

WEEK - 2

• Multi-Variate Linear Regression

\textcircled{n} = number of features

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

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

∴ Hypothesis:

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

$$x_0 = 1$$

↳ for convenience of notation

vector

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

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

∴
$$h_{\theta}(x) = \theta^T \cancel{x}$$

• Gradient descent for Multiple Variables

Hypothesis:

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

parameters: $\theta \leftarrow n+1$ dimensional vector

cost function:

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

m = Training ex.
 n = features

Gradient descent: ($n \geq 1$)

Repeat {

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

until convergence

}

\Downarrow o/p

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

example : (for $n=2$)

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

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

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

* • Feature Scaling : To run gradient descent much faster

Idea : Make sure features are on a similar scale

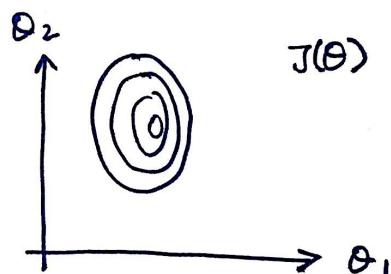
e.g. $x_1 = \text{size (0-2000 ft}^2)$

$x_2 = \text{no. of bedrooms (0-5)}$



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

$$0 \leq x_1 \leq 1$$



$$x_2 = \frac{\text{no. of bedrooms}}{5}$$

$$0 \leq x_2 \leq 1$$

→ We can speed up gradient descent by having each of our input values in roughly the same range. This is because θ will descend quickly on small ranges and slowly on large ranges, and so will oscillate down to the optimum when the variables are uneven ***

(15)

General Idea: Feature Scaling

Get every feature into approximately

$$[-1 \leq x_i \leq 1] \text{ range}$$

This is subjective, but the idea is they should be close enough

Mean Normalization

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

Eg.

$$x_1 = \frac{\text{size} - 1000}{2000}, \quad x_2 = \frac{\# \text{bedrooms} - 2}{5}$$

or

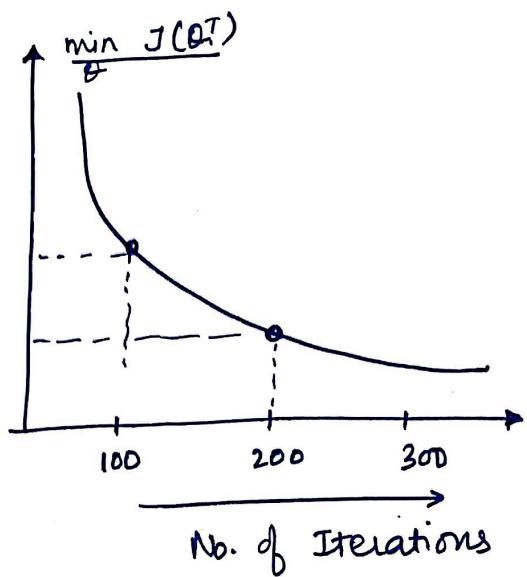
$$x_i \leftarrow \frac{x_i - \mu_i}{s_i}$$

Avg. value
of x in training set

range (max-min) of feature set/training
or standard deviation

(16)

- Gradient descent : Learning Rate α



+ $J(\theta^T)$ should dec.
after every iteration

← * look at such plots
to check how α is
working

- * Example autodecet :

