect a resulting matrix of (3 by 2)*(2 by 1) = (3 by 1)
mult_AB = A*B

% Make sure you understand why we got that result
A =
    1     2
    3     4
    5     6

B =
    1
    2

mult_AB =
    5
    11
    17

```

## Matrix Multiplication Properties

---

**Not commutative**

$$\underbrace{3 \times 5}_{\text{=}} = 5 \times 3 \quad \text{"Commutative"}$$

Let  $A$  and  $B$  be matrices. Then in general,  
 $A \times B \neq B \times A$ . (not commutative.)

E.g.

$\begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 2 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix}$	$\cancel{\text{A} \times \text{B}}$	$\begin{array}{c} \text{A} \times \text{B} \\ \text{m} \times \text{n} \end{array} \quad \begin{array}{c} \text{B} \times \text{A} \\ \text{n} \times \text{m} \end{array}$
$\begin{bmatrix} 0 & 0 \\ 2 & 0 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 2 & 2 \end{bmatrix}$		$\begin{array}{c} \text{A} \times \text{B} \\ \text{m} \times \text{n} \end{array} \quad \begin{array}{c} \text{B} \times \text{A} \\ \text{n} \times \text{m} \end{array}$

### Associative

$$(A * B) * C = A * (B * C)$$

$$\begin{array}{l} \cancel{3 \times 5 \times 2} \\ 3 \times 10 = 30 = 15 \times 2 \end{array} \quad 3 \times (5+2) = (3 \times 5) + 2 \quad \boxed{\text{"Associative"}}$$

$\begin{array}{c} \uparrow \\ A \times (B \times C) \end{array}$ 
  
 $\begin{array}{c} \leftarrow \\ (\underline{A \times B}) \times C \end{array}$

$$A \times B \times C.$$

$$\text{Let } \underline{D = B \times C}. \text{ Compute } A \times D.$$

$$\text{Let } \underline{E = A \times B}. \text{ Compute } E \times C.$$

$\begin{array}{c} \text{A} \times (\text{B} \times \text{C}) \\ (\text{A} \times \text{B}) \times \text{C} \end{array}$ 
  
 $\Rightarrow$  Some answer.

### Identity Matrix:

The identity matrix of size  $n$  is the  $n \times n$  square matrix with ones on the main diagonal and zeros elsewhere. Details of Identity Matrix:

## Identity Matrix

$1$  is identity.

$$\boxed{1 \times z = z \times 1 = z}$$

↗ for any  $z$

Denoted  $I$  (or  $I_{n \times n}$ ).

Examples of identity matrices:

$$\begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad 1 \times 1$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad 3 \times 3$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad 4 \times 4$$

Informally:

$$\begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix}$$

For any matrix  $A$ ,

$$A \cdot I = I \cdot A = A$$

$\begin{matrix} \nearrow m \times n & \nearrow n \times n & \nearrow m \times m & \nearrow m \times n & \nearrow m \times n \end{matrix}$

$$I_{n \times n}$$

Note:

$$AB \neq BA \text{ in general}$$

$$AI = IA \quad \checkmark$$

Andrew Ng

## Inverse and Transpose

Start by how it relates to real numbers.

Inverse Matrix:

I = "identity."

$$3 \underbrace{\left[ \begin{smallmatrix} 3^{-1} \\ \hline 1 \end{smallmatrix} \right]}_{\frac{1}{3}} = 1$$

$$12 \underbrace{\left( 12^{-1} \right)}_{\frac{1}{12}} = 1$$

Not all numbers have an inverse.

**Matrix inverse:**

If  $A$  is an  $m \times m$  matrix, and if it has an inverse,  $A^{-1}$

$$\rightarrow A(A^{-1}) = A^{-1}A = I.$$

E.g.  $\underbrace{\begin{bmatrix} 3 & 4 \\ 2 & 16 \end{bmatrix}}_A \underbrace{\begin{bmatrix} 0.4 & -0.1 \\ -0.05 & 0.075 \end{bmatrix}}_{A^{-1}} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I_{2 \times 2}$

Matrices that don't have an inverse are "singular" or "degenerate"

Andrew Ni

- Matrices that don't have an inverse are "**singular**" or "**degenerate**"

So how can we get the inverse matrix?

Sometimes you can compute inverses by hand but almost no one does that those day. It turns out there is very good **numerical software** for taking a matrix and computing.

**Matrix Transpose:**

The transpose of a matrix is an operator which flips a matrix over its diagonal.

## Matrix Transpose

Example:

$$\underline{A} = \begin{bmatrix} 1 & 2 & 0 \\ 3 & 5 & 9 \end{bmatrix}_{2 \times 3}$$

$$\underline{B} = \underline{A}^T = \begin{bmatrix} 1 & 3 \\ 2 & 5 \\ 0 & 9 \end{bmatrix}_{3 \times 2}$$

Let  $A$  be an  $\underline{m \times n}$  matrix, and let  $B = A^T$ .

Then  $B$  is an  $\underline{n \times m}$  matrix, and

$$B_{ij} = A_{ji}.$$

$$B_{12} = A_{21} = 2$$

$$B_{32} = 9 \quad A_{23} = 9.$$

Matlab Command

```

% Initialize matrix A
A = [1,2,0;0,5,6;7,0,9]

% Transpose A
A_trans = A'

% Take the inverse of A
A_inv = inv(A)

% What is A^(-1)*A?
A_invA = inv(A)*A

A =
    1      2      0
    0      5      6
    7      0      9

A_trans =
    1      0      7
    2      5      0
    0      6      9

A_inv =
    0.3488   -0.1395    0.0930
    0.3256    0.0698   -0.0465
   -0.2713    0.1085    0.0388

A_invA =
    1.0000   -0.0000    0.0000
    0.0000    1.0000   -0.0000
   -0.0000        0    1.0000

```

## FAQ for Week1:

**Q1: In the cost function, why don't we use absolute value(or mod(), or some other function) instead of the squared error?**

The absolute value has some bad characteristics for minimization.

- The gradient is not continuous because the absolute value function is not differentiable at its minimum point.
- It does not emphasize the correction of large errors.

- The `abs()` function is also not very mathematically efficient.

However, the squared error cost function has some very desirable characteristics:

- The cost can be computed very efficiently.
- Its partial derivative is easily computed.
- Its partial derivative is continuous.

## **Q2: How does the Cocktail Party problem work?**

It isn't clustering.

Clustering would not preserve the time-sequence of the sound samples, so the results would not sound like speech.

The method used is Principal Component Analysis. This method is a mathematical trick that takes two sets of correlated data, and returns two new sets of data that are not correlated. The sequence of the data is retained.

If the data is audio recordings, the result has the effect to the human ear of sounding like two separate audio tracks.