


Multivariate Linear Regression

We have seen single variable (x) linear regression which is

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

Now, we have multivariate linear Regression:

Ex

Size (feet ²)	# of bedrooms	# of floors	Age of home	Price (\$100)
2104	5	1	45	460
1416	3	2	40	232
1532	2	2	30	315
852	2	1	76	178
\vdots			\vdots	\vdots
x_1	x_2	x_3	x_4	

Now we multiple variables of x and one y .

We call all \bar{x} 's $\Rightarrow n = 4$ (in above ex)

Training set size $\Rightarrow m = 4$ (4 rows consider)

$x^{(i)}$ The input feature's of training set (i) row

$$\text{Ex: } x^2 = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix} \quad x_{(i)}^2$$

$x_j^{(i)}$ Value of feature of m training set (i) row and j column.

$$\text{Ex: } x_{(4)}^2 = 40$$

Multivariable Hypothesis :-

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4.$$

also $x_0 = 1$ (reason:- vector multiplication)

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_{n+1} \end{bmatrix}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{n+1} \end{bmatrix}$$

we can do

$$\theta^T x$$

$$h_{\theta}(x) = [\theta_0, \theta_1, \theta_2, \dots, \theta_{n+1}] \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_{n+1} \end{bmatrix} = \theta^T x$$

Gradient Descent for multiple Variable :-

for single variable $\hat{\theta} = 1$

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

repeat until converge

$$\theta_0 := \theta_0 - \alpha \frac{d}{d\theta_0} J(\theta) := \theta_0 - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \right]$$

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta) := \theta_1 - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \underbrace{x^{(i)}}_{\frac{d}{d\theta_1} J(\theta_0, \theta_1)} \right]$$

(Simultaneously update θ_0, θ_1)

$$\theta_j := \theta_j - \alpha \frac{d}{d\theta_j} J(\theta) \quad \text{- common algo for Gradient Descent.}$$

Multiple variable Gradient descent algorithm

($D \geq 1$):

repeat until converges

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

for $j = 0, 1, 2, \dots, n$

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

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

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

⋮

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

Gradient Descent in practice 1 :- Feature Scaling

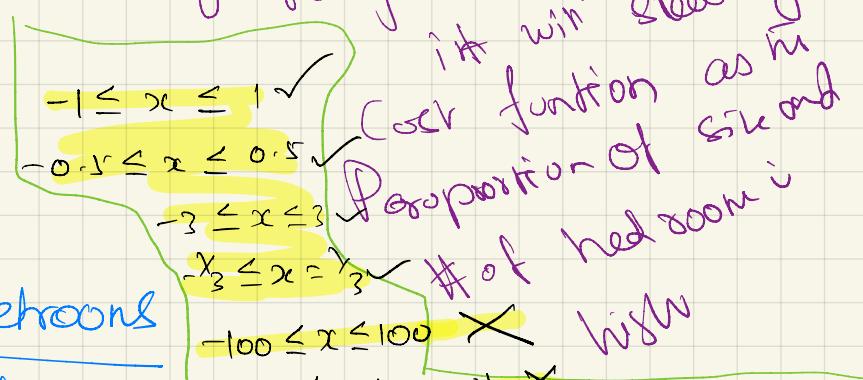
Ex: $x_1 = \text{Size (feet}^2)(1-2000)$
 $x_2 = \# \text{of bedrooms (1-5)}$

So you can scale in feature

like $x_1 = \frac{\text{Size}}{2000}$

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

This way the values will come



$$0 \leq x_1 \leq 1 \quad \text{and} \quad 0 \leq x_2 \leq 1$$

Mean Feature scaling:-

x_1 = Size feet (2000)

x_2 = # of bedrooms (1-5)

You can mean value of \bar{x} feature

you have \bar{x} formula

for mean feature scaling.

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

x_i = Feature

μ_i = Mean value of \bar{x} feature

s_i = range of \bar{x} feature or
(standardized deviation of \bar{x} feature)

Ex house above $\mu_1 = 1000 \quad \mu_2 = 2$

$$x_1 := \frac{\text{Size(feet)} - 1000}{2000} \quad (2000-1) = 1999.99$$

$$x_2 := \frac{\# \text{ of beds} - 2}{5 - 1} \quad (\text{it can } 5-1=4)$$

Gradient Descent in practice II - Learning Rate

You can choose your α -learning rate like:

$$\begin{matrix} 0.001 & , & 0.01 & , & 0.1 & , & 1 & \dots \\ \downarrow & & \downarrow & & \downarrow & & \dots & \\ 0.003 & & 0.03 & & 0.3 & & 3 & \dots \end{matrix}$$

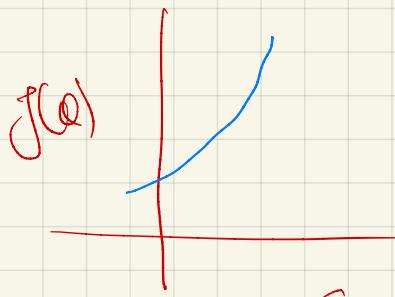
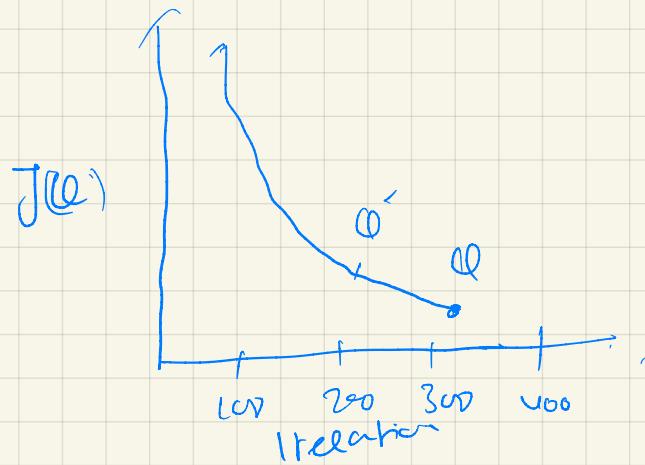
if your α is too small: slow convergence

