

## Week 2 Lecture Notes

### ML: Linear Regression With Multiple Variables

Linear regression with multiple variables is also known as "multiple linear regression".

We will introduce notation for equations where we can have any number of input variables.

$x_j^{(i)}$  = value of feature  $j$  in the  $i^{\text{th}}$  training example

$x^{(i)}$  = The feature column vector of all the feature inputs of the  $i^{\text{th}}$  training example

$m$  = The number of training examples

$n = |x^{(i)}|$ ; (The number of features)

Now define the multi-variable form of the hypothesis function as follows, accommodating these multiple features:

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

In order to develop intuition about this function, we can think about  $\theta_0$  as the basic price of a house,  $\theta_1$  as the price per square meter,  $\theta_2$  as the price per floor, etc.  $x_1$  will be the number of square meters in the house,  $x_2$  the number of floors, etc.

Using the definition of matrix multiplication, our multi-variable multi-variable hypothesis function can be concisely represented as:

$$h_{\theta}(x) = [\theta_0 \quad \theta_1 \quad \dots \quad \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} = \theta^T x$$

This is a vectorization of our hypothesis function for one training example; the learning algorithm is now more.

**Remark:** Note that for convenience reasons in this code Mr. Ng assumes  $x_0^{(i)} = 1$  for  $i \in [1, \dots, m]$

[Note: So that we can do matrix operations with theta and  $X$ , we will set  $x_0^{(i)} = 1$ , for all values of  $i$ . This makes the two vectors 'theta' and  $x^{(i)}$  match each other element-wise (that is, have the same number of elements:  $N+1$ ).]

The training examples are stored in  $X$  row - wise, like such:

$$X = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} \\ x_0^{(2)} & x_1^{(2)} \\ x_0^{(3)} & x_1^{(3)} \end{bmatrix}, \quad \Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$$

We can calculate the hypothesis as a column vector of size  $(m \times 1)$  with:

$$h_{\Theta}(x) = X\Theta$$

For the rest of these notes, and other lecture notes,  $X$  will represent a matrix of training examples  $x^{(i)}$  stored row-wise (or) column-wise.

For the parameter vector  $\Theta$  (of type  $\mathbb{R}^{n+1}$  or in  $\mathbb{R}^{(n+1) \times 1}$ ),

The cost function is:

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

The vectorized version is:

$$J(\Theta) = \frac{1}{2m} (X\Theta - \vec{y})^T (X\Theta - \vec{y})$$

Where  $\vec{y}$  denotes the vector of all  $y$  values.

(2)

## - Gradient Descent for Multiple Variables:

The gradient descent equation itself is generally the same form;  
we just have to repeat it for our 'n' features:

Repeat until convergence: {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_0^{(i)}$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_1^{(i)}$$

$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_2^{(i)}$$

...

}

In other words:

Repeat until convergence: {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \text{ for } j := 0..n$$

}

## Matrix Notation

The gradient descent rule can be expressed as:

$$\theta := \theta - \alpha \nabla J(\theta)$$

Where  $\nabla J(\theta)$  is a column vector of the form:

$$\nabla J(\theta) = \begin{bmatrix} \frac{\partial J(\theta)}{\partial \theta_0} & \frac{\partial J(\theta)}{\partial \theta_1} & \vdots & \frac{\partial J(\theta)}{\partial \theta_n} \end{bmatrix}$$

The  $j$ -th component of the gradient is the summation of the product of this term:

$$\begin{aligned}\frac{\partial J(\theta)}{\partial \theta_j} &= \frac{1}{m} \sum_{i=1}^m \left( h_\theta(x^{(i)}) - y^{(i)} \right) \circ x_j^{(i)} \\ &= \frac{1}{m} \sum_{i=1}^m \cancel{x_j^{(i)}} \cdot \left( h_\theta(x^{(i)}) - y^{(i)} \right)\end{aligned}$$

Sometimes, the summation of the product of these terms can be expressed as the product of two vectors.

