

$m \rightarrow$  # of training examples

$x \rightarrow$  input variable / feature

$y \rightarrow$  output variable / target

$(x, y)$  - one training example

$(x^i, y^i)$  - a specific training example

$(x^2, y^2) \rightarrow$  2nd example

| size of House ( $x$ ) | Price \$ ( $y$ ) |  |
|-----------------------|------------------|--|
| 2104                  | 460              |  |
| 1416                  | 232              |  |
| 1534                  | 315              |  |
| 852                   | 178              |  |
| :                     | :                |  |

$m = 47$

Use training set get a hypothesis function, which is created by a learning algorithm.



The hypothesis can then predict an output ( $y$ ) for a certain input ( $x$ ).

$$\begin{matrix} \text{size of house} \\ x \end{matrix} \rightarrow \boxed{\text{hypo}} \rightarrow \begin{matrix} \text{est. value} \\ y \end{matrix}$$



$h$  is represented by a straight line  $\rightarrow h_0(x) = \theta_0 + \theta_1 x$   $\Rightarrow$  same as function for a straight line

This is linear regression with one variable  
The variable is  $x$  (Univariate linear regression)

$$\begin{aligned} h(x) &= \theta_0 + \theta_1 x \\ y &= a + b x \end{aligned}$$

↓  
y-int      ↓  
              slope

## Cost Function

- helps us figure out the best possible straight line to fit our data

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

$\theta_i$

are parameters

→ how do we choose those values??

Depending on the value, the line will change

① We want to choose values for  $\theta_0$  and  $\theta_1$  so that  $h(x)$  is close to the  $y$  for the training examples

② We want to minimize  $\theta_0$  and  $\theta_1$ . The difference between  $h_{\theta}(x)$  and  $y$  to be small  
 & the estimate so the actual value

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

minimize the sum of the squares  
of  $h(x) - y$

$h_{\theta}(x) - y$  should be small

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

Take the average

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

note later with easier

minimize the average of

use constant  $\frac{1}{2}$   $\Rightarrow \frac{1}{2m}$

gives same values of  $\theta_0, \theta_1$  as minimizing  
the function ??

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

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

find me the values of  $\theta_0, \theta_1$  that  
causes the hypothesis function  
to be minimized.

$\frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$  is expressed as  $J(\theta_0, \theta_1)$  so we are minimizing that  $\Rightarrow$

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

$\frac{1}{2} \bar{x}$  → average of the squares of  
predicted val. - actual

↓ cost function or squared error function (most commonly used for linear  
regression, but there are others)  
cost is higher for greater differences

## Gradient of Descent

Test hypothesis for different values of  $\theta$  ( $\theta_0, \theta_1$ ) to reach the minimum the quickest as possible

$$\left\{ \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad \text{for } j=0, j=1 \right\} \text{ repeat until convergence}$$

Assignment  
↓  
overwrite  $\theta_j$   
with new value

$\alpha \rightarrow$  Alpha  $\rightarrow$  learning rate, Controls the rate at which we descend (the step we take)

$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \rightarrow$  derivative term (explained later)  $\rightarrow$  It is the slope of the tangent line at a certain point of  $J(\theta_0, \theta_1)$



$\theta_0$  and  $\theta_1$  should be updated simultaneously

①  $\text{temp}0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$   $\triangleright$  first calculate both

②  $\text{temp}1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$   $\star$  NOT  $① \rightarrow ③ \rightarrow ② \rightarrow ④$

③  $\theta_0 := \text{temp}0$   $\triangleright$  then assign

④  $\theta_1 := \text{temp}1$

for  $\alpha$ , if it's too small, the steps will take a long time

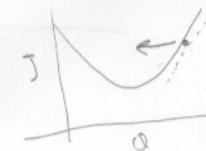


If its too large, it could overshoot the minimum, It could fail to converge or even diverge



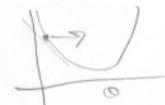
When the derivative has a positive slope,  $\theta_j := \theta_j - \alpha^{(\text{pos. num})}$

so  $\theta_j$  will be updated to a smaller number

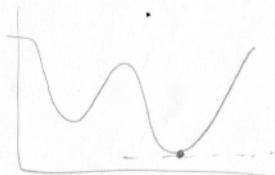


when derivative has a negative slope,  $\theta_j := \theta_j - \alpha^{(\text{neg. num})}$ , so  $\theta_j + (\alpha \cdot \text{num})$

$\theta_j$  will increase



## Gradient of Descent (cont)



If  $\theta$  is already at a minimum  
the slope of the derivative will be 0,

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \rightarrow \text{is } 0$$

$$\theta_j := \theta_j - \alpha \times 0$$

so  $\theta$  will not change!

