

[178]

★ Week 2 ★

--@★ Multiple features (variables) --

Size (feet ²)	No of bedroom	No 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

$m \rightarrow$ no of entries

$n \rightarrow$ no of features

$x^{(i)}$ \rightarrow input (features) of i^{th} training examples.

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

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

$x_3^{(2)} = 2$

Hypothesis

previously: $h_0(x) = \theta_0 + \theta_1 x$ only 1 feature.

Now: $h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$

$$(100 - (100 - 20)) = 180 + 0.1x_1 + 0.01x_2 + 3x_3 - 2x_4$$

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

$$\begin{aligned} \text{Let } x_0 &= 1 \\ (x_0^{(i)} &= 1) \end{aligned}$$

$$h_0(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

$$\boxed{h_0(x) = \theta^T X} \quad \text{Multivariate linear regression}$$

$$X = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \begin{array}{c} \text{---} \\ \theta^T \end{array} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad \begin{array}{c} \text{---} \\ x \end{array}$$

(n+1) x 1

--- @ Gradient descent for Multiple Variables ---

$$\text{Hypothesis: } h_0(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

Parameters: $\theta_0, \theta_1, \dots, \theta_n$ θ (n+1) dimensional vector

$$\text{Cost function: } \underbrace{J(\theta_0, \theta_1, \dots, \theta_n)}_{\text{vector}} = \frac{1}{2m} \sum_{i=1}^n (h_0(x^{(i)}) - y^{(i)})^2$$

Gradient descent:

$$\text{Repeat } \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1, \dots, \theta_n) \quad J(\theta)$$

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

Gradient Descent

- Previously ($n=1$)

Repeat {

$$Q_0 := Q_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})$$

$$Q_1 := Q_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

(simultaneously update Q_0, Q_1)

}

- New algorithm ($n \geq 1$)

Repeat {

$$Q_j := Q_j - \alpha \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

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

}

$$Q_0 := Q_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$Q_1 := Q_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

$$Q_2 := Q_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

-- @ Gradient Descent in Practice I --

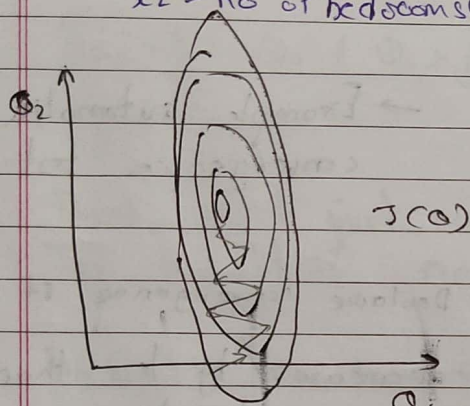
Tips to run gradient descent faster

1. feature scaling in gradient descent algo.

Idea: Make sure features are on similar scale.

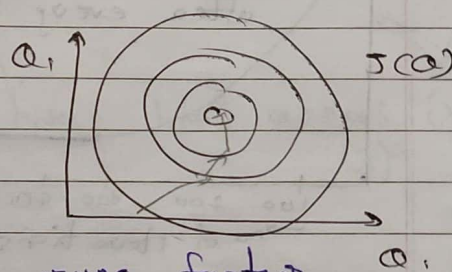
eg. $x_1 = \text{size (0-2000 feet}^2)$ $x_1 = \frac{\text{size (feet}^2)}{2000}$

$x_2 = \text{no of bedrooms (1-5)}$ $x_2 = \frac{\text{no of bedrooms}}{5}$



$0 \leq x_1 \leq 1$

$0 \leq x_2 \leq 1$



* Feature scaling

Get every feature into approximately

$x_0 = 1$ ← always right

$0 \leq x_1 \leq 3$ ✓

$-2 \leq x_2 \leq 0.5$ ✓

$-100 \leq x_3 \leq 100$ ✗

$-0.0001 \leq x_4 \leq 0.0001$ ✗

$-1 \leq x_i \leq 1$ range

↑ -3 fine 3 ↑

$\frac{1}{3}$ fine $\frac{1}{3}$

• Mean normalization:

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

$x_1 = \frac{\text{size} - 1000}{2000}$

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

$-0.5 \leq x_1 \leq 0.5$

$-0.5 \leq x_2 \leq 0.5$

$x_i = \frac{x_i - \mu_i}{s_i}$ ← average value of x_i in training set

range (max-min) / standard deviation

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

-- @ Gradient Descent in Practice II --

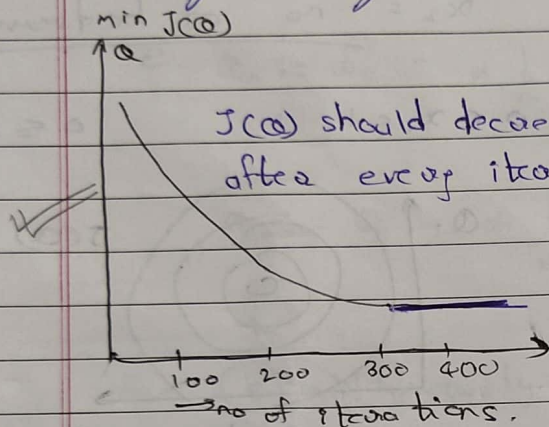
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

2>> Debugging: How to make sure gradient descent is working correctly.

ch

How to choose learning rate α

Making sure gradient descent is working correctly.