Here,  $x_j^{(i)}$ , for  $i = 1, \dots, m$ , represents the  $m$  elements of the  $j$ -th column,  $\vec{x}_j$ , of the training set  $X$ .

The other term  $(h_\theta(x^{(i)}) - y^{(i)})$  is the vector of the residuals between the predicted value  $h_\theta(x^{(i)})$  and the true value  $y^{(i)}$ .

Re-writing  $\frac{\partial J(\theta)}{\partial \theta_j}$ , we have:

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \vec{x}_j^T (\vec{x}_\theta - \vec{y})$$

$$\nabla J(\theta) = \frac{1}{m} X^T (X\theta - \vec{y})$$

Finally, the matrix notation (Notation) of the gradient descent rule is

$$\theta := \theta - \frac{\alpha}{m} X^T (X\theta - \vec{y})$$

We can speed up gradient descent by having each of our input values in roughly the same range. This is because  $\theta$  will descend quickly on small ranges and slowly on large ranges, and is ~~W-N~~ oscillate inefficiently between the optimum when the variables are very ~~unconstrained~~ <sup>designed</sup>.

The way to prevent this is to modify the range of the input values so that they are all roughly the same. Ideally:

$$-1 \leq x_{(i)} \leq 1$$

or

$$-0.5 \leq x_{(i)} \leq 0.5$$

There aren't hard requirements; we are simply trying to speed up things. The goal is to get all input variables into roughly one of these ranges, give or take a few.

This techniques to help with this are feature scaling and mean normalization. Feature scaling involves dividing the input values by the range (i.e. the maximum value minus the minimum value) of the input variable, resulting in a new range of just 1. Mean normalization involves subtracting the average value for an input variable, resulting in a new average value for the input variable of just zero. To implement both of these techniques, adjust your input values as shown in this formula:

$$x_i := \frac{x_i - \bar{x}_i}{s_i}$$

Where  $\bar{x}_i$  is the average of all the values of feature  $(i)$  and  $s_i$  is the range of values (max-min), or  $s_i$  is the standard deviation.

Note that dividing by the range, or dividing by the standard deviation,

Give different results. The zig-zag in their course will happen. The  
programming exercise has standard deviation.

Example:  $x_i$  is housing price with range of 100 to 2000, with  
a mean value of 1000. Then,  $X := \frac{\text{Price} - 1000}{1000}$ .

### - Gradient Descent Tipps:

Debugging gradient descent. Make a plot with number of iterations  
on the x-axis. Now plot the cost function,  $J(\theta)$  after the  
number of iterations of gradient descent. If  $J(\theta)$  ~~is~~ <sup>ever</sup> increasing  
then you probably have to decrease  $\alpha$ .

Automatic Convergence Test. Check convergence if  $J(\theta)$  decreases  
by less than  $E$  in one iteration, where  $E$  is some small value  
such as  $10^{-3}$ .

It has been proven that if learning rate  $\alpha$  is sufficiently small  
then  $J(\theta)$  will decrease on every iteration. Andrew Ng recommends  
decreasing  $\alpha$  by multiple of 3.

④

## - Features and Polynomial Regression:

We can improve our feature and the form of our hypothesis function in a couple different ways.

We can combine multiple features into one. For example, we can combine  $X_1$  and  $X_2$  into a new feature  $X_3$  by taking  $X_1 \cdot X_2$ .

## Polynomial Regression

Our hypothesis function need not be linear (a straight line) if that does not fit the data well.

We can change the behavior or curve of our hypothesis function by making it a quadratic, cubic or square root function (or many other form).

For example, if our hypothesis function is  $h_\theta(x) = \theta_0 + \theta_1 x_1$ , then we can create additional features based on  $x_1$ , to get the quadratic function  $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$  or the cubic function  $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3$

In the cubic version, we have created new features  $X_2$  and  $X_3$  where  $X_2 = x_1^2$  and  $X_3 = x_1^3$ .

