

House sizes:

2104
1416
1534
852

Have 3 computing hypothesis

1. $h_0(x) = -40 + 2.5x$

2. $h_0(x) = 200 + 0.1x$

3. $h_0(x) = -150 + 0.4x$

$$\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix} \times \begin{bmatrix} -40 & 200 & -150 \\ 0.25 & 0.1 & 0.4 \end{bmatrix} = \begin{bmatrix} 486 & 410 & 692 \\ 314 & 342 & 416 \\ 344 & 353 & 464 \\ 173 & 285 & 191 \end{bmatrix}$$

Properties for matrix multiplication

→ A and B, $AB \neq BA$ (not commutative)

→ $3 \times (5 \times 2) = (3 \times 5) \times 2$, $A \times (B \times C) = (A \times B) \times C$
(Associative) (matrix mult is associative)

→ Identity $I (I_{n \times n})$ $I \times A = A \times I = A$

→ Inverse $A A^{-1} = A^{-1} A = I$ (A must be a square matrix)

(singular, degenerate matrices) → don't have inverse.

⇒ Multiple features

earlier, we had only size of house on the basis of which we had to predict price.

Size (x_1)	No. of bedrooms (x_2)	No. of floors (x_3)	Age of house (x_4)	Price (1000\$) (y)	} $n = 47$
2104	5	1	45	460	
1416	3	2	40	232	
1534	3	2	30	315	
852	2	1	36	178	

Notation

n = number of features $n = 4$ here.

$x^{(i)}$ = input (features) of i^{th} training example

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

$$\rightarrow X^{(2)} = \begin{bmatrix} 1416 \\ 2 \\ 2 \\ 40 \end{bmatrix} \rightarrow \underbrace{\hspace{1cm}}_{\text{vector}}$$

4 features that determine price of 2nd house.

$X^{(i)} \rightarrow n$ dimensional vector.

$$\rightarrow X_3^{(2)} = 2 \quad X_1^{(2)} = 1416$$

Hypothesis

Previously $h_0(x) = \theta_0 + \theta_1 x$

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

eg. $h_0(x) = 80 + 0.1x_1 + 0.01x_2 + 3x_3 - 2x_4$

$\underbrace{80}_{\text{Price of house}} + \underbrace{0.1x_1}_{\text{effect of size on price}} + \underbrace{0.01x_2}_{\text{effect of bedrooms on price}} + \underbrace{3x_3}_{\text{effect of floors on price}} - \underbrace{2x_4}_{\text{effect of age on price}}$

$\rightarrow n$ features

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

we define $x_0 = 1$.

$$(x_0^{(i)} = 1)$$

\rightarrow we have basically defined an additional feature that always has value of 1 for any training set.

So $X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \rightarrow (n+1) \text{ dimensional vector}$

$\in \mathbb{R}^{n+1}$

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

\rightarrow 0 indexed vector.

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

\rightarrow we can write this as $\theta^T x$.

$$= \underbrace{\begin{bmatrix} \theta_0 & \theta_1 & \dots & \theta_n \end{bmatrix}}_{\text{parameter vector}} \underbrace{\begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}}_{\text{feature vector}} = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

This is called multivariate linear regression.

⇒ Gradient Descent for multiple features

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

Parameters: $\theta_0, \theta_1, \dots, \theta_n \rightarrow \theta \rightarrow n+1$ dimensional vector.

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

we'll call this as $J(\theta)$

where θ is vector $(n+1)$ dim.

Gradient Descent: Repeat {
 $\theta_j := \theta_j - \alpha \frac{d}{d\theta_j} J(\theta_0, \dots, \theta_n)$ $\nearrow J(\theta)$
} (simultaneously update for every $j = 0, 1, \dots, n$.)

Previously ($n=1$):

Repeat { $\theta_0 := \theta_0 - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})}_{\frac{d}{d\theta_0} J(\theta)}$

$\rightarrow \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)}) \cdot x_1^{(i)}$ $\xrightarrow{\text{now } x_1^{(i)}}$
(simultaneously update θ_0, θ_1)
}

Now ($n \geq 1$): Repeat {

$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)}) x_j^{(i)}$

(simultaneously update θ_j for $j = 0, 1, \dots, n$.)
}

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_0 x^{(i)} - y^{(i)})^2$$

$$= \frac{1}{2m} \left((h_0 x^{(1)} - y^{(1)})^2 + (h_0 x^{(2)} - y^{(2)})^2 + \dots \right)$$

$$= \frac{1}{2m} \left(\theta_0 x_0^{(1)} + \theta_1 x_1^{(1)} + \dots + \theta_n x_n^{(1)} - y^{(1)} \right)^2 + \left(\theta_0 x_0^{(2)} + \theta_1 x_1^{(2)} + \dots + \theta_n x_n^{(2)} - y^{(2)} \right)^2$$