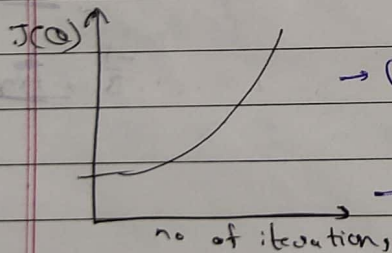


$J(\theta)$ should decrease after every iteration

→ Example automatic convergence test.

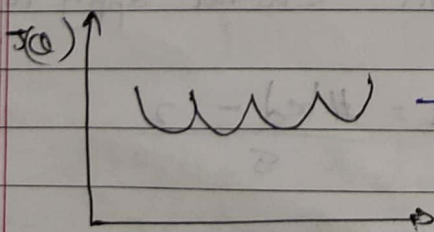
→ Declare convergence if $J(\theta)$ decreases by less than 10^{-3} in one iteration

→ Here we are very near to minimum.



→ Gradient descent is not working properly

→ use smaller α to work properly.



→ use smaller α .

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

To choose λ try.

0.0001, 0.01, 0.1, 1
0.003 0.03 0.3

-- @ Features and Polynomial Regression --

House price prediction

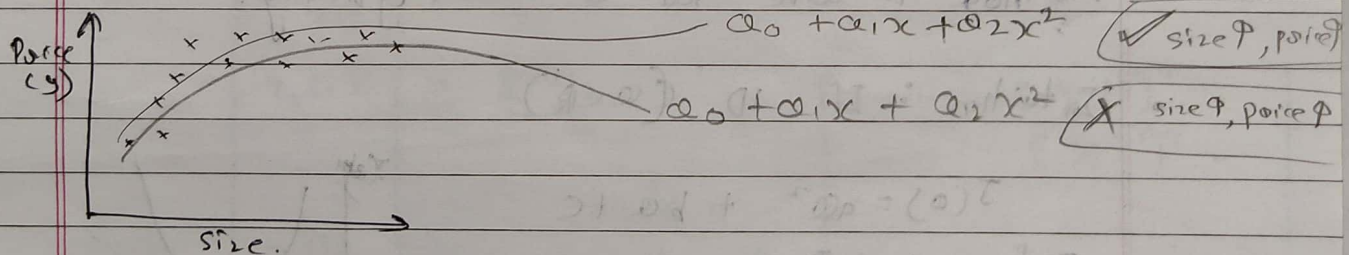
$$h_0(x) = \alpha_0 + \alpha_1 \underbrace{x \text{ footage}}_{x_1} + \alpha_2 \underbrace{x \text{ depth}}_{x_2}$$

Think, don't just use given,
Try to obtain new feature like area: $(x_1) * (x_2)$
(x)

$$h_0(x) = \alpha_0 + \alpha_1 x(\text{area})$$

This may make model accuracy \uparrow .

Polynomial regression: ~~single linear regression~~



$$h_0(x) = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3$$

$$= \alpha_0 + \alpha_1(\text{size}) + \alpha_2(\text{size}) + \alpha_3(\text{size})$$

$$x_1 = \text{size}$$

$$x_2 = \text{size}^2$$

$$x_3 = \text{size}^3$$

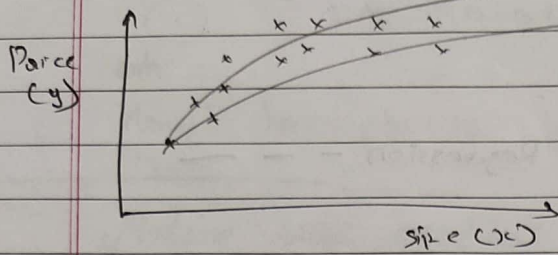
feature scaling gets importance.

$$x_1 : 1 - 1000 \text{ (size)}$$

$$x_2 : 1 - 1000000 \text{ (size)}^2$$

$$x_3 : 1 - 10^9 \text{ (size)}^3$$

choice of features.



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

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

que

model: $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 \sqrt{x}$

$x \rightarrow 1-1000 \text{ feet}^2$

$\sqrt{1000} = 32$

Ans

$x_1 = \frac{\text{size}}{1000}, x_2 = \frac{\sqrt{\text{size}}}{32}$

- @ Normal Equation --
- * Normal equation \rightarrow Without iterating, (running gradient descent) calculate θ_j in one go.
- Method to solve for θ analytically.

