

# Multiple Feature

- Single var:

Size (feet <sup>2</sup> ) $\xrightarrow{x}$	Price (\$1000) $y \leftarrow$
2104	460
1416	232
1534	315
$h_{\theta}(x) = \theta_0 + \theta_1 x$	178

Size (feet <sup>2</sup> )	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
$x_1$	$x_2$	$x_3$	$x_4$	$y$
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...	...	...	...	...

$$x_3^{(i)} = 2$$

$$x^{(i)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$$

$n$ : num of features

$x^{(i)}$ : input training examples  $\xrightarrow{n}$  dimension.

$y^{(i)}$ : value of feature  $j$  in  $i$ th training example.

## Question

Size (feet) <sup>2</sup>	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...	...	...	...	...

In the training set above, what is  $x_1^{(4)}$ ?

- The size (in feet<sup>2</sup>) of the 1<sup>st</sup> home in the training set
- The age (in years) of the 1<sup>st</sup> home in the training set
- The size (in feet<sup>2</sup>) of the 4<sup>th</sup> home in the training set
- The age (in years) of the 4<sup>th</sup> home in the training set

✓ Correct

Hypothesis:

$$\text{One Var: } h_{\theta}(x) = \theta_0 + \theta_1 x.$$

null:

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

$$\text{eg: } h_{\theta}(x) = 80 + 0.1 x_1 + 0.01 x_2.$$

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

For convenience, define  $x_0 = 1$  ( $x_0^{(i)} = 1$ )

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n.$$

$$= \theta^T x.$$

$$\underbrace{[\theta_0, \theta_1, \dots, \theta_n]}_{\theta^T}$$

$(n+1) \times 1$  matrix,

Multivariate linear regression

## Multiple Features

**Note:** [7:25 -  $\theta^T$  is a 1 by  $(n+1)$  matrix and not an  $(n+1)$  by 1 matrix]

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

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

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

$x^{(i)}$  = the input (features) of the  $i^{th}$  training example

$m$  = the number of training examples

$n$  = the number of features

The multivariable form of the hypothesis function accommodating these multiple features is as follows:

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

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 multivariable hypothesis function can be concisely represented as:

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

This is a vectorization of our hypothesis function for one training example; see the lessons on vectorization to learn more.

Remark: Note that for convenience reasons in this course we assume  $x_0^{(i)} = 1$  for  $(i \in 1, \dots, m)$ . This allows us to do matrix operations with theta and x. Hence making the two vectors ' $\theta^T$ ' and ' $x^{(i)}$ ' match each other element-wise (that is, have the same number of elements:  $n+1$ ).]

# Gradient Descent for multiple variables

Hypothesis:  $h_{\theta}(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$

parameters:  $\theta_0, \dots, \theta_n$  ( $n+1$  dimension)

Cost function:

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

Gradient descent:

Repeat {  
     $\rightarrow \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n)$   $J(\theta)$   
}

(simultaneously update for every  $j = 0, \dots, n$ )

## Question

When there are  $n$  features, we define the cost function as

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

For linear regression, which of the following are also equivalent and correct definitions of  $J(\theta)$ ?

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

Correct

$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left( \left( \sum_{j=0}^n \theta_j x_j^{(i)} \right) - y^{(i)} \right)^2$  (Inner sum starts at 0)

Correct

$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left( \left( \sum_{j=1}^n \theta_j x_j^{(i)} \right) - y^{(i)} \right)^2$  (Inner sum starts at 1)

$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left( \left( \sum_{j=0}^n \theta_j x_j^{(i)} \right) - \left( \sum_{j=0}^n y_j^{(i)} \right) \right)^2$

## Gradient Descent

Previously (n=1):

Repeat {

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

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

(simultaneously update  $\theta_0, \theta_1$ )

}

New algorithm ( $n \geq 1$ ):

Repeat {

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

(simultaneously update  $\theta_j$  for  $j = 0, \dots, n$ )

}

*gradient descent*

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

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

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

...

## Gradient Descent For Multiple Variables

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

The following image compares gradient descent with one variable to gradient descent with multiple variables:

### Gradient Descent

Previously (n=1):

Repeat {

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

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

}

New algorithm ( $n \geq 1$ ):

Repeat {

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

(simultaneously update  $\theta_j$  for  $j = 0, \dots, n$ )

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

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

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

...

