

Machine Learning

Week 1: Introduction, Linear Regression, and Linear Algebra

23-12-2021

Introduction

Machine learning is the science of getting computers to learn, without being explicitly programmed.

Machine learning

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

Examples:

- Database mining
 - o Large datasets from growth of automation/web
 - o Eg: web click data, medical records, biology, engineering
- Applications can't program by hand
 - o Eg: autonomous helicopter, handwriting recognition, most of Natural Language Processing (NLP), computer vision
- Self customizing programs
 - o Eg: amazon, netflix product recommendations
- Understanding human learning (brain, real AI)

What is Machine Learning?

- Arthur Samuel (1959): 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 .
 - o Eg: playing chess
 - E - the experience of playing many games of chess
 - T - the task of playing chess
 - P - the probability that the program will win the next game

Machine learning algorithms:

- Supervised learning
- Unsupervised learning

Others: reinforcement learning, recommender systems

Also talk about: practical advice for applying learning algorithm

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.

Regression: predict continuous valued output

Classification: discrete valued output

Unsupervised Learning

Cocktail party problem algorithm: allows to find structure in a chotic environment

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

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.

There is no feedback based on the prediction results.

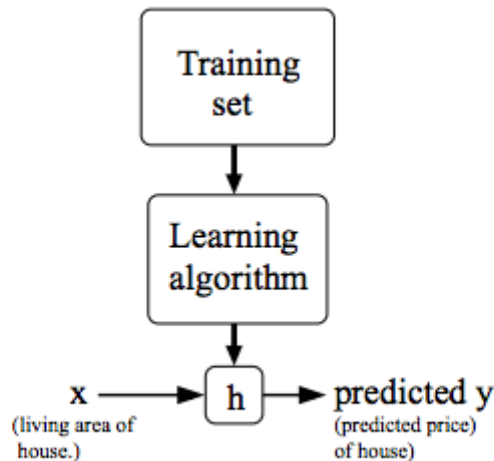
Clustering

Non-clustering

24-12-2021

Linear Regression with One Variable

Model Representation



How do we represent h (hypothesis)?: $h(x)$ or $h_{\theta}(x) = \theta_0 + \theta_1 x$

Linear regression with one variable or Univariant linear regression model is a straight line best-fit for a two variable dataset.

A pair $(x^{(i)}, y^{(i)})$ is called a training example, and the dataset that we'll be using to learn – a list of m training examples $(x^{(i)}, y^{(i)})$; $i = 1, \dots, m$ – is called a training set.

Cost Function

We can measure the accuracy of our hypothesis function by using a cost function. This function is also called Squared error function or Mean squared error.

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

For different θ_0 and θ_1 the hypothesis function changes.

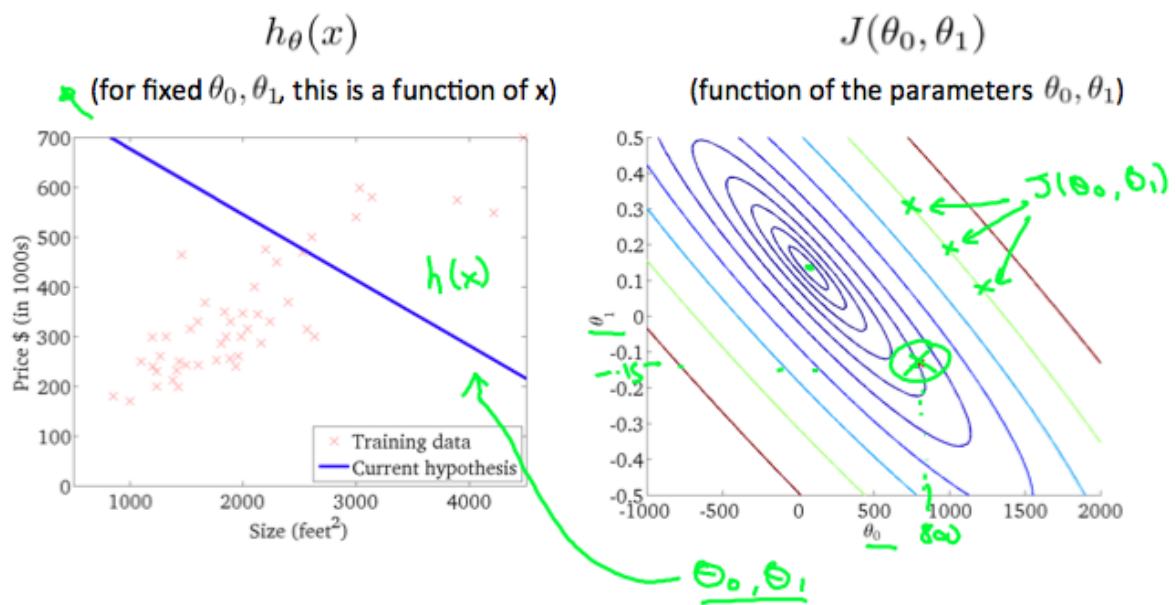
In linear regression we want to find values of θ_0 , θ_1 such that the straight line formed is the best fit for the data set.