Gradient of Descent can converge to local minimum even if  $\alpha$  is constant, because

the slope of the derivative naturally decreases,  
thus the descent also naturally takes smaller  
and smaller steps.  $\rightarrow$  No need to decrease  $\alpha$   
over time

$$\theta_j := \theta_j - \alpha (2)$$

$$\theta_j := \theta_j - \alpha (1)$$

$$\theta_j := \theta_j - \alpha (0.5)$$

Decreases gradually  
become smaller



## Gradient of Descent for Linear Regression p4

- Putting together Gradient of Descent with our hypothesis and Cost function

First we need to figure out the partial derivative term

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

↑ cost function      ↑ hypothesis

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

$$\textcircled{3} \quad \theta_0 := \theta_0 - \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \rightarrow \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

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

↑ Getting from  $\textcircled{2}$  to  $\textcircled{3}$  requires multivariate calculus but don't need to worry now

## Gradient descent algorithm for Linear Regression (cont.)

repeat until convergence {

$$\theta_0 := \theta_0 - \alpha \left( \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \right) \quad \xrightarrow{\text{derivative } \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)}$$

$$\theta_1 := \theta_1 - \alpha \left( \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \right) \quad \xrightarrow{\text{derivative } \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)}$$

}

} update  $\theta_0$  and  $\theta_1$ ,  
simultaneously!

For Linear Regression, the cost function is always a bowl-shaped function (convex function)  
 → No local optimum except for the global optimum

Using this to get the optimum hypothesis is called "Batch" Gradient Descent because it looks at the entire batch of training examples at each step.  
 (some algorithms use a sample of batches)

## Multivariate linear regression

- Predict price with multiple features, e.g. size of house, #bedrooms, #floors, Age, etc

| size          | # Bed | # floor | Age   | # f   | Y |
|---------------|-------|---------|-------|-------|---|
| 2104          | 5     | 1       | 45    | 460   |   |
| 1416          | 3     | 2       | 40    | 232   |   |
| :             | :     | :       | :     |       |   |
| feature $x_1$ | $x_2$ | $x_3$   | $x_4$ | $x_5$ | Y |

$$x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ \textcircled{2} \\ 40 \end{bmatrix} \quad x_3^{(2)} = 3$$

$n = \# \text{ of features}$

$m = \# \text{ of examples}$

$x^{(i)} = \text{input (e.g. vector) of all features for the } i\text{th training set} \rightarrow i = \text{index of set}$

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

features

$$\text{Hypothesis} = h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n \rightarrow \text{For convenience, define an extra feature} \rightarrow x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\text{Example } h_{\theta}(x) = 80 + 0.1x_1 + 0.01x_2 + 3x_3 - 2x_4$$

size      #bedrooms    #floor    Age

$$x_0 = 1 \text{ so } x_0^{(i)} = 1$$

$$\text{parameters}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

With the convenience feature, we can do matrix calculations!

↓

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

always

1 is same  
as above hypothesis

$$= \theta^T x$$

transpose

$$\begin{bmatrix} \theta_0 & \theta_1 & \dots & \theta_n \end{bmatrix} \underbrace{\begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}}_{\theta^T}$$

$1 \times (n+1)$  matrix  
or row matrix

→ that is the inner product between feature vector and parameter vector

## Gradient of Descent for Multiple Variables

- Basically the same as for one variable, except we have to repeat for "n" features!

repeat until convergence {

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

this is one  
for  $x_0$

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

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

:

}

↓

repeat until convergence {

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

}

### \* Feature scaling

- Dividing input values by the range ( $\max.\text{val} - \min.\text{val}$ )

$$\text{E.g. } x_1 = \frac{\text{size}}{2000} \quad \text{ends up } 0 \leq x_1 \leq 1$$

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

### \* Mean Normalization

- Subtracting Average for an input value resulting in a new average of just 0 for that input variable.

replace  $x_i$  with  $x_i - \bar{\mu}_i$  → Average value  
(Don't apply to  $x_0$ )

### Getting Gradient of Descent to Work Well

- Speed Gradient Descent up by making each input variable roughly in the same range ( $\theta$  descends quickly on small ranges but slowly for large ranges)