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

* for sufficiently small α , $J(\theta)$ should decrease on every iteration

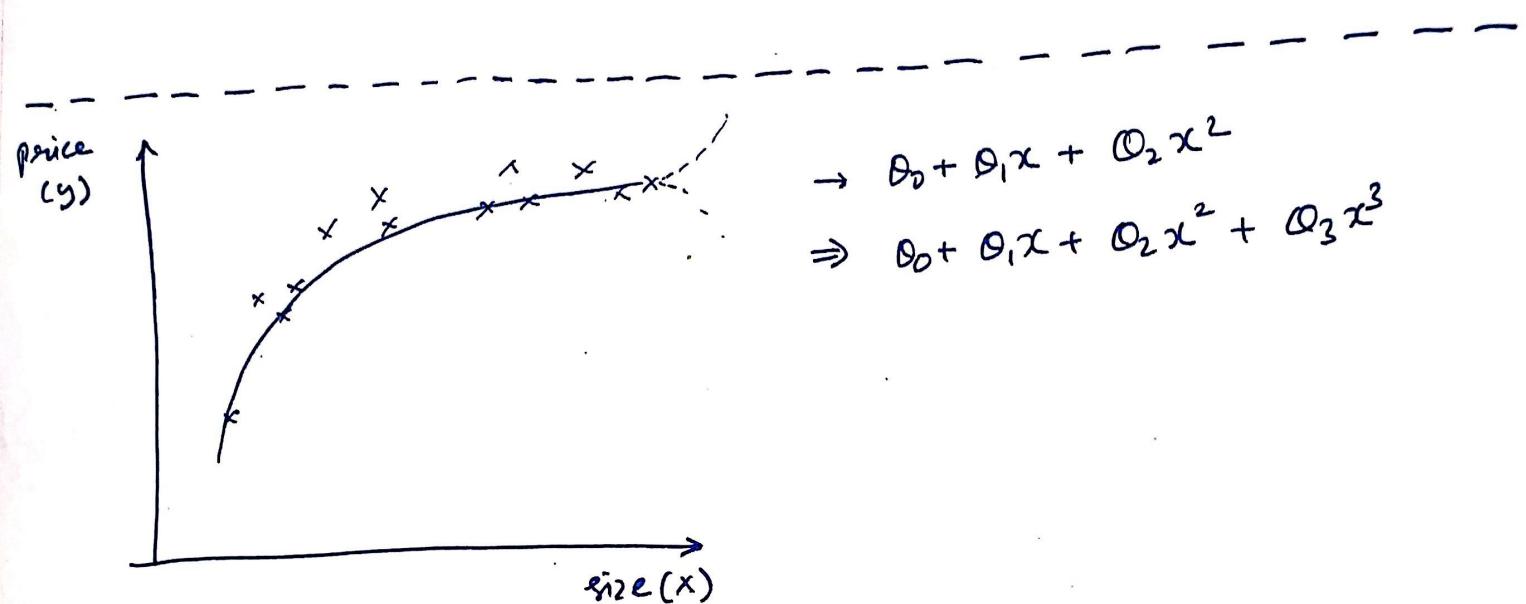
* But if α is too small, gradient descent can be slow to converge

• Features and Polynomial regression:

e.g. $h_{\theta}(x) = \theta_0 + \underbrace{\theta_1 \times \text{frontage}} + \underbrace{\theta_2 \times \text{depth}}$

$\frac{\text{Area}}{x} = \text{frontage} \times \text{depth}$

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



$$\begin{aligned} h_{\theta}(x) &= \theta_0 + \theta_1 x_1^2 + \theta_2 x_2^2 + \theta_3 x_3^3 \\ &= \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3 \end{aligned}$$

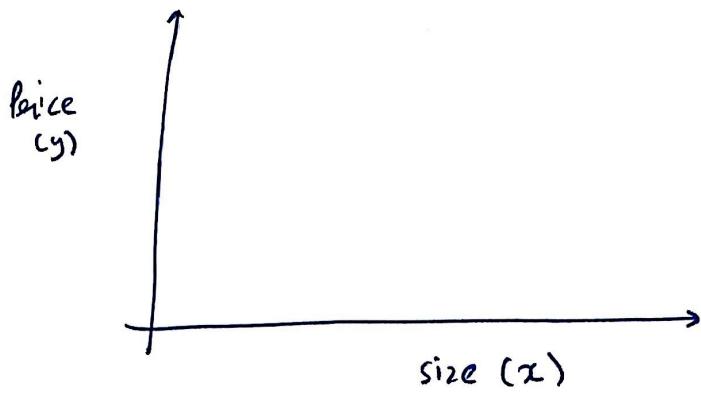
feature scaling becomes very important

$$\text{size} \rightarrow 1-1000$$

$$\text{size}^2 \rightarrow 1 - 10^6$$

$$\text{size}^3 \rightarrow 1 - 10^9$$

choices of feature:



$$h_{\theta}(x) = \theta_0 + \theta_1 \underbrace{\text{size}}_{x_1} + \theta_2 \underbrace{\sqrt{\text{size}}}_{x_2}$$

e.g.: if size range $\rightarrow 1-1000$

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

$$x_2 \leftarrow \sqrt{\text{size}} = \sqrt{\frac{\text{size}}{32}}$$

} feature scaling
without mean
normalization

● Computing Parameters Analytically :-

(Normal Equation) : solve for θ analytically

Intuition: if 1D ($\theta \in \mathbb{R}$)

