

## ⇒ Week - 2 : Multivariable linear regression

Linear regression with multiple variables is known as multivariable linear regression.

### → Notation :

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

$x^{(i)}$  = the input (features) of the  $i^{\text{th}}$  training sample → Generally a vector

$m$  = the number of training samples

$n$  = the number of features

eg: Size (feet<sup>2</sup>)

	<u>Size (feet<sup>2</sup>)</u>	<u>No. of bedrooms</u>	<u>No. of floors</u>	<u>Age of home (years)</u>	<u>Price (\$1000)</u>
	$x_1$	$x_2$	$x_3$	$x_4$	$y$
$x^1$	2104	5	1	45	460
$x^2$	1416	3	2	40	232
$x^3$	1534	3	2	30	315
$x^4$	852	2	1	36	178

$m = 4$  ,  $n = 47$  (say) ,  $x_3^1 = 1$

$$x^1 = \begin{bmatrix} 2104 \\ 1416 \\ 1534 \\ 852 \end{bmatrix}$$

→ Hypothesis :

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

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_m x_m$$

→ Vectorization of hypothesis function

$$h_{\theta}(x^i) = [\theta_1 \quad \theta_2 \quad \dots \quad \theta_m]$$

$$= \theta^T x^i$$

$$\begin{bmatrix} x_0^i \\ x_1^i \\ \vdots \\ x_m^i \end{bmatrix}$$

\* For convenience reasons, we assume  $x_0^{(i)} = 1$

for  $(i \in 1, \dots, m)$

## ⇒ Gradient descent for multiple variables

Hypothesis :  $h_{\theta}(x) = \theta^T x$

$$= \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_m x_m$$

$\downarrow$   
 $x_0 = 1$

Parameters :  $\theta_0, \theta_1, \dots, \theta_m$

Cost function :

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

→ Gradient descent algo :

Repeat { until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_m)$$

$$\forall j = 0, \dots, m$$

}

$$\Rightarrow \frac{\partial J(\theta)}{\partial \theta_j} = \frac{2}{2m} \sum_{i=1}^m (h_{\theta} x^i - y_i) \frac{\partial}{\partial \theta_j} (h_{\theta} x^i - y_i)$$

$$h_{\theta} x^i = \theta_0 \underset{\downarrow}{x_0^i} + \theta_1 x_1^i + \theta_2 x_2^i + \dots + \theta_m x_m^i$$

$$\Rightarrow \frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta} x^i - y_i) \left[ \frac{\partial}{\partial \theta_j} [\theta_0 x_0^i + \dots + \theta_m x_m^i] \right]$$

$$\therefore \frac{\partial J(\theta)}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (h_{\theta} x^i - y_i) \left[ \frac{\partial}{\partial \theta_0} \theta_0 x_0^i \right]$$

As  $x_0^i = 1$  ~~⊗~~

$$\Rightarrow \frac{\partial J(\theta)}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (h_{\theta} x^i - y_i)$$

Similarly >

$$\frac{\partial J(\theta)}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m (h_{\theta} x^i - y_i) x_{1,i} \frac{\partial}{\partial \theta_1} \theta_1$$

$$\frac{\partial J(\theta)}{\partial \theta_2} = \frac{1}{m} \sum_{i=1}^m (h_{\theta} x^i - y_i) x_{2,i}$$

⋮

⋮

⋮



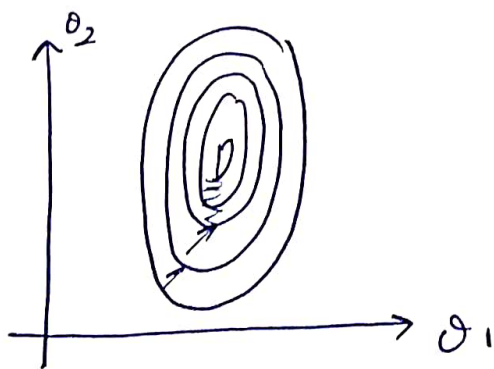
\* Notice ~~the~~ that  $\frac{\partial J(\theta)}{\partial \theta_0}$  &  $\frac{\partial J(\theta)}{\partial \theta_1}$  are same for both multivariate & linear regression.