To make it a square root function, we could do:  $h_\theta(x) = \theta_0 + \theta_1 \sqrt{x_1}$

In the 2:52 and through 6:22 in the "Features and Polynomial regression" video, the slide that Prof Andrew Ng discusses how doesn't mean you can't come back down it in reference to the hypothesis function that uses the `sqrt()` function (shown by the solid purple line), and the one that uses `size2` (shown with the dotted blue line). The quadratic form of the hypothesis function would have the shape shown with the blue dotted line if  $\theta_2$  was negative.

The important thing to keep in mind is, if you choose your feature this way then A scaling becomes very important.

Eg. If  $x_1$  has range  $1 - 1000$  Then range of  $x_1^2$  becomes  $1 - 1000000$   
and that of  $x_1^3$  becomes  $1 - 1000000\dots$ .

## Normal Equation:

⑤

The "normal Equation" is a method of finding the optimum theta without iteration.

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

There is no need to do feature scaling with the normal equation.

Mathematical proof of the Normal equation requires knowledge of linear ~~regression~~ algebra and is fairly involved so you do not need to worry about the details.

The following is a comparison of gradient descent and the normal equation:

### Gradient Descent

Need to choose alpha

Needs many iterations

$O(n^2)$

Works well when N is large also if n is very large

### Normal Equation

No need to choose alpha

No need to iterate

$O(n^3)$ , need to calculate  $X^T X$

and to calculate inverse.

With the normal equation, computing the inversions has complexity  $O(n^3)$ . So if we have very large number of features, the normal equation will be slow.

In practice, when n exceeds 10,000 it might be a good idea to go from a normal solution to an iterative one.

### Normal Equation Non-invertibility

When implementing the normal equation in a code we want to use the 'pinv' function rather than 'inv'.

$X^T X$  may be non-invertible. The common cases are:

- \* Redundant feature, When two features are ~~redundant~~ closely related (i.e. They are linearly dependent)
- \* Too many features (e.g.  $m \leq n$ ). In this case, delete some features or use "soft Nearest-neighbor" (to be explained in a later lesson).

Solutions to the above problem include deleting a feature that is linearly dependent with another or deleting one or more features when there are too many features.

India: Normal Equation:

Gradient Descent

[Note]

Normal equation: Method to solve  
for  $\theta$  analytically



With given  $x, y$   
find  $\theta$  of the minimum of  
the cost function

Ex.  $m = 4$

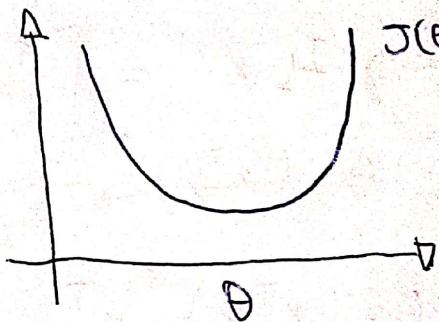
$$\theta = (X^T X)^{-1} X^T y$$

Intuition: If  $J(\theta \in \mathbb{R})$

$$J(\theta) = a\theta^2 + b\theta + c$$

$$\frac{\partial}{\partial \theta} J(\theta) = \dots \stackrel{\text{set}}{=} 0$$

Solve for  $\theta$



$$\theta \in \mathbb{R}^{n+1} \quad J(\theta_0, \theta_1, \dots, \theta_m) =$$

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

Partial derivative  $\frac{\partial}{\partial \theta_j} J(\theta) = \dots = 0$  (for every  $j$ )

Solve for  $\theta_0, \theta_1, \dots, \theta_n$

x.  $m = 4$ .