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

$$\frac{\partial J(\theta)}{\partial \theta} = 0 \quad \left. \right\} \text{solve for } \theta$$

$\theta \in \mathbb{R}^{n+1} \rightarrow n\text{-dimensional}$

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

e.g:

$m=4$

x_1 size (ft ²)	x_2 No. of bedrooms	x_3 floors	x_4 Age of home	y Price (\$) ⁰⁰⁰
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178

Design Matrix

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

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

$m \times (n+1)$

$m \times 1$

$$\hat{\theta} = (X^T X)^{-1} X^T y \xrightarrow{\text{pinv } (X^T X)^* X^T * y} \quad \text{Octave.} \quad (20)$$

$$x^{(i)} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} \quad x \in \mathbb{R}^{n+1}$$

↓
Design Matrix

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

e.g.

$$x^{(i)} = \begin{bmatrix} 1 \\ x_1 \end{bmatrix} \quad X = \begin{bmatrix} 1 & x_1^{(1)} \\ 1 & x_2^{(1)} \\ \vdots & \ddots \\ 1 & x_m^{(1)} \end{bmatrix} \quad y = \begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^m \end{bmatrix}$$

* Feature scaling is not required in the normal-equation method.

choice between Gradient descent & Normal Equation

'm' training examples

'n' features

Gradient

Normal eqⁿ

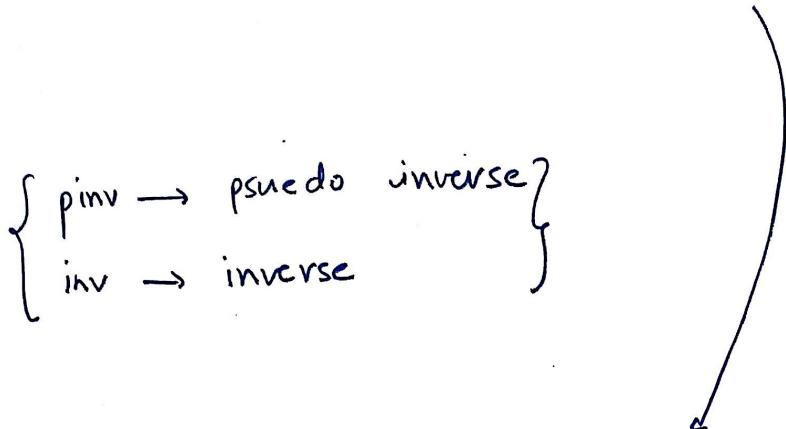
- Need to choose α .
 - Needs many iterations
 - Work well even when
n is large
- No need to choose α
 - Don't need to iterate.
 - Need to compute
 $(X^T X)^{-1}$
 - slow if n is very large

 $O(Kn^2)$ $O(n^3)$

• Non-Equation Noninvertibility

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

- what if $x^T x$ is non-invertible (singular / generate)



✗ Redundant features (linear dependence)

e.g. x_1 = size in ft^2

x_2 = size in m^2

✗ Too many features ($m \leq n$)

- delete some features, or use regularization.