→ As no. of feature = 1 in linear regression, so  $x_{1,i} = x^i$  {or  $x_i$ }

⇒ Gradient descent in practice : Feature scaling

\* If there is a huge range difference b/w the features, ~~it can take~~ running gradient descent on this kind of cost function can take a long time to find the global minimum.

eg :  $x_1 =$  Size (0 - 2000 ft)  
 $x_2 =$  No. of bedrooms (1-5)



\* 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 & slowly on large ranges, &

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,

$$-1 \leq x_i \leq 1 \quad \text{or} \quad -0.5 \leq x_i \leq 0.5$$

or  $0 \leq x_i \leq 2$

The goal is to get all input variables into roughly one of these ranges, give or take a few <sup>for this</sup> techniques:

1) Feature scaling: Divide input values by the range (i.e. the max value minus the min value) of the input variable.

2) Mean normalization: 
$$x_i := \frac{x_i - \mu_i}{\Delta_i}$$

$\mu_i$  = average / mean of all values for feature  $i$

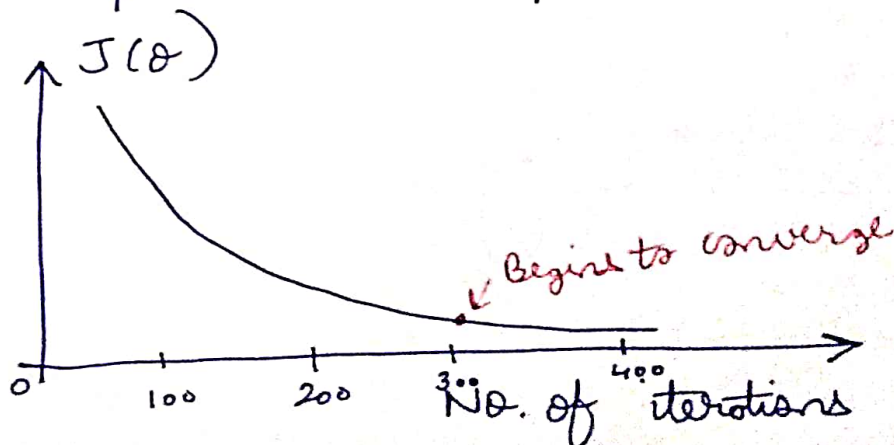
$\Delta_i$  = range of values (max - min), or it can be also be standard deviation.

⇒ Gradient descent in practice - Learning rate

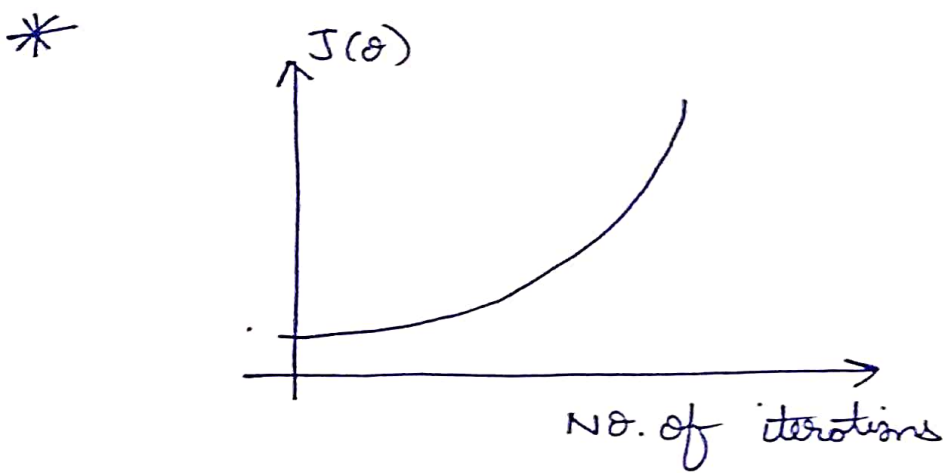
→ Debugging gradient descent: Make a plot with number of iterations <sup>of gradient descent</sup> on the  $x$  axis & the Cost function  $J(\theta)$  on the  $y$  axis.