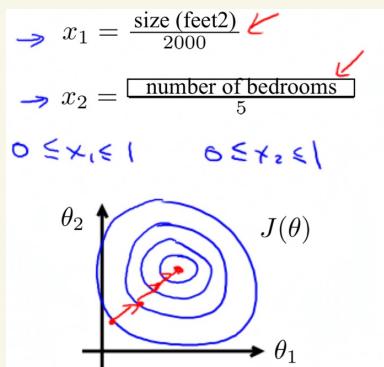
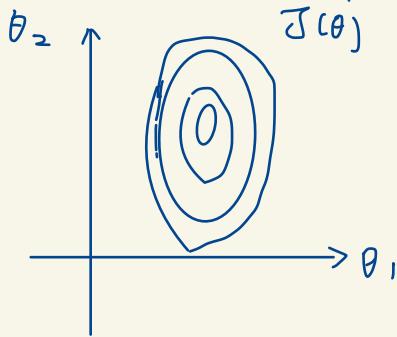
# Feature Scaling

- Feature Scaling

Idea: Make sure features are on a similar scale.

E.g.:  $x_1 = \text{size (0~2000 feet)}$

$x_2 = \text{num of bedrooms (1~5)}$



## Feature Scaling

Get every feature into approximately a  $-1 \leq x_i \leq 1$  range.

$$x_0 = 1$$

$$0 \leq x_1 \leq 3 \quad \checkmark$$

$$-2 \leq x_2 \leq 0.5 \quad \checkmark$$

range too big:  
 $100 \leq x_3 \leq 100$   $\times$

range too small:  
 $-0.0001 \leq x_4 \leq 0.0001$   $\times$

$$-3 \rightarrow 3 \quad \checkmark$$

$$-\frac{1}{2} \rightarrow \frac{1}{2} \quad \checkmark$$

## Mean normalization

Replace  $x_i$  with  $x_i - \mu_i$  to make features have approximately zero mean  
(Do not apply to  $x_0 = 1$ ).

E.g.  $\rightarrow x_1 = \frac{\text{size} - 1000}{2000}$

Average  $\text{size} = 1000$

$$x_2 = \frac{\#\text{bedrooms} - 2}{5}$$

$\frac{1-\Sigma}{5}$  bedrooms

$$\rightarrow [-0.5 \leq x_1 \leq 0.5] \quad [-0.5 \leq x_2 \leq 0.5]$$

$$x_1 \leftarrow \frac{x_1 - \mu_1}{\sigma_1}$$

*avg value of  $x_1$  in training set*

*range (max-min) (or standard deviation)*

$$x_2 \leftarrow \frac{x_2 - \mu_2}{\sigma_2}$$

Ans

### Question

Suppose you are using a learning algorithm to estimate the price of houses in a city. You want one of your features  $x_i$  to capture the age of the house. In your training set, all of your houses have an age between 30 and 50 years, with an average age of 38 years. Which of the following would you use as features, assuming you use feature scaling and mean normalization?

- $x_i = \text{age of house}$
- $x_i = \frac{\text{age of house}}{50}$
- $x_i = \frac{\text{age of house} - 38}{50}$
- $x_i = \frac{\text{age of house} - 38}{20}$

✓ Correct

# Gradient Descent in Practice I - Feature Scaling

**Note:** [6:20 - The average size of a house is 1000 but 100 is accidentally written instead]

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 so will oscillate inefficiently down to the optimum when the variables are very uneven.

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

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

or

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

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

Two 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 from the values for that 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:

$$x_i := \frac{x_i - \mu_i}{s_i}$$

Where  $\mu_i$  is the **average** of all the values for 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 quizzes in this course use range - the programming exercises use standard deviation.

For example, if  $x_i$  represents housing prices with a range of 100 to 2000 and a mean value of 1000, then,  $x_i := \frac{price - 1000}{1900}$ .

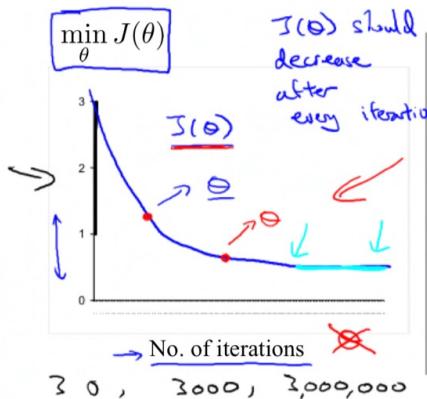
# Learning Rate

## Gradient descent

$$\rightarrow \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

- “Debugging”: How to make sure gradient descent is working correctly.
- How to choose learning rate  $\alpha$

Making sure gradient descent is working correctly.