MATLAB/OCTAVE

TUTORIAL

→ MATLAB / OCTAVE TUTORIAL :

①

- `disp(a)`
- `disp (sprintf ('2 decimals : %0.2f', a))` } string
- `disp (sprintf ('6 decimals : %0.6f', a))`

②

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

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}$$

3×2 matrix

$$A = [1 \ 2 \ 3] \quad \text{row vector}$$

$$v = [1; 2; 3] \quad \text{or}$$

$$v = 1:0.1:2$$

↓
start stop

③

$$c = 2 * \text{ones}(2,3)$$

$$\begin{bmatrix} 2 & 2 & 2 \\ 2 & 2 & 2 \end{bmatrix}$$

$$z = \text{zeros}(1,3)$$

$$\begin{bmatrix} 0 & 0 & 0 \end{bmatrix}$$

$$\underline{\text{rand}(3,3)} \leftarrow \text{range}(0,1) \rightarrow \underline{\frac{\text{randn}(1,3)}{\sigma}}$$

$$\begin{bmatrix} 0.08 & - & - \\ - & - & - \\ - & - & - \end{bmatrix}$$

$$\begin{bmatrix} -1.44 & 1.55 & - \end{bmatrix}$$

Identity matrix:

$I = \text{eye}(3)$

$$\rightarrow \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

✓ `help eye`

✓ `hist` → histogram function

Moving data around

* $\text{size}(A) \rightarrow [3 \downarrow \quad 2] \xrightarrow{\text{columns}}$
rows

* $\text{size}(A, 1) \rightarrow 3$

* $\text{size}(A, 2) \rightarrow 2$

* $\text{length}(A) \rightarrow 3$ highest dimension

(`pwd`) → gives the current direct.

(`cd 'C:\Users\Naren.mandolia\Desktop'`)



change directory

- * load features X.dat
- * load ('featuresX.dat') → filename
- * load ('Price Y.dat')

- * who → shows me variables in memory currently
- * who's → gives detail view
- * clear _____ } clears variable
- * $v = \text{Price Y}(1:10)$
- * save hello.txt Var -ascii : save as ASCII format

- * $A(3,2)$: element of Matrix
- $\begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}$

 $A(2,:) \Rightarrow [3 \ 4]$
 $A(:,2) \Rightarrow \begin{bmatrix} 2 \\ 4 \\ 6 \end{bmatrix}$
 $A([1\ 3],:) \Rightarrow \begin{bmatrix} 1 & 2 \\ 5 & 6 \end{bmatrix}$
- can also be used
for assignments.

- * $A = [A, [100; 101; 102]]$ → append another column vector there. to right.
- $\begin{bmatrix} 1 & 2 & 100 \\ 3 & 4 & 101 \\ 5 & 6 & 102 \end{bmatrix}$

(26)

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

$$B = [7 \ 8 ; 9 \ 10 ; 11 \ 12]$$

$$C = [A \ B] \rightarrow \begin{bmatrix} 1 & 2 & 7 & 8 \\ 3 & 4 & 9 & 10 \\ 5 & 6 & 11 & 12 \end{bmatrix} \quad [A, B]$$

↓
Treats
as
column

$$C = [A ; B] \rightarrow \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \\ 7 & 8 \\ 9 & 10 \\ 11 & 12 \end{bmatrix}$$

↓
Treats
as
row

Computing on data:

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

$$B = [11 \ 12 ; 13 \ 14 ; 15 \ 16]$$

$$C = [11 ; 2 \ 2]$$

operations:

$$\textcircled{1} \quad A^*B \quad [3 \times 2]. [3 \times 2]$$

$$\textcircled{2} \quad A \cdot ^*B \rightarrow \text{element wise multiplication}$$

③ $A.^2 \hookrightarrow$ element wise square.

④ $1.^A \hookrightarrow$ element wise invert

⑤ $\text{abs}(v) \hookrightarrow$ absolute value element wise

⑥ $v + \text{ones}(\text{length}(v), 1)$ or v+1

\hookrightarrow increments each element by 1

⑦ A'
 \hookrightarrow gives transpose

⑧ $\exp(v) -$

$\text{abs}(v)$

$\log(v) -$

⑨ $a = [1 \quad 15 \quad 2 \quad 0.5]$

$$\max(a) = 15$$

$$[\text{val}, \text{ind}] = \max(a)$$

$$\text{val} = 15$$

$$\text{ind} = 2$$

⑩ $a < 3 \rightarrow$ does element wise operation check.