E.g. size of house 0 - 2000 → will slow down grad. des.  
but rooms 1 - 5

- Ideally, modify input variables so they range:

$$-1 \leq x_{(i)} \leq 1 \quad \text{or} \quad -0.5 \leq x_{(i)} \leq 0.5$$

\* As long as it's close or small range, doesn't have to break

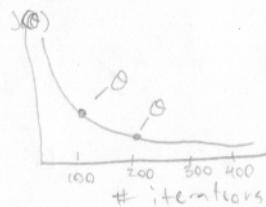
- Can use feature scaling or mean normalization

combine Average value of  $x$   
in training set

$$x_1 \rightarrow \frac{x_1 - \bar{\mu}_1}{\frac{\text{range (max-min)}}{\text{standard deviation}}}$$

Making sure Gradient of Descent is working properly

- Look at the cost  $J(\theta)$  after several iterations
- $J(\theta)$  should decrease after every iteration

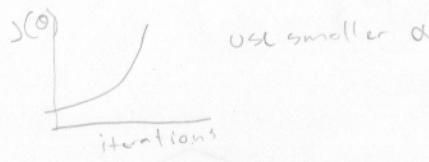


hard to predict how many iterations it will take. ~3000 in 3 mil?

You can also use automatic convergence test

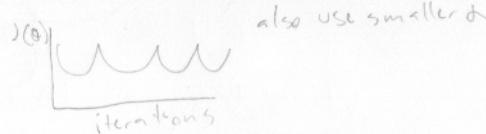
e.g. Declare convergence if  $J(\theta)$  decreases less than  $10^{-3}$  in one iteration

Not Working:



use smaller  $\alpha$

\* If  $\alpha$  is small enough, should decrease on every iteration. But too small is slow.



also use smaller  $\alpha$

Choosing  $\alpha$ :

$$\text{Try: } \dots, 0.001, 0.01, 0.1, 1 \dots$$

$\uparrow \quad \uparrow \quad \uparrow$   
0.001 0.01 0.1

Find value too small, too large

Features & Polynomial Regression

- We can combine features into one  
e.g. House frontage  $\times$  depth = area  
 $x_1 \quad x_2 \quad x_3$

• Hypothesis does not need to be linear!  
→ can change the hypothesis function by making it a quadratic, cubic, or square root function

E.g. If original hypothesis is

$h_{\theta}(x) = \theta_0 + \theta_1 x_1$ , then we can create additional features based on  $x_1$ ,

Quadratic

$$\hookrightarrow h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$$

Cubic

$$\hookrightarrow h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3$$

Sq.Rt

$$\hookrightarrow h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 \sqrt{x_1}$$

\* Feature Scaling is important here

because e.g.  $x_1 \rightarrow 1-1000$  range,  $x_1^2 \rightarrow 1-1,000,000$  range,  
 $x_1^3 \rightarrow 1-10^9$  range

## Normal Equation - an alternative to Gradient Descent

- solves for  $\theta$  analytically without taking several steps
- runs in one step!
- Does not need feature scaling!

| (Exm) | $(x_1)$<br>size | $(x_2)$<br>Bedrooms | $(x_3)$<br>Floors | $(x_4)$<br>Age | Value ( $y$ ) |
|-------|-----------------|---------------------|-------------------|----------------|---------------|
| 1     | 2104            | 5                   | 1                 | 45             | 460           |
| 1     | 1416            | 3                   | 2                 | 40             | 232           |
| 1     | 1534            | 3                   | 2                 | 30             | 315           |
| 1     | 852             | 2                   | 1                 | 36             | 178           |

Create M-matrix  
with inputs  
( $m \times (n+1)$  matrix)

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

create vector with  
outputs  
( $n$ -dimensional  
vector)

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

$$\theta = (X^T X)^{-1} X^T y \text{ gives you } \theta \text{ that minimizes cost}$$

$X^T$  transpose of  $X$

$(X^T X)^{-1}$  → inverse of full matrix

Octave:  $\text{pinv}(X^T * X) * X^T * y$

↑ takes inverse of matrix

in training ex,  $n$  features

p9

## Gradient Desc.

- \* Need to choose  $\alpha$
- \* Need many iterations
- \* Works well even w/  
a large  $n$