$$\frac{dJ(\theta)}{d\theta_j} = \frac{1}{2m} \left(2(h_0 x^{(1)} - y^{(1)}) \cdot x_j^{(1)} + 2(h_0 x^{(2)} - y^{(2)}) \cdot x_j^{(2)} + \dots \right)$$

$$= \frac{1}{m} \sum_{i=1}^m (h_0 x^{(i)} - y^{(i)}) \cdot x_j^{(i)}$$

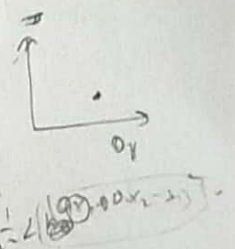
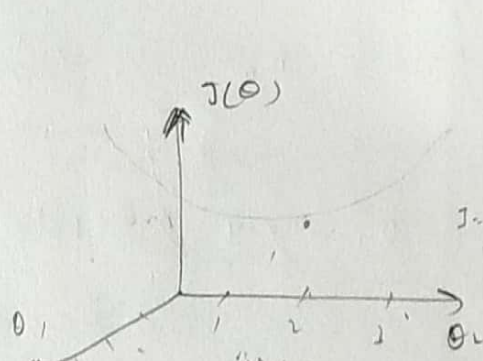
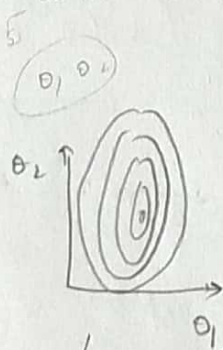
⇒ 4D in practice : feature scaling

Idea: Make sure features are on a similar scale.

{ Eg. $x_1 = \text{Size (0-2000 feet}^2\text{)}$ } makes 4D converge more quickly.
 $x_2 = \text{No}^{\circ} \text{ of bedrooms (1-5)}$

↳ let's draw contours

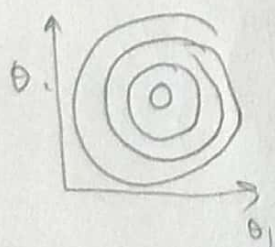
$$h(x) = \theta_1 x_1 + \theta_2 x_2$$



as x_1 and x_2 have varied ranges contours will be distorted. , here we can scale the features

$$x_1 = \frac{\text{Size (feet}^2\text{)}}{200}$$

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



Now, contours would be circular.

4D would work much easily
 here is the trajectory would be easy

→ So generally, get every feature into approx. $-1 \leq x_i \leq 1$ range

also like $\left. \begin{array}{l} 0 \leq x_1 \leq 3 \\ -2 \leq x_2 \leq 0.5 \end{array} \right\} \rightarrow$ These would also work as range is not far however.

$-1000 \leq x_3 \leq 1000 \} X \rightarrow$ This wouldn't work.

$-0.00001 \leq x_4 \leq 0.00001 \} \rightarrow$ not workable.

→ Mean normalisation

Replace x_i with $\frac{x_i - \mu_i}{s_i}$ to make features have approx. zero mean (not for $\mu_i = 1$)

eg. $x_1 = \frac{\text{size} - 1000}{2000} \rightarrow \text{Avg} = 1000$
 $x_2 = \frac{\text{bedrooms} - 2}{5} \rightarrow \text{Avg} = 2$

so $-0.5 \leq x_1 \leq 0.5$ and $-0.5 \leq x_2 \leq 0.5$.

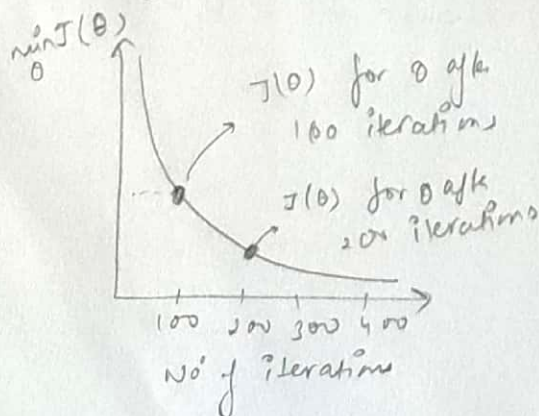
General rule: replace x_i with $\frac{x_i - \mu_i}{s_i}$
 where $\mu_i \rightarrow$ average
 $s_i \rightarrow$ range = (max - min value)

⇒ Gradient Descent in Practice 2: Learning rate (α)

$$\text{GD: } \theta_j^0 := \theta_j^0 - \alpha \frac{d}{d\theta_j} J(\theta)$$

- 'Debugging': How to make sure GD is working correctly.
- How to choose learning rate α .