If  $J(\theta)$  keeps decreasing, you're good to go! But <sup>if</sup>  $J(\theta)$  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. So, ~~for~~ whenever possible, just plot it!



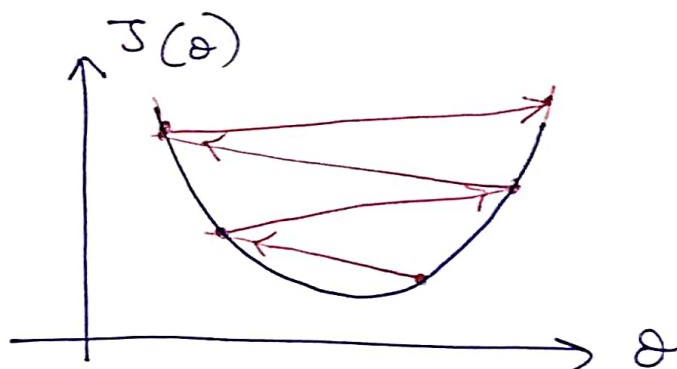




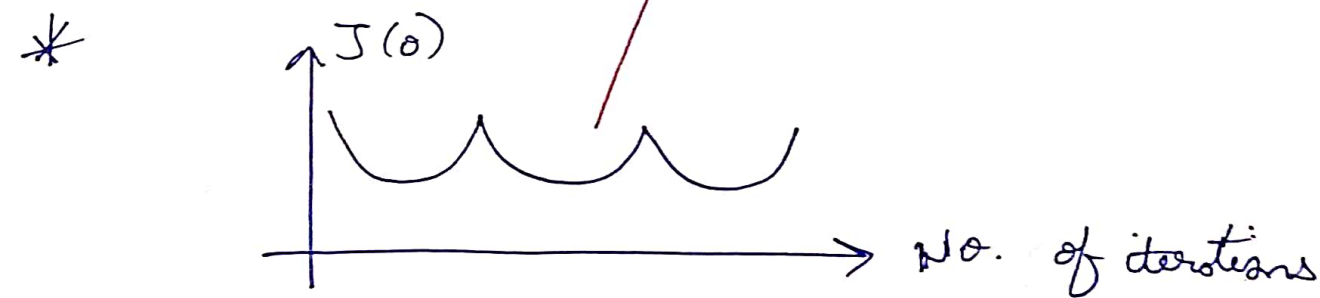
Back

Behind  
the  
scenes

$\Rightarrow$



So, choose a smaller  $\alpha$  in such case!



$\rightarrow$  It's proven that for sufficiently small  $\alpha$ ,  $J(\theta)$  should decrease on every iteration.

$\rightarrow$  But if  $\alpha$  is too small, gradient descent can be slow to converge.

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



## ⇒ Creating features

→ We can improve our features & the form of our hypothesis function in a couple different ways.

We can combine multiple features into one.

eg: House price prediction:

→ 2 features: frontage ( $x_1$ ) & depth ( $x_2$ )

→ You might decide that an imp. feature is the land area.

→ So create a new feature  $x_3 = x_1 \times x_2$

→  $h(x) = \theta_0 + \theta_1 x_3$

→ Area is indeed a better indicator

→ Often by defining new features, you may get a more efficient model.

## ⇒ Polynomial regression

\* Our hypothesis  $\text{func}^m$  need not be linear (a straight line) if that does not fit the data well.