## Normal Eq.

- \* No need for  $\alpha$
- \* No iterations
- \* Need to compute  
 $(X^T X)^{-1}$   
↳ slow if  $n$  is  
very large  
(over 10,000)

constructing  $X$ :  $\theta = 135.1$

each training example looks like

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ x_3^{(i)} \\ x_4^{(i)} \end{bmatrix} \text{ e.g. } \begin{bmatrix} 1 \\ 2104 \\ 5 \\ 1 \\ 45 \end{bmatrix}$$

each example is  
transposed into  
a row of  $X$

$$X = \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \\ \vdots \\ -(x^{(n)})^T \end{bmatrix} = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$$

making  $m \times (n+1)$  matrix

|                                      | Midterm Exam           | Final Exam |
|--------------------------------------|------------------------|------------|
| Midterm Exam                         | (Midterm) <sup>2</sup> | Final      |
| Midterm Exam                         | 7921                   | 96         |
| Stogot of 100%                       | 5184                   | 74         |
| 72 ( $x_1^2$ )                       | 8836                   | 87         |
| Mean 79<br>Std Dev 9.4<br>(800, 9.4) | 4761                   | 78         |
| 69                                   |                        |            |

| (x) value | (n) | (x) | (y) | (x)   | (n) |
|-----------|-----|-----|-----|-------|-----|
| 0.01      | 20  | 0.1 | 2   | 0.013 | 1   |
| 0.03      | 0.0 | 0.5 | 8   | 0.011 | 1   |
| 0.06      | 0.0 | 0.8 | 6   | 0.021 | 1   |
| 0.07      | 30  | 1.0 | 3   | 0.008 | 1   |

Use polynomial regression so the model fits  $h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$ , where  $x_1$  is midterm and  $x_2$  is mid-term squared. Want to use feature scaling and mean normalization

What is the normalized feature  $x_2^{(2)}$ ?

$$(7921 + 5184 + 8836 + 4761) / 4$$

$$\frac{(x_2^{(2)} - \text{mean of all } x_2^{(2)})}{\text{range of } x_2}$$

$$8836 - 4761$$

$$\frac{4761 - 6675.5}{4075} = -0.37$$

$x_1$  has  $(1, n)$  shape

$$X = \begin{bmatrix} 1 & 0.1 & 0.5 & 0.8 \\ 20 & 0.0 & 0.0 & 0.0 \\ 20 & 0.1 & 0.5 & 0.8 \\ 20 & 0.5 & 0.8 & 0.8 \end{bmatrix}$$

$$Y = X \cdot \theta$$

return to previous slide

## Normal Equation - an alternative to Gradient Descent

- solves for  $\theta$  analytically without taking several steps
- runs in one step!
- Does not need feature scaling!

| (Exm) | $(x_1)$<br>size | $(x_2)$<br>Bedrooms | $(x_3)$<br>Floors | $(x_4)$<br>Age | Value ( $y$ ) |
|-------|-----------------|---------------------|-------------------|----------------|---------------|
| 1     | 2104            | 5                   | 1                 | 45             | 460           |
| 1     | 1416            | 3                   | 2                 | 40             | 232           |
| 1     | 1534            | 3                   | 2                 | 30             | 315           |
| 1     | 852             | 2                   | 1                 | 36             | 178           |

Create M-matrix  
with inputs  
( $m \times (n+1)$  matrix)

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

create vector with  
outputs  
( $n$ -dimensional  
vector)

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

$$\theta = (X^T X)^{-1} X^T y \text{ gives you } \theta \text{ that minimizes cost}$$

$X^T$  transpose of  $X$

$(X^T X)^{-1}$  → inverse of full matrix

Octave:  $\text{pinv}(X^T * X) * X^T * y$

↑ takes inverse of matrix

in training ex,  $n$  features

p9

## Gradient Desc.

- \* Need to choose  $\alpha$
- \* Need many iterations
- \* Works well even w/  
a large  $n$

## Normal Eq.

- \* No need for  $\alpha$
- \* No iterations
- \* Need to compute  
 $(X^T X)^{-1}$   
↳ slow if  $n$  is  
very large  
(over 10,000)

constructing  $X$ :  $\theta = 135.1$

each training example looks like

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ x_3^{(i)} \\ x_4^{(i)} \end{bmatrix} \text{ e.g. } \begin{bmatrix} 1 \\ 2104 \\ 5 \\ 1 \\ 45 \end{bmatrix}$$

each example is  
transposed into  
a row of  $X$

$$X = \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \\ \vdots \\ -(x^{(n)})^T \end{bmatrix} = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$$

making  $m \times (n+1)$  matrix