Making sure GD is working correctly



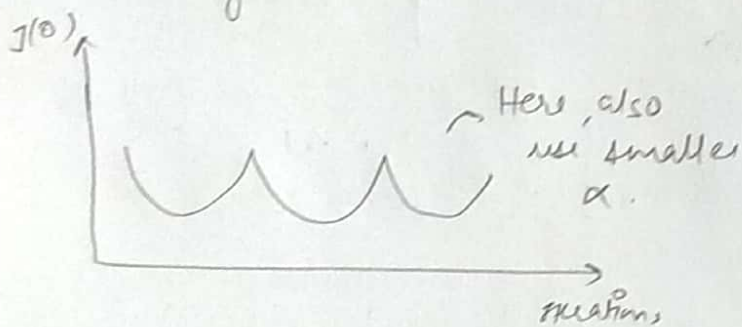
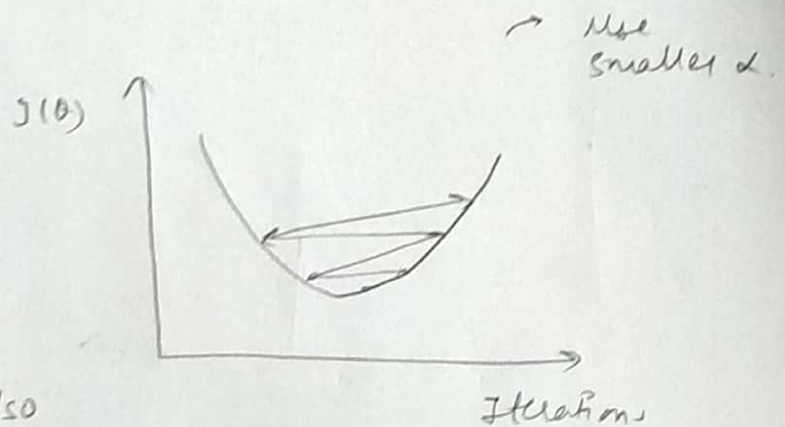
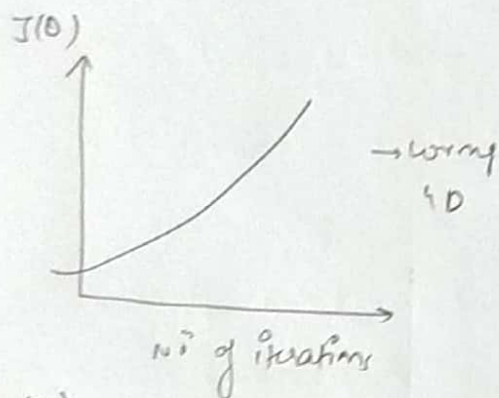
Thus, if GD is working correctly, then $J(\theta)$ decreases with number of iterations

Also, from the curve, it is clear that $J(\theta)$ decreases slowly as iterations increase.

We can also do automatic convergence test

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

However, guessing the number 10^{-3} beforehand is difficult.



→ For sufficiently small α , $J(\theta)$ should decrease in every iteration

→ But if α is too small, α can be slow to converge.

Summary :

- If α is too small : slow convergence
- If α is too large : $J(\theta)$ may not decrease in every iteration, may not converge (slow converge also possible)

Now, to debug all this, graph can help

To choose, α , try 0.003, 0.01, 0.1
0.001, 0.01, 0.1, 1 and choose the α that seems to cause $J(\theta)$ to decrease rapidly.

⇒ Features and Polynomial Regression

Housing price prediction.

$$h_0(x) = \theta_0 + \theta_1 \times \underbrace{\text{frontage}}_{x_1} + \theta_2 \times \underbrace{\text{depth}}_{x_2}$$

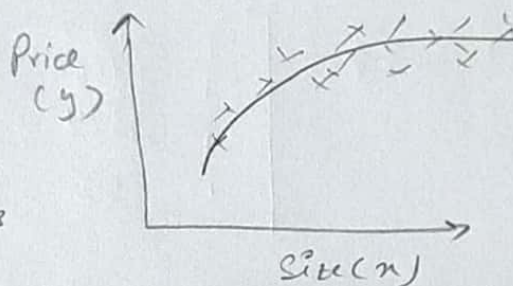
$$\text{Area} = X = \text{frontage} \times \text{depth}$$

$$\text{let } h_0(x) = \theta_0 + \theta_1 X \rightarrow \text{land area}$$

Polynomial Regression

$$\text{may be } \theta_0 + \theta_1 x + \theta_2 x^2 + \dots$$

$$\text{or } \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$



$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$= \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2 + \theta_3(\text{size})^3 \quad \left. \begin{array}{l} \text{size} \\ \text{size}^2 \\ \text{size}^3 \end{array} \right\} \text{ set features}$$

$$x_1 = \text{size} \quad x_2 = \text{size}^2 \quad x_3 = \text{size}^3$$

choice of features is important.