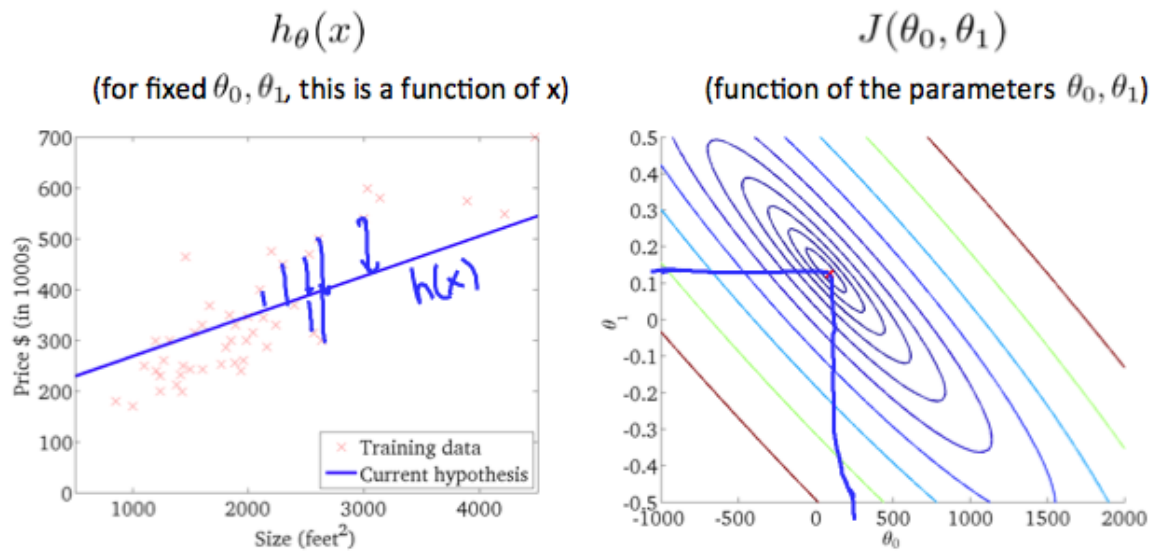
We would like to minimize the difference between output of the hypothesis and the actual value observed or got.

Minimize error over $\theta_0, \theta_1 = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2$

	<u>Simplified</u>
Hypothesis: <u>$h_{\theta}(x) = \theta_0 + \theta_1 x$</u>	$h_{\theta}(x) = \theta_1 x$ $\theta_0 = 0$
Parameters: <u>θ_0, θ_1</u>	<u>θ_1</u>
Cost Function: $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$	$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$
Goal: minimize $J(\theta_0, \theta_1)$ θ_0, θ_1	minimize $J(\theta_1)$ θ_1

Contour plot: is a graph that contains many contour lines. A contour line of a two variable function has a constant value at all points of the same line.