\* We can change the behaviour or curve of our hypothesis  $\text{func}^m$  by making it quadratic, cubic, square root  $\text{func}^m$  (or any other form).

eg:  $h_\theta(x) = \theta_0 + \theta_1 x_1$

we can ~~eg~~ create additional features based on  $x_1$ , to get: ~~the~~

the quadratic  $\text{func}^m$ :  $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$

or the cubic  $\text{func}^m$ :  $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3$

or the square root  $\text{func}^m$ :  $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 \sqrt{x_1}$

→ How do we fit the model to this data?

\* Set,

$$x_1 = x$$

$$x_2 = x^2$$

$$x_3 = x^3, \quad x_4 = \sqrt{x}$$

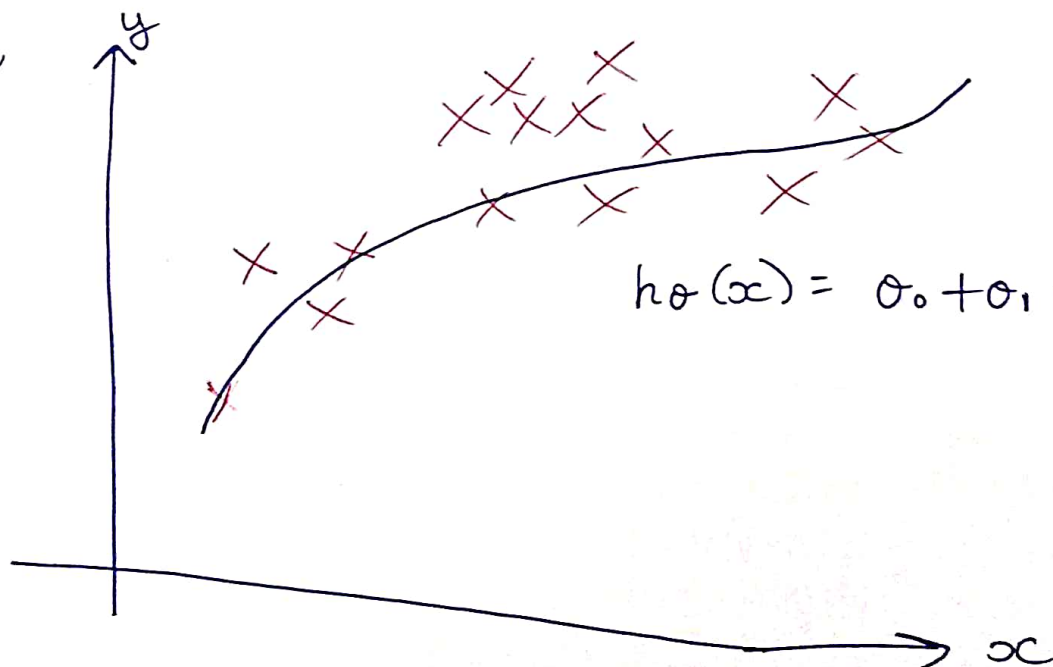
⋮

~~\*~~

\*

By selecting features like this & applying the linear regression algo, you can do polynomial linear regression.

→



Polynomial regression



\* Remember, feature scaling becomes ~~even~~ even imp. here.

eg: if  $x_i \in [1, 1000]$

then  $x_i^2 \in [1, 1,000,000]$

$\Delta$   $x_i^3 \in [1, 1,000,000,000]$





⇒ How does it work?

→ Supp. our cost func<sup>n</sup> depends ~~on~~ on only 1 variable,

Let  $J(\theta) = a\theta^2 + b\theta + c$ ;  $\theta \rightarrow$  real number

\*  $J(\theta) \rightarrow$  quadratic func<sup>n</sup>

\* How do we minimize this?

Ans: Put  $\frac{d}{d\theta} J(\theta) = 0$  & <sup>find</sup> ~~solve~~ <sup>values of  $\theta$ .</sup> ~~for the~~

~~value of~~

\* ~~Also~~ Hence, we find value of  $\theta$  which minimizes  $J(\theta)$ .