	$Sig(fit^2)$	$N^2$	$N^o$	Agl	Price
	$x_1$	$x_2$	$x_3$	$x_4$	$y$
0					
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	148

$$X = \begin{bmatrix} 1 & 2104 & 5 & 7 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix} \quad m \times (n+1) \quad m - \text{dimensional vector}$$

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 148 \end{bmatrix}$$

n examples  $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})$  in feature space

$$x^{(\cdot)} = \begin{bmatrix} x_0^{(\cdot)} \\ x_1^{(\cdot)} \\ x_2^{(\cdot)} \\ \vdots \\ x_n^{(\cdot)} \end{bmatrix}$$

$\in \mathbb{R}^{n+1}$

$$X = \begin{bmatrix} - (x^{(1)})^T - \\ - (x^{(2)})^T - \\ - (x^{(3)})^T - \\ \vdots \\ - (x^{(n)})^T - \end{bmatrix} m \times (n+1)$$

E.g. If  $x^{(\cdot)} = \begin{bmatrix} 1 \\ x_1^{(\cdot)} \end{bmatrix}$

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

$$X = \begin{bmatrix} 1 & x_1^{(1)} \\ 1 & x_1^{(2)} \\ \vdots & \vdots \\ 1 & x_1^{(n)} \end{bmatrix} m \times 2$$

$$y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

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

$\rightarrow$  is inverse of Matrix  $x^T x$

$$S\!N A = x^T x$$

$$(x^T x)^{-1} = A^{-1}$$

Goal:  $\min(x^T x) \cdot x^T \cdot y$

$$\theta = (x^T x)^{-1} x^T y \quad \min_{\theta} J(\theta)$$

Feature Scaling

$$-1 \leq x_i \leq 1$$

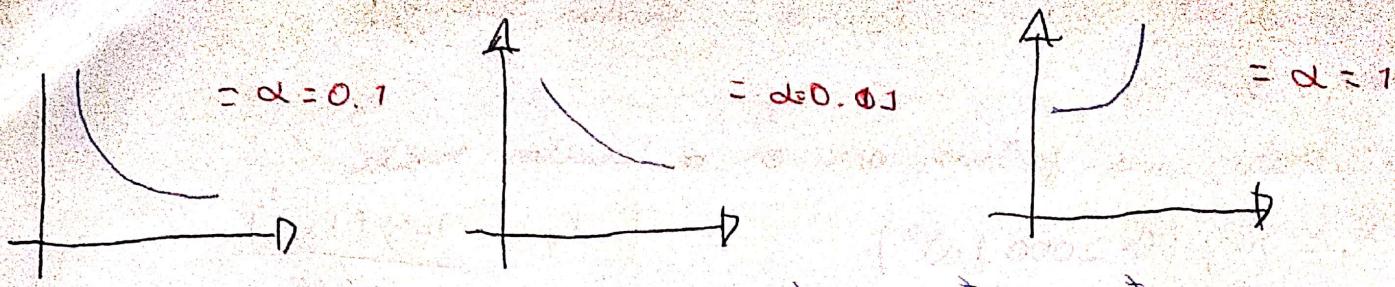
$$0 \leq x_i \leq 100$$

$$0 \leq x_i \leq 10^{-6}$$

Gradient Descent  $\nabla_{\theta} J(\theta) = \text{Normal Equation}$

= Week 2 Lecture Notes

Composition



$\alpha = 0.01$

$\alpha = 0.1$

$\alpha = 1$

In graph C, the cost function is increasing as the learning rate is set too high. Both graphs A and B converge to an optimum of the cost function, but the graph B does so very slowly, as its learning rate is set too little.

Graph A lies between the two.