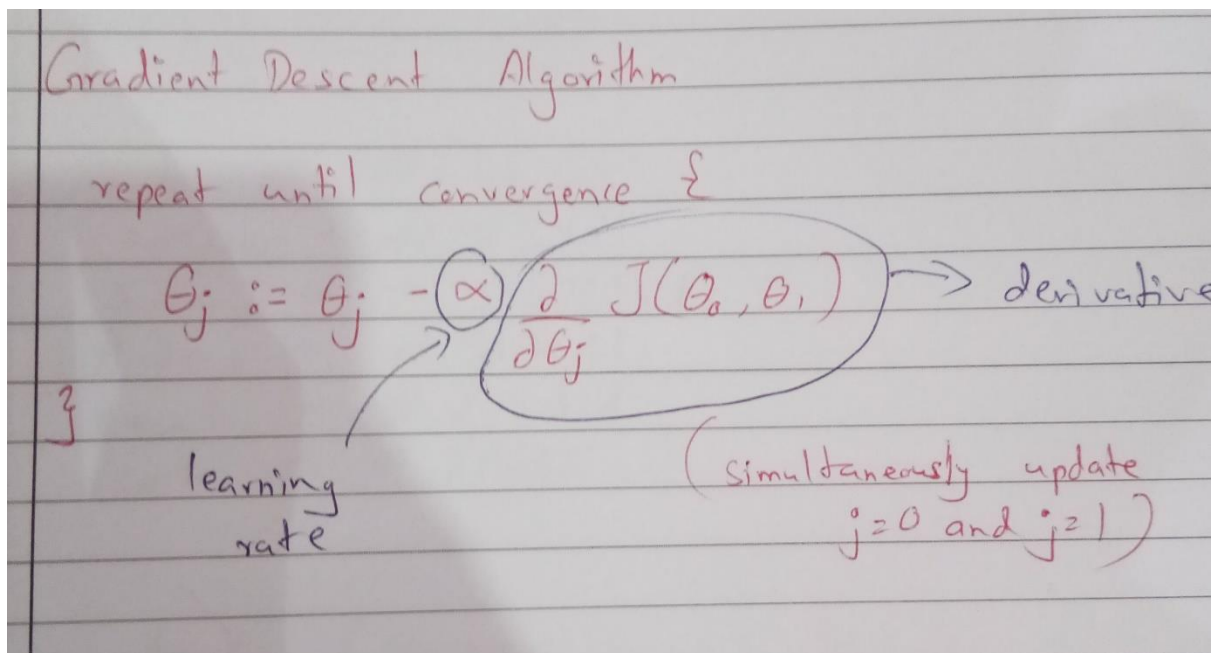


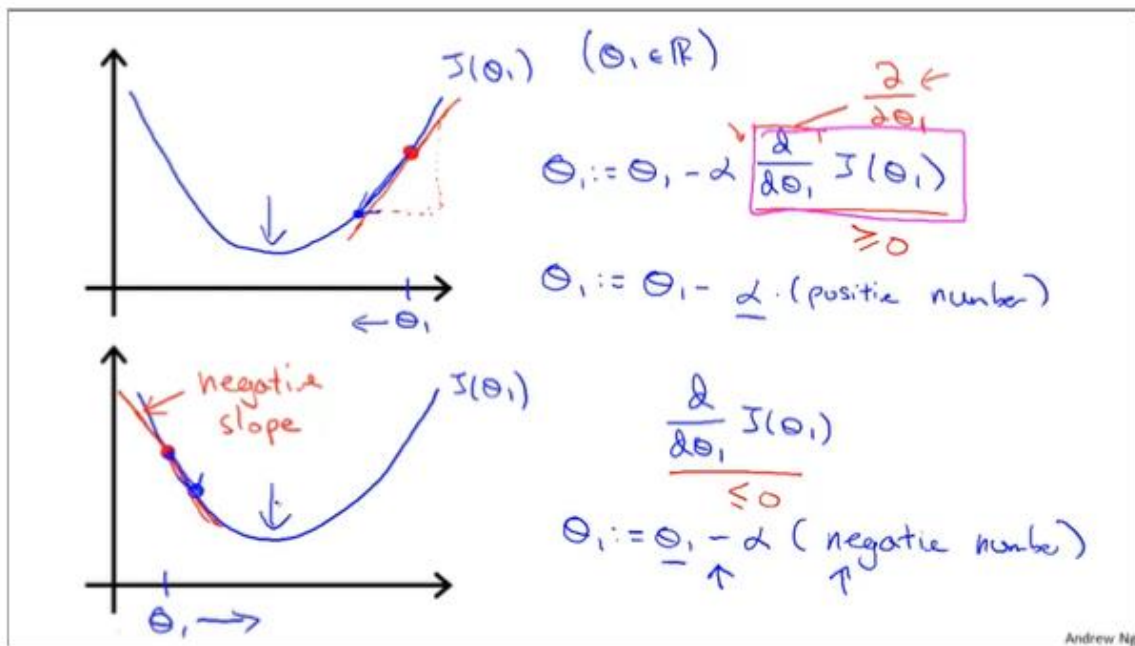
Gradient Descent

For minimizing cost function.