⑪ $\text{find}(a < 3)$: $[r, c] = \text{find}(A \geq 7)$ ⑫

↳ returns the index, of values which are less than 3

$$r = \begin{matrix} \frac{1}{2} \\ \frac{2}{3} \\ 3 \end{matrix} \quad c = \begin{matrix} \frac{1}{2} \\ \frac{1}{2} \\ 3 \end{matrix}$$

⑫ $A = \text{magic}(3)$

↳ creates a magic matrix of 3×3

↖ sum = same

⑬ $\text{sum}(a)$

↳ sum of all the elements.

⑭ $\text{prod}(a) \rightarrow \text{prod. of all elements.}$

$\text{floor}(a) \rightarrow$

$\text{ceil}(a) \rightarrow$

⑮ $\text{rand}(3)$

↓
 3×3

$$\begin{bmatrix} - & - & - \\ - & - & - \\ - & - & - \end{bmatrix}$$

⑯ $\max(A, [], 1/2)$

column

row

→ Each matrix
operations are column
first.

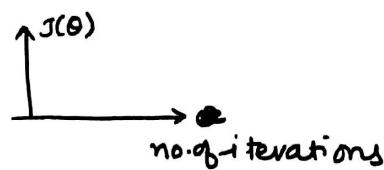
Plotting data :

(23)

```
t = [0 : 0.01 : 0.98]
```

```
y1 = sin(2 * pi * 4 * t);  
y2 = cos(2 * pi * 4 * t);
```

```
plot(t, y1)
```



```
plot(x, y) **
```

`hold on;` → holds on the previous graph.

`plot(t, y2, 'r');` → plots in red color

```
xlabel('time');  
ylabel('value');  
legend('sin', 'cos');  
title('my plot');
```

~~at print~~ `print -dpng 'myplot.png'`

↓ saves in directory.

```
close;
```

```
figure(1); plot(t, y1);  
figure(2); plot(t, y2);
```

```
subplot(1, 2, 1); % Divides plot a 1x2 grid, access first element
```

```
clf; → clears a figure/plot
```

◆: control statements : for, while, if

① $v = \text{zeros}(10, 1)$

FOR

```
for i = 1:10,
    v(i) = 2^i ;
end;
```

② **while**

i = 1

```
while i <= 5,
    v(i) = 100
    i = i + 1;
end;
```

③ if $i == 6$,

$\text{disp}(i)$;

elseif $i == 8$,

$\text{disp}(i)$;

else

$\text{disp}(i)$;

end;

◆ functions:

(31)

→

function $y = \text{squareThisNumber}(x)$
 $y = x^2$

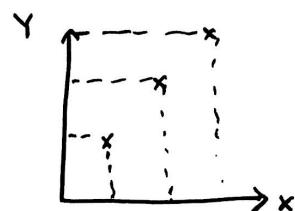
→

function $[y_1, y_2] = \text{sqAndCube}(x_1, x_2)$
 $y_1 = x_1^2$
 $y_2 = x_2^3$

* returns
multiple
values.

design matrix
 $\rightarrow \mathbb{X} = [1 \ 1; 1 \ 2; 1 \ 3]$

$$y = [1; 2; 3]$$



Goal: Define a fn. to compute the cost fn. $J(\theta)$

$$\text{theta} \Theta = [0; 1]$$

Sol: fn.

→ function $J = \text{cost Function}(X, y, \text{theta})$

$$m = \text{size}(X, 1);$$

$$\text{predictions} = X^* \text{theta};$$

$$\text{square error} = (\text{predictions} - y)^2;$$

$$J = \frac{1}{2m} * (\text{sum}(\text{square errors}));$$

$\xleftarrow{\text{Vectorization}}$

$$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j \Rightarrow \theta^T x$$



Unvectorized implementation

$$\text{pred} = 0$$

for $j = 1:n+1$

$$\text{predict} = \text{pred} + \theta_j * x$$

end;

vectorized implementation

$$\text{prediction} = \theta^T * x$$

vectorized implementation of GRADIENT DESCENT

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} \left[\frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \right]$$

$$\boxed{\theta := \theta - \alpha \delta}$$

$n+1$
dimensional
vector
 (R^{n+1})

$$\theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix} \quad \delta = \begin{bmatrix} \delta_0 \\ \vdots \\ \delta_n \end{bmatrix}$$