Summary:

- If  $\alpha$  is too small: slow convergence.
- If  $\alpha \approx 1$ :  $J(\theta)$  may not decrease or may increase.

choose  $\alpha$ , try:

$$\dots, 0.003, 0.005, 0.01, 0.03, 0.1, 0.3, \dots$$

$$\dots, \frac{\partial}{\partial \theta}, \frac{\partial}{\partial \theta}, \frac{\partial}{\partial \theta}, \frac{\partial}{\partial \theta}, \frac{\partial}{\partial \theta}, \dots$$

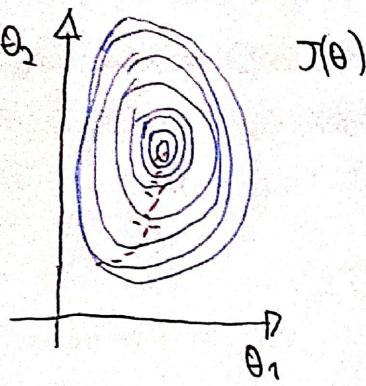
# Videos: Gradient Descent in Practice I - Feature Scaling:

## -Feature Scaling

Idea: Make sure features are on a similar scale.

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

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



$$x_1 = \frac{\text{size (ft}^2)}{2000}$$

$$x_2 = \frac{\text{number of bedrooms}}{5}$$



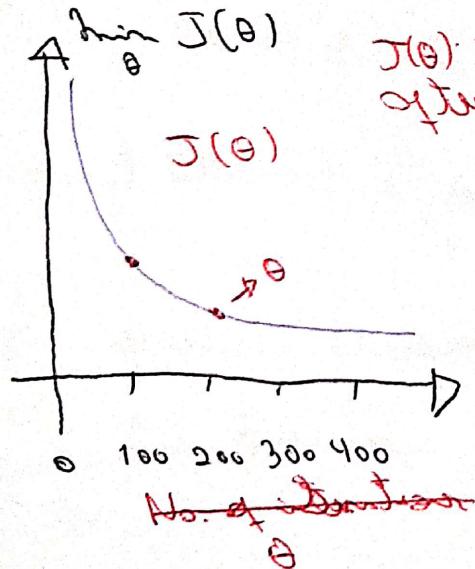
## Feature Scaling

$$0 \leq x_i \leq 1$$

(But many features have approximately a  $-1 < x_i \leq 1$  range.)

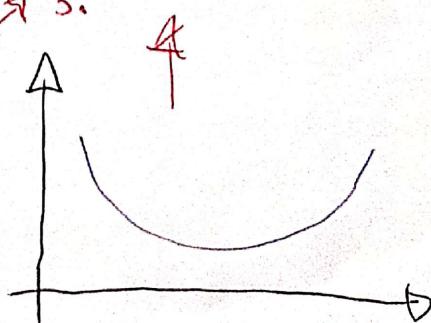
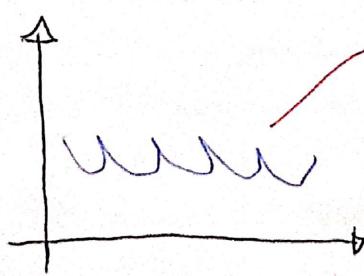
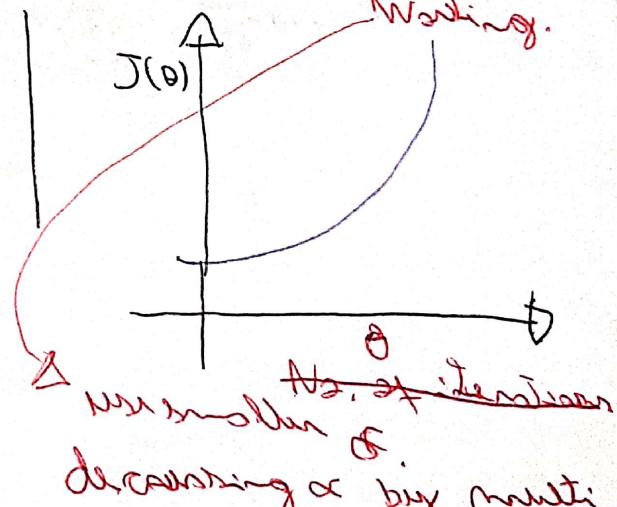
## -D. in Practice II - Learning Rate: Videos

Working when gradient descent is Working correctly.