α is too high: may not decrease on every iteration and

so use α small enough } thus may not converge
to converge soon enough

) from this you can flat $J(\alpha)$ vs iterations
it took to

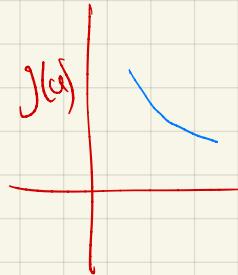
see which is
a good α .



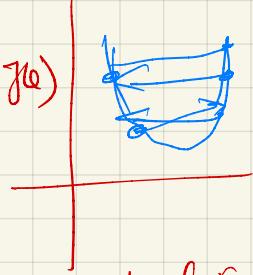
$\alpha = \text{too dirn}$



$\alpha = \text{Middle}$

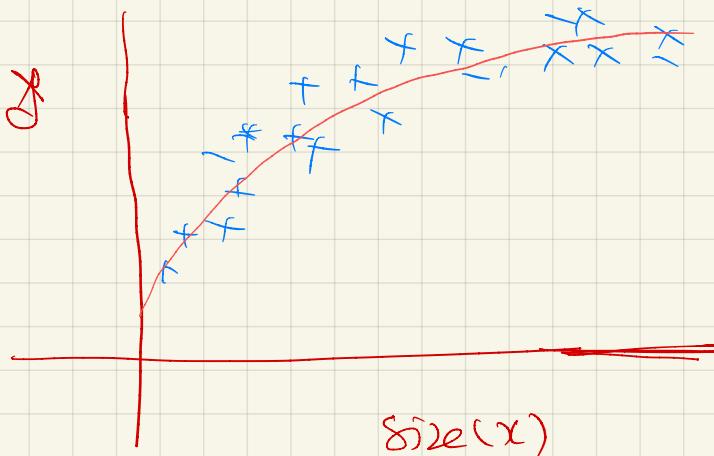


$\# \text{ of I}$
 $\alpha = \text{too slow}$, $\alpha = \text{too high}$



$\# \text{ of I}$

Feature and Polynomial Regression:-



So you use polynomial regression :-

if you use $h_0(x) = \theta_0 x_0 + \theta_1(x_1) + \theta_2(x_1)^2$
 quadratic function there is a change
 that it will curve down

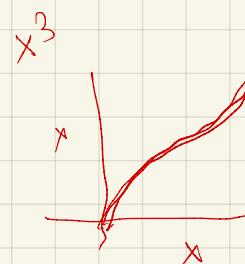
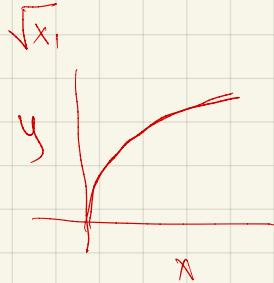
so use

$$h_0(x) = \theta_0 x_0 + \theta_1(x_1) + \theta_2(\sqrt{x_1})$$

square root function or

$$h_0(x) = \theta_0 x_0 + \theta_1(x_1) + \theta_2(x_1)^2 + \theta_3(x_1)^3$$

Cube root function. It will not curve down again



But in Feature scaling is very important when you do polynomial regression

Feature Scaling for Polynomial

size of house (1 to 1000)

then

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

$$\sqrt{x_1} = \frac{\sqrt{\text{Size}}}{\sqrt{1000}}$$

$$x_1^3 = \frac{(\text{Size})^3}{1000^3}$$

We can change behavior of hypothesis function by making quadratic, square root, cube root.

Normal equation: (2nd algorithm to solve for parameters Q)

Method

Only For specific model of linear regression

Gradient descent is one way of minimize J. and

Normal equation is the second way -

Examples: $m = 4$.

x_0	Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
	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

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

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

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

Ex shows $n = 4$

and we just make x_0 as all 1's
and

X of the matrix
 y of the matrix.

$$\text{Normal equation} = (X^T X)^{-1} X^T y$$

No need to feature scaling for normal equation

Gradient Descent

1. need to choose α
2. needs many iteration
3. $O(kn^2)$
4. works well with large n
5. choose if $n \geq 10^6$

vs

Normal Equation

1. no need to compute α
2. no need to iterate
3. $O(n^3)$, need to calculate $(X^T X)^{-1}$
4. slower if n is large
5. choose if $n \leq 10^5$

In future when we get to more complex algorithms, like classification and logistic algo, the normal equation method don't work and we have to choose Gradient descent. Very useful algo Gradient descent

Normal equation Noninvertibility

normal when a matrix can't be inverted

if it is linearly dependent (singular).

In normal equation we do

$(X^T X)^{-1} \Rightarrow$ There are some common issues.

1. Redundant features.

Ex: size of the house in foot^2 x_1 ,

size of the house in m^2 x_2

then you will end up with linearly dependent matrix which can't be invertible. So delete one of the feature.

2. $m \leq n$, if $m=10$ and $n=101$, then

it is too many features for m

so we have to find and remove some of the unnecessary features

Or use "regularization". (will come in future)

in octave

$\text{pinv}(X^T X) \Rightarrow$ solve the issue itself for you

Q When you compute $\text{inv}(X^T X)$ will return non-invertible because you have to manually correct it.