→ In more complex problems:

\*  $\theta \rightarrow 1 \times (n+1)$  vector of real numbers

\*  $J(\theta) \rightarrow$  func<sup>n</sup> of vector value

\* How do we minimize this?

Ans: → Calculate  $\frac{\partial J(\theta)}{\partial \theta_j} \forall j \in [0, n]$  & set to 0

→ Do that & solve for every value from  $[\theta_0, \theta_n]$ .

→ Gives values of  $\theta$  which minimize  $J(\theta)$ .

## ⇒ Normal equation

\* Instead of minimizing the cost function iteratively (through gradient descent), we can do it analytically - in other words, a mathematical eq<sup>n</sup> that gives the result directly. This is called the Normal eq<sup>n</sup>.

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

$\hat{\theta}$  → value of  $\theta$  that minimizes the cost function

$y$  → vector of target values containing  $y^1$  to  $y^m$ .

$X$  → "design matrix" — contains all the training data features in an  $[m \times m+1]$  matrix

$n$  → No. of features

$m$  → No. of training samples

## ⇒ Workflow of normal eq<sup>n</sup> method

eg:	Size (feet <sup>2</sup> ) $x_1$	NO. of bedrooms $x_2$	NO. of floor $x_3$	Age of home (years) $x_4$	Price (\$ 1000) $y$
$x^1$	2104	5	1	45	460
$x^2$	1416	3	2	40	232
$x^3$	1534	3	2	30	315
$x^4$	852	2	1	36	178

→  $m = 4, n = 4$

→  $X \equiv$  design matrix =

$$\begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}_{m \times (n+1)}$$

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

→ Compute  $(X^T X)^{-1} X^T y$  to get the value of  $\theta$  which minimize the cost function.



→ ~~X is cov~~

→ For  $m$  training samples &  $n$  features, ~~each~~ each training sample is a  $(n+1)$  dimensional feature column vector.

$$x^i = \begin{bmatrix} x_0^i \\ x_1^i \\ \vdots \\ x_m^i \end{bmatrix}_{1 \times (n+1)} \in \underset{\substack{\uparrow \\ \text{vector space}}}{\mathbb{R}^{n+1}}$$

→ ~~X~~ Design matrix  $X$  is constructed by taking each training sample, determining its transpose & using it ~~for~~<sup>as</sup> a row for itself.

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

→ \*\* If you're using the normal eq<sup>n</sup> method then there's no need for feature scaling!



→ Gradient descent

\* Need to choose  $\alpha$

\* Needs many iterations

\* Time complexity  $\rightarrow O(m^2)$

\* Work well when  $m$  is large ( $> 10^4$ )

Normal eq<sup>n</sup>

\* NO need to choose  $\alpha$

\* NO need to iterate

\* Time complexity  $\rightarrow O(m^3)$ , may need to calculate  $X^T X$

\* Slow when  $m$  is large ( $> 10^4$ )

⇒ Normal eq<sup>n</sup> non-invertibility

⇒ When computing  $(X^T X)^{-1} X^T y$ , what if  $X^T X$  is non-invertible (singular / degenerate)?

→ Octave / Matlab can invert ~~or~~ such invertible matrices with the `pinv()` function (pseudo inverse of  $\text{func}^m$ ). So use `pinv()` instead of `inv()`.

→ What does it mean for  $X^T X$  to be non-invertible?

2 causes generally:

1) Redundant features: where 2 features are very closely related (i.e. they are linearly dependent)

eg:  $x_1$  = size in feet

$x_2$  = size in meters

$$\therefore x_2 \approx 3.2 x_1$$

2) Too many features: (eg:  $m \leq n$ ,  $m=10$ ,  $m=100$ )

Not enough data. So delete some features or use "regularization".

\* Sol<sup>n</sup> to above problems include deleting a feature that is linearly dependent with another or deleting one or more features when there are too many features.