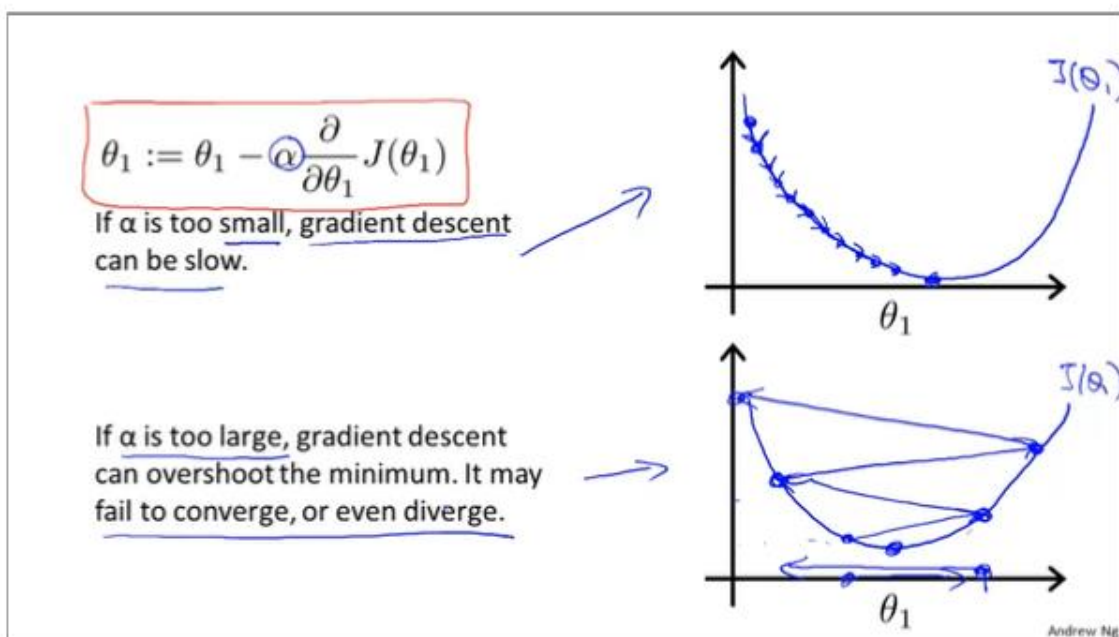
17-01-2022

Gradient Descent Intuition



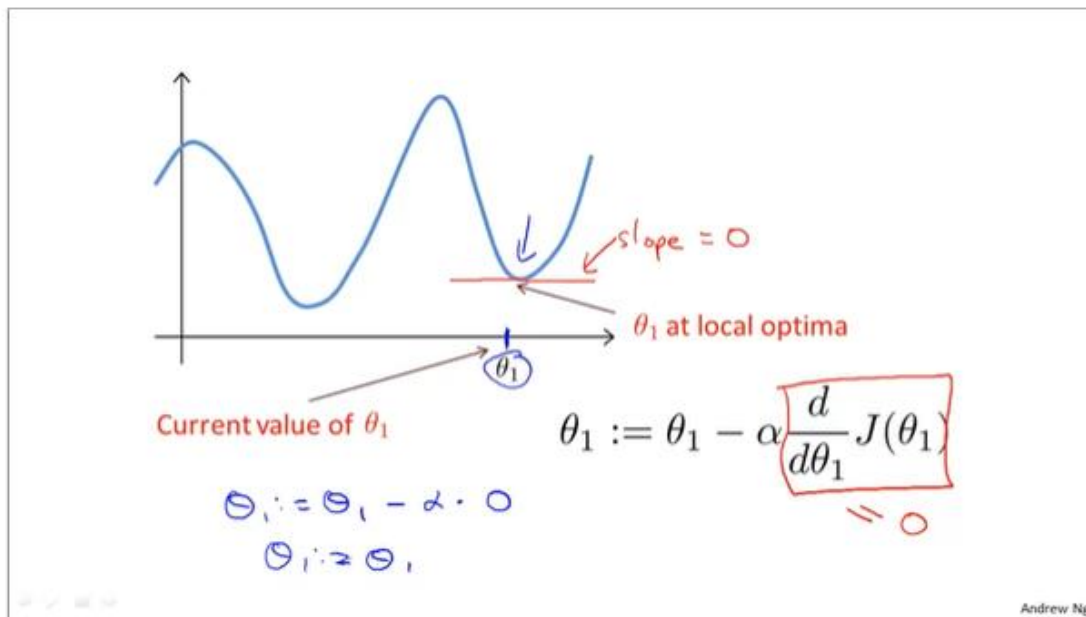


So we can see that the derivative term is shifting our $J(\theta_1)$ term closer to the minimum.