J(theta) should decrease after every iteration.

G.D. Not Working.



using prior prediction: Linear, Decision Tree and Polynomial Regression

$$h_{\theta}(x) = \theta_0 + \theta_1 \times \underbrace{\text{fringe}}_{x_1} + \theta_2 \times \underbrace{\text{depth}}_{x_2}$$

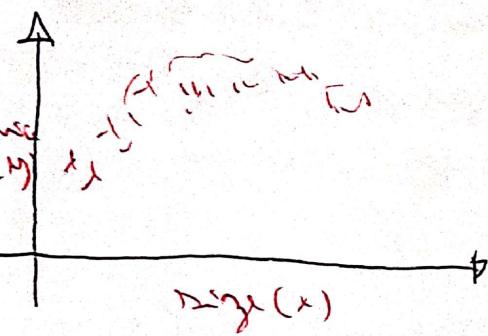
Area

$$X = \text{fringe} * \text{depth}$$

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

As Land Area

Polynomial Regression



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

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

$$= \theta_0 + \theta_1 (\text{fringe}) + \theta_2 (\text{fringe})^2 + \theta_3 (\text{fringe})^3$$

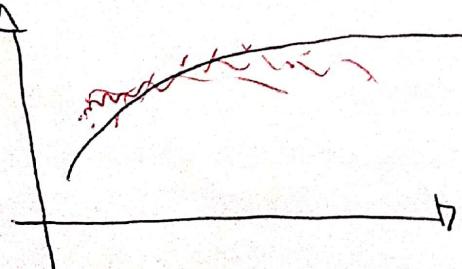
$$x_1 = (\text{fringe})$$

$$x_2 = (\text{fringe})^2$$

$$x_3 = (\text{fringe})^3$$

$$\left| \begin{array}{l} \text{fringe: } 1 - 1000 \\ \text{fringe}^2: 1 - 1000,000 \\ \text{fringe}^3: 1 - 10^9 \end{array} \right.$$

price of houses:



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

Normal equation

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

- What if  $X^T X$  is non-invertible? (Singularity) degenerate
- Define:  $\text{P}_{\text{rank}}(X^T \cdot X) \cdot X^T \cdot y$  [P\_{\text{rank}}] invert

- Redundant features (linearly dependent).

Eg  $x_1 = \text{Area in ft}^2$   $b_m = 3.28 \text{ feet}$   
 $x_2 = \text{Diag in m}^2$

$$x_1 = (3.28)^2 x_2 \quad n = 10$$

- Too many features (e.g.  $m \leq n$ ).  $n = 100$
- Delete some features, or use regularization

(Gauss)

$$x_1 = 99 - 69 = 29$$

$$x_2 = 8836 - 4161 = 4025$$

$$y_5 = 90 - 74 = 22$$

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$$x_1^{(1)} = ?$$

$$x_i := \frac{x_i - \bar{M}_i}{S_i}$$

$x_1$  = mid term score

$x_2 = (\text{mid term score})^2$

$$h_0 = \theta_0 + \theta_1 x_1 + \theta_2 x_2 +$$

$$\frac{\text{Score} - \bar{M}_i}{S_i}$$

49.15

The mean of  $x_1$  is ~~66.15~~

and the standard deviation  $94 - 69 = 25$

$$x_1^{(1)} = \frac{(99 - 49.15)}{25} = 0.34$$

$$\begin{aligned}
 & \left. \begin{array}{c} x_1 \\ x_2 \end{array} \right\} \quad \left. \begin{array}{c} 89 \\ 12 \\ 94 \\ 69 \end{array} \right\} \quad \left. \begin{array}{c} 1921 \\ 5181 \\ 8836 \\ 4161 \end{array} \right\} \\
 & 19.15 \div 4 = 66.155
 \end{aligned}$$

2.

$$\alpha = 0.3$$