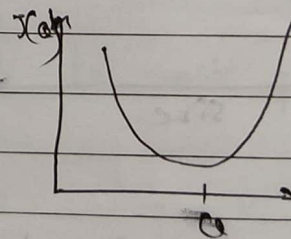
- Intuition : If $1D (\theta \in \mathbb{R}^1)$

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

To make $J(\theta)$ minimum

$$\frac{d}{d\theta} J(\theta) = \text{set } 0$$

solve for θ .



- $\theta \in \mathbb{R}^{n+1}$ $J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$
- $\frac{d}{d\theta_j} J(\theta) = \text{set } 0$ (for every j)
- solve for $\theta_0, \theta_1, \dots, \theta_n$

e.g. $m=4$

We have added this

	size (sq ft)	# bed	# floors	Age	Price
x_0	x_1	x_2	x_3	x_4	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

Design Matrix

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

$m \times (n+1)$ matrix

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

m -dimensional vector

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

General case

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

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

$(n+1)$ dimensional vector

$X =$
Design matrix

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

e.g. If $x^{(i)} = \begin{bmatrix} 1 \\ x_{1,i} \end{bmatrix}$

$$X = \begin{bmatrix} 1 & x_{1,1} \\ 1 & x_{1,2} \\ \vdots & \vdots \\ 1 & x_{1,m} \end{bmatrix}$$

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

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

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

Feature scaling is not necessary.

ex, $x_1 \leq 1$
 $0 \leq x_2 \leq 1000$
 $0 \leq x_3 \leq 10^{-5}$
are okay.

octave: ~~pinv(X)~~ $\text{pinv}(X' * X) * X' * y$

$$X^T \longrightarrow X'$$

m training examples, n features

Gradient Descent

Normal Eq.

- | | |
|------------------------------------|---|
| • Need to choose α | no need |
| • Need many iterations | no need. |
| • Works well even when n is large. | $(X^T X)^{-1}$
slow to invert if n is very large |

$$n = 10^6$$

$$n = 100, 1000, 10000$$

~~n = m~~ here, n \rightarrow no of features. Ignore

-- (a) Normal Equation Noninvertibility --

* Normal eqⁿ invertibility. (converge)

$$\text{Normal eq}^n: \theta = (X^T X)^{-1} X^T y$$

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

Issue \rightarrow what if $(X^T X)^{-1}$ is non-invertible? (singular / degenerate)

octave has two fun to inverse \therefore i.e. pinv inv

works ^{even} with $X^T X$ as singular

It computes θ .

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

\rightarrow where it can't be invertible

- Redundant features (linearly dependent)

- e.g. $x_1 = \text{size in feet}^2$

$x_2 = \text{size in m}^2$ $1\text{m} = 3.28\text{feet}$

dependent

features make

$(X^T X)$ non-invertible.

$$x_1 = (3.28)^2 x_2$$

So, keep only one among x_1, x_2 .

- Too many features ($m \leq n$)

Delete some features, or use regularization.

latter

-- (a) Basic operations --

* Octave

$5+6 \rightarrow \text{ans}=11$

$1/2 \rightarrow 0.5000$

$2^6 \rightarrow 64$

$1==2$ % false $\rightarrow \text{ans}=0$

$1 \sim 2 \rightarrow 1$

$1 \& 0 \rightarrow 0$

$1 \mid 0 \rightarrow 1$

$\text{xor}(1,0) \rightarrow 1$

$\text{PSK} \Rightarrow ') ;$

$a=3 \rightarrow a=3$

$a=3 ; \rightarrow$ nothing prints

$a \rightarrow 3$

$b = 'hi'$

$a = \pi ;$

$a \rightarrow 3.1416$

$\text{disp}(a) ; \rightarrow \text{print}(a)$

$\text{disp}(\text{sprintf}('2 \text{ decimals} : \%0.2f', a))$
 $\rightarrow 2 \text{ decimals} : 3.14$

format short

format long

$A = [1 \ 2 ; 3 \ 4 ; 5 \ 6]$

$A =$

1 2

3 4

$A = [1 \ 2 ;$

3 4 ;

5 6] same.

start inc step
N = 1: 0.1: 2

V =

columns 1 through 7:

1	1.1	1.2	1.3	1.4	1.5	1.6
1.7	1.8	1.9	2			

v = 1: 6

v = 1 2 3 4 5 6

ones(2,3)

1	1	1
1	1	1

C = 2 * ones(2,3)

2	2	2
2	2	2

w = zeros(1,3)

rand(2,3)

-	-	-	0 to 1 rand
-	-	-	
-	-	-	

w = randn(1,3)

w = -0.33 1.26 -0.2

w = -6 + sqrt(10) * (randn(1,10000))

hist(w)

hist(w, 50) → 50 bins.

eye(4)

1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1

help eye