So we see that alpha/learning rate defines how fast we move towards the minimum which if too large can cause completely missing the minimum and shifting forward or backward from it or if too small can be very slow.

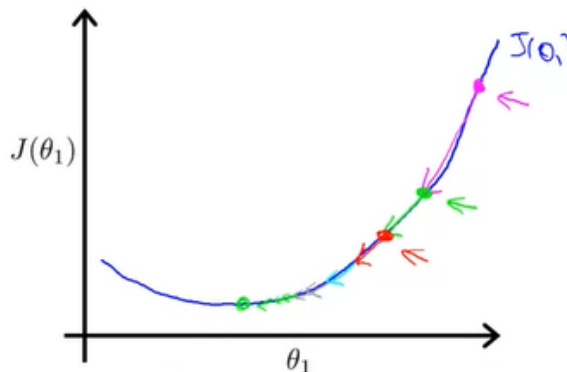
Since derivative for a local minimum will be zero so if we are at a local optima then the value of θ_1 remains unchanged.



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.



So even if we fix the learning rate the derivative will gradually tend to zero and thus we will approach the local minimum.

Gradient Descent for Linear Regression

Gradient descent algorithm	Linear Regression Model
$\text{repeat until convergence } \left\{ \begin{array}{l} \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \\ \text{(for } j = 1 \text{ and } j = 0) \end{array} \right. \}$	$h_{\theta}(x) = \theta_0 + \theta_1 x$ $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$ $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

Andrew Ng

Finding the derivative:

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) &= \frac{\partial}{\partial \theta_j} \left(\frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \right) \\ &= \frac{\partial}{\partial \theta_j} \left(\frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2 \right) \end{aligned}$$
$$\begin{aligned} \theta_0, j = 0 : \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) &= \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \\ \theta_1, j = 1 : \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) &= \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \end{aligned}$$

Andrew Ng

So our gradient descent algorithm converts to:

Gradient descent algorithm

repeat until convergence {

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

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

}

update
 θ_0 and θ_1
simultaneously

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

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

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

Linear Algebra Review

Matrices and Vectors

Matrix: Rectangular array of numbers

Dimension: number of rows * number of columns

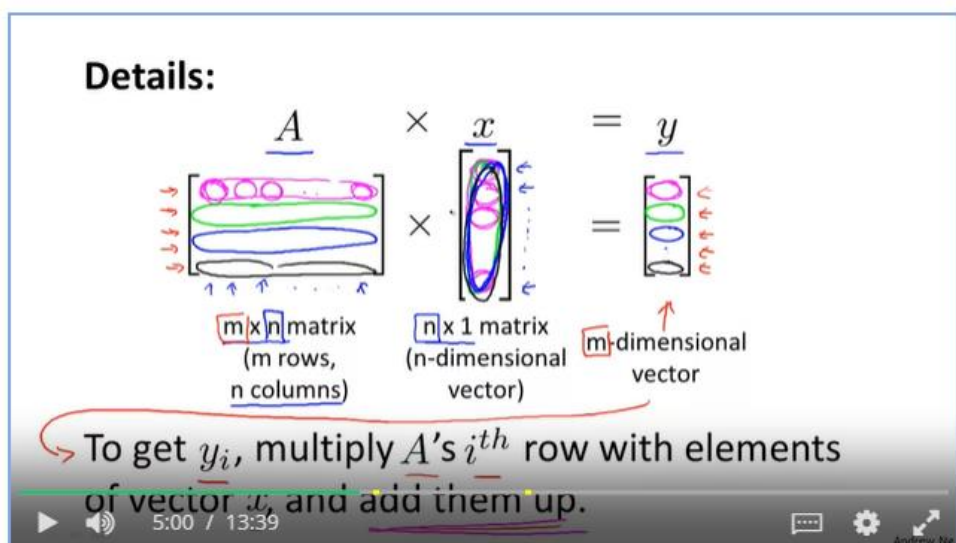
For a matrix A , A_{ij} = "i,j th entry" in the i^{th} row and j^{th} column

Vector: An $n \times 1$ matrix

Indexing of vector elements can start from both 1 and 0.