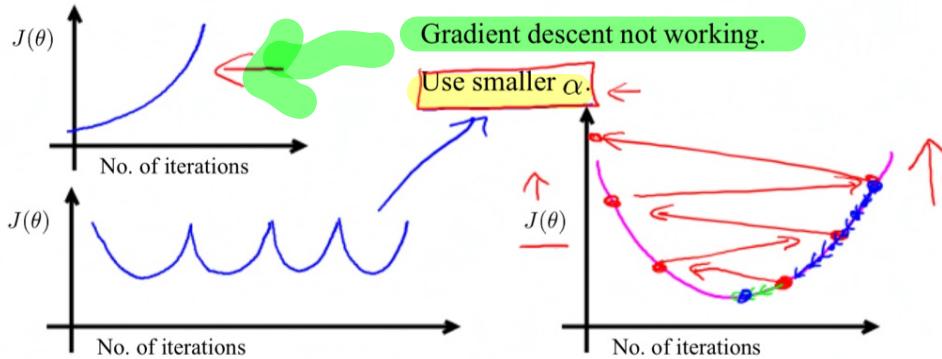


→ Example automatic convergence test:

→ Declare convergence if  $J(\theta)$  decreases by less  $10^{-3}$  than in one iteration.

And

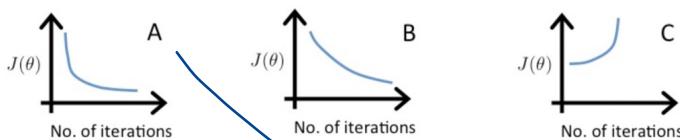
## Making sure gradient descent is working correctly.



- For sufficiently small  $\alpha$ ,  $J(\theta)$  should decrease on every iteration.
- But if  $\alpha$  is too small, gradient descent can be slow to converge.

### Question

Suppose a friend ran gradient descent three times, with  $\alpha = 0.01$ ,  $\alpha = 0.1$ , and  $\alpha = 1$ , and got the following three plots (labeled A, B, and C):



Which plots corresponds to which values of  $\alpha$ ?

- A is  $\alpha = 0.01$ , B is  $\alpha = 0.1$ , C is  $\alpha = 1$ .
- A is  $\alpha = 0.1$ , B is  $\alpha = 0.01$ , C is  $\alpha = 1$ .
- A is  $\alpha = 1$ , B is  $\alpha = 0.01$ , C is  $\alpha = 0.1$ .
- A is  $\alpha = 1$ , B is  $\alpha = 0.1$ , C is  $\alpha = 0.01$ .

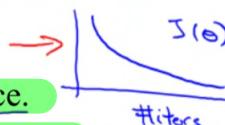
✓ Correct

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

too small  $\alpha \Rightarrow$  too slow  
too big  $\alpha \Rightarrow$  prob won't converge

## Summary:

- If  $\alpha$  is too small: slow convergence.
- If  $\alpha$  is too large:  $J(\theta)$  may not decrease on every iteration; may not converge. (Slow converge also possible)



To choose  $\alpha$ , try

$$\dots, 0.001, \underbrace{0.003}_{\approx 2x}, \underbrace{0.01}_{\approx 2x}, \underbrace{0.03}_{\approx 3x}, \underbrace{0.1}_{\approx 3x}, \underbrace{0.3}_{\approx 3x}, 1, \dots$$

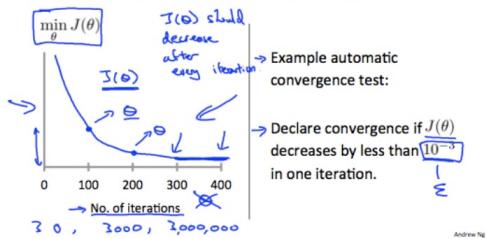
## Gradient Descent in Practice II - Learning Rate

**Note:** [5:20 - the x-axis label in the right graph should be  $\theta$  rather than No. of iterations]

**Debugging gradient descent.** Make a plot with *number of iterations* on the x-axis. Now plot the cost function,  $J(\theta)$  over the number of iterations of gradient descent. If  $J(\theta)$  ever increases, then you probably need to decrease  $\alpha$ .

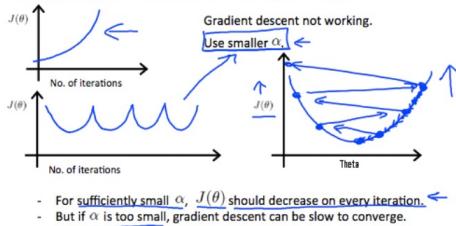
**Automatic convergence test.** Declare convergence if  $J(\theta)$  decreases by less than  $E$  in one iteration, where  $E$  is some small value such as  $10^{-3}$ . However in practice it's difficult to choose this threshold value.

**Making sure gradient descent is working correctly.**



It has been proven that if learning rate  $\alpha$  is sufficiently small, then  $J(\theta)$  will decrease on every iteration.

**Making sure gradient descent is working correctly.**



To summarize:

If  $\alpha$  is too small: slow convergence.

If  $\alpha$  is too large:  $J(\theta)$  may not decrease on every iteration and thus may not converge.

# Features & Polynomial Regression

## Housing prices prediction

$$h_{\theta}(x) = \theta_0 + \theta_1 \times \boxed{\text{frontage}} + \theta_2 \times \boxed{\text{depth}}$$

$x_1$        $x_2$



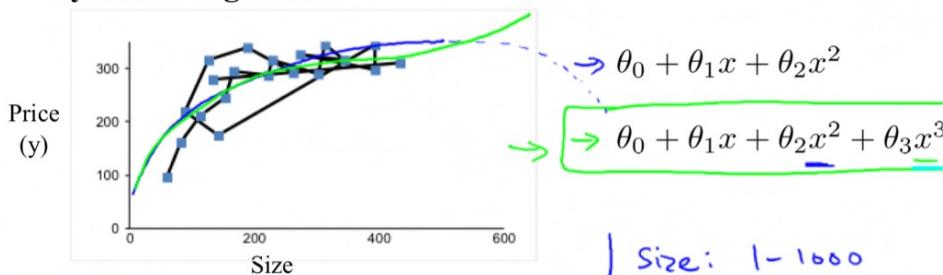
Area

$$\times = \underline{\text{frontage} * \text{depth}}$$

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

↑ land area

## Polynomial regression



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

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

→  $x_1 = (\text{size})$

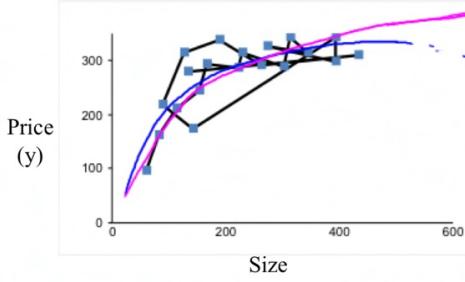
→  $x_2 = (\text{size})^2$

→  $x_3 = (\text{size})^3$

Size:	1 - 1000
Size <sup>2</sup> :	1 - 1000,000
Size <sup>3</sup> :	1 - 10 <sup>9</sup>

And

## Choice of features



K ↴

### Question

Suppose you want to predict a house's price as a function of its size. Your model is

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

Suppose size ranges from 1 to 1000 (feet<sup>2</sup>). You will implement this by fitting a model

$$h_{\theta}(x) = \theta_0 + \theta_1x_1 + \theta_2x_2.$$

Finally, suppose you want to use feature scaling (without mean normalization).

Which of the following choices for  $x_1$  and  $x_2$  should you use? (Note:  $\sqrt{1000} \approx 32$ .)

- $x_1 = \text{size}, x_2 = 32\sqrt{(\text{size})}$
- $x_1 = 32(\text{size}), x_2 = \sqrt{(\text{size})}$
- $x_1 = \frac{\text{size}}{1000}, x_2 = \frac{\sqrt{(\text{size})}}{32}$
- $x_1 = \frac{\text{size}}{32}, x_2 = \sqrt{(\text{size})}$ .

To feature scaling.

$$x_1 = \frac{\text{size}}{1000}$$

$$x_2 = \sqrt{\frac{\text{size}}{1000}}$$

$$= \frac{\sqrt{\text{size}}}{32}$$

✓ Correct

# Features and Polynomial Regression

We can improve our features 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 any 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  $x_3 = x_1^3$ .

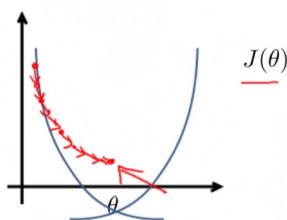
To make it a square root function, we could do:  $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 \sqrt{x_1}$

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

e.g. if  $x_1$  has range 1 - 1000 then range of  $x_1^2$  becomes 1 - 1000000 and that of  $x_1^3$  becomes 1 - 1000000000

# Normal Equation

Gradient Descent

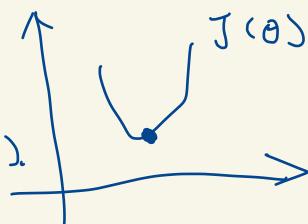


Normal equation: Method to solve for  $\theta$  analytically.

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

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

Solve for  $\theta$ . get  $\frac{\partial}{\partial \theta} J(\theta)$ .



$$\theta \in \mathbb{R}^{n+1} \quad J(\theta_0, \theta_1, \dots, \theta_m) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \dots \stackrel{\text{set}}{=} 0 \quad (\text{for every } j)$$

Solve for  $\underline{\theta_0, \theta_1, \dots, \theta_n}$

Examples:  $m = 4$ .