In linear algebra 1 indexed notation is more commonly used and in ML applications 0 indexed.

- Matrix addition works by adding the corresponding elements, given that both the matrices are of same order/dimensions.
- Multiplying a matrix by scalar results in a matrix with all the elements multiplied by the scalar. The dimensions of the matrix does not changes.
- Multiplication of matrix by a vector
 - Multiplication of a matrix by a vector will result in a vector.
 - Multiplying a matrix A with vector V of dimensions:
 - A : n rows * m columns
 - V : m rows * 1 column
 - Resultant matrix: n rows * 1 column
 - No of columns of matrix should be equal to no of rows in the vector
 - Multiplication works by multiplying a row to the vector with corresponding elements multiplying and then their sum is the corresponding result for that element position in resultant vector



Example

$$\begin{bmatrix} 1 & 2 & 1 & 5 \\ 0 & 3 & 0 & 4 \\ -1 & -2 & 0 & 0 \end{bmatrix}_{3 \times 4} \times \begin{bmatrix} 1 \\ 3 \\ 2 \\ 1 \end{bmatrix}_{4 \times 1} = \begin{bmatrix} 14 \\ 13 \\ -7 \end{bmatrix}_{3 \times 1}$$

$$\begin{aligned}
 1 \times 1 + 2 \times 3 + 1 \times 2 + 5 \times 1 &= 14 \\
 0 \times 1 + 3 \times 3 + 0 \times 2 + 4 \times 1 &= 13 \\
 -1 \times 1 + (-2) \times 3 + 0 \times 2 + 0 \times 1 &= -7
 \end{aligned}$$

Andrew Ng

House sizes:

$\rightarrow 2104$
 $\rightarrow 1416$
 $\rightarrow 1534$
 $\rightarrow 852$

Matrix

$$\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix}_{4 \times 2}$$

$$h_{\theta}(x) = -40 + 0.25x$$

$h_{\theta}(x)$

2x1 Vector

$$\times \begin{bmatrix} -40 \\ 0.25 \end{bmatrix}$$

4x1 matrix

$$\begin{bmatrix} -40 \times 1 + 0.25 \times 2104 \\ -40 \times 1 + 0.25 \times 1416 \\ \vdots \end{bmatrix}_{4 \times 1}$$

$$\text{prediction} = \text{Data Matrix} \times \text{Parameters}$$

for $i = 1, \dots, 1000$,
prediction(i) = ...

Andrew Ng

Matrix-Matrix Multiplication

Details:

$$\begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}_{m \times n \text{ matrix (m rows, n columns)}} \times \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}_{n \times o \text{ matrix (n rows, o columns)}} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}_{m \times o \text{ matrix}}$$

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

Example

$$\begin{bmatrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 5 \end{bmatrix} \begin{bmatrix} 3 \\ 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 11 & 10 \\ 9 & 14 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 0 \\ 5 \end{bmatrix} = \begin{bmatrix} 11 \\ 9 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 3 \\ 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 10 \\ 14 \end{bmatrix}$$

Andrew Ng

○

House sizes:

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

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

$$\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix} \times \begin{bmatrix} -40 & 200 & -150 \\ 0.25 & 0.1 & 0.4 \end{bmatrix} = \begin{bmatrix} 486 & 410 & 692 \\ 314 & 342 & 416 \\ 344 & 353 & 464 \\ 173 & 285 & 191 \end{bmatrix}$$

Prediction of first h_{θ}

Predictions of 2nd h_{θ}

Andrew Ng

○

- Matrix multiplication properties
 - Non-commutative $A \times B \neq B \times A$
 - Associative $A \times (B \times C) = (A \times B) \times C$
 - Identity Matrix I or $I_{n \times n}$ has 1's along the diagonals and all non-diagonal elements are zero.
 - $A_{m \times n} \cdot I_{n \times n} = I_{m \times m} \cdot A_{m \times n} = A_{m \times n}$
- Inverse and Transpose
 - If A is an $m \times m$ matrix, and if it has an inverse, then $A \cdot A^{-1} = A^{-1} \cdot A = I$
 - Matrices that don't have an inverse are "singular" or "degenerate"
 - If A is an $n \times m$ matrix, then A^T is the transpose matrix of dimension $m \times n$, where $A^T_{ij} = A_{ji}$