$x_0$	Size (feet <sup>2</sup> )	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
$x_1$	$x_2$	$x_3$	$x_4$		$y$
2104	5	1	45	460	
1416	3	2	40	232	
1534	3	2	30	315	
852	2	1	36	178	

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$$

$M \times (n+1)$

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

$m$  examples  $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$ ;  $n$  features.

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$X = \begin{bmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix}$$

training example  $x_i$

$$\text{eg.: if } x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & x_1^{(1)} \\ 1 & x_1^{(2)} \\ \vdots & \vdots \\ 1 & x_1^{(m)} \end{bmatrix} \quad \vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

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

$M \times 2$ .

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

$(X^T X)^{-1}$  is inverse of  $X^T X$

Octave:  $\text{pinv}(X' * X) * X' * y \quad (X' = X^T)$

$$\begin{array}{c} \text{pinv} \\ \hline X^T * X \\ \downarrow \\ (X^T * X)^{-1} \end{array}$$

- Gradient descent : Feature scaling needed

Normal equation : No need.

m training examples, n features.

### Gradient Descent

- Need to choose  $\alpha$ .
- Needs many iterations.
- Works well even when n is large.

### Normal Equation

- No need to choose  $\alpha$ .
- Don't need to iterate.
- Need to compute  $(X^T X)^{-1}$   $O(n^3)$
- Slow if  $n$  is very large.

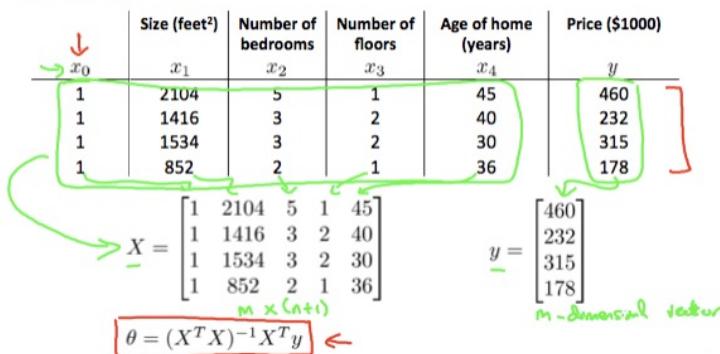
# Normal Equation

**Note:** [8:00 to 8:44 - The design matrix X (in the bottom right side of the slide) given in the example should have elements x with subscript 1 and superscripts varying from 1 to m because for all m training sets there are only 2 features  $x_0$  and  $x_1$ . 12:56 - The X matrix is m by (n+1) and NOT n by n.]

Gradient descent gives one way of minimizing J. Let's discuss a second way of doing so, this time performing the minimization explicitly and without resorting to an iterative algorithm. In the "Normal Equation" method, we will minimize J by explicitly taking its derivatives with respect to the  $\theta_j$ 's, and setting them to zero. This allows us to find the optimum theta without iteration. The normal equation formula is given below:

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

Examples: m = 4.



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

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

Gradient Descent	Normal Equation
Need to choose alpha	No need to choose alpha
Needs many iterations	No need to iterate
$O(kn^2)$	$O(n^3)$ , need to calculate inverse of $X^T X$
Works well when n is large	Slow if n is very large

With the normal equation, computing the inversion has complexity  $\mathcal{O}(n^3)$ . So if we have a very large number of features, the normal equation will be slow. In practice, when n exceeds 10,000 it might be a good time to go from a normal solution to an iterative process.

# Normal Equation Noninvertibility

Normal equation

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

$X^T X$

- What if  $X^T X$  is non-invertible? (singular/  
degenerate)  $\rightarrow$  Non-invertible
- Octave:  $\text{pinv}(X' * X) * X' * y$   
 $\ominus$  This will do right thing
  - ↑
  - $\rightarrow \text{pinv}$
  - $\rightarrow \text{inv}$

What if  $X^T X$  is non-invertible?

- Redundant features (linearly dependent).  
E.g.  $x_1 = \text{size in feet}^2$        $1m = 3.28 \text{ feet}$   
 $x_2 = \text{size in m}^2$   
 $x_1 = (3.28)^2 x_2$        $\rightarrow m = 10 \leftarrow$   
 $\rightarrow n = 100 \leftarrow$   
 $\Theta \in \mathbb{R}^{101}$
- Too many features (e.g.  $m \leq n$ ).
  - Delete some features, or use regularization.  
 $\downarrow$  later