For ex.

$$h_0(x) = \theta_0 + \theta_1(\text{size}) + \theta_2(\sqrt{\text{size}})$$

feature scaling becomes important

size	-	1-1000
size ²	-	1-1,000,000
size ³	-	10 ⁹

⇒ Normal equation → helps to solve parameters.



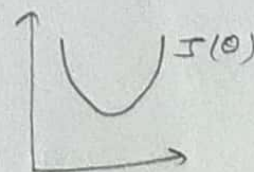
→ Iterative algorithm,

but now, we solve for $J(\theta)$ min analytically

Intuition: If $JD(\theta \in \mathbb{R})$

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

$$\frac{dJ(\theta)}{d\theta} = 0 \rightarrow \text{and get } \theta$$



but for us

$$(\theta \in \mathbb{R}^{n+1}) \quad J(\theta_0, \theta_1, \theta_2, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m ((h_0(x^{(i)}) - y^{(i)})^2)$$

$$\text{so here, } \frac{\partial}{\partial \theta_j} J(\theta) = 0 \dots = 0 \text{ (for every } j)$$

then, solve for $\theta_0, \theta_1, \dots, \theta_n$

Example $m = 4$

x_0	Size (feet ²) x_1	No. of bedrooms x_2	No. of floors x_3	Age of house x_4	Price ($\times 1000$) 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

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

$m \times (n+1)$

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

$(m \times 1)$

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

θ are the parameters for minimum $J(\theta)$ (cost function)

\rightarrow m examples $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^{(m)}, y^{(m)})$; n features

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

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

(design matrix)

eg. if $x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \end{bmatrix} \rightarrow X = \begin{bmatrix} 1 & x_1^{(1)} \\ 1 & x_1^{(2)} \\ \vdots & \vdots \\ 1 & x_1^{(m)} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$

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

Also, feature scaling is unnecessary in this method, or eventually we'll get some value of θ

GD

- \rightarrow Need to choose α
- \rightarrow Needs many iterations
- \rightarrow Works well even when n is large

Normal equation

- \rightarrow No need to choose α
- \rightarrow Don't need to iterate
- \rightarrow Need to compute $(X^T X)^{-1}$ slow if n is large.

$n = 1000 \rightarrow$ works fine for normal.

$n = 10000 \rightarrow$ normal.

$n = 106 \rightarrow$ GD.

⇒ Normal equations and non-invertibility

$\theta = (X^T X)^{-1} X^T y$ → what if $X^T X$ is non invertible.
(singular / degenerate)?

$(X^T X) \rightarrow$ is non invertible.

↳ this can be because of → (i) Redundant features (linearly dep.)

(ii) Too many features (eg. $m \geq n$)

eg. $x_1 = \text{cirt in ft}^2$
 $x_2 = \text{size in m}^2$

- Delete some features, or use regularisation.
later.

⇒ MATLAB onramp

→ ; , suppress command output is suppressed

→ MATLAB does not calculate all variables if some previous one is changed.

→ save filename - saves workspace variables

→ load filename - loads workspace variable

→ clear - clears workspace
clear - clears command window

→ $x = 1:4 \rightarrow x = [1 \ 2 \ 3 \ 4]$

→ $x = [1; 2; 3] \rightarrow x = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$

→ $x = 3:2:13 \rightarrow x = [3 \ 5 \ 7 \ 9 \ 11 \ 13]$

→ $x = \text{linspace}(\text{first}, \text{last}, \text{number of elements})$

$x = \text{linspace}(1, 10, 5) \rightarrow x = [1 \ 3.25 \ 5.5 \ 7.75 \ 10]$

→ $x = x'$ (transpose operator)

→ $x = \text{rand}(3)$ → $\begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$ random numbers.
random number matrix

→ $x = \text{rand}(2, 3) \rightarrow 2 \times 3$ matrix is formed.

→ $x = \text{zeros}(2, 3) \rightarrow \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$

→ $\text{size}(x) \rightarrow 2 \times 3$.

data is a 7×4 matrix

→ $x(3)$ index of array $x(\text{row}, \text{col})$

→ $x = \text{data}(6,3) \rightarrow$ element in 6th row, 3rd column

→ $x = \text{data}(\text{end}, 3) \rightarrow \text{last row, 3rd colm}$

→ $x = \text{data}(11) \rightarrow$ traverses columnwise

→ $x = \text{data}(:, 2)$ → give 2nd column all rows

→ $\begin{bmatrix} 3 & 4 \end{bmatrix} * \begin{bmatrix} 10 & 20 \end{bmatrix} \rightarrow \text{Error. (Matrix multiplication)}$

$$\rightarrow [3 \ 4] \cdot [10 \ 20] = [30 \ 80] \quad (\text{element wise multiplication})$$

→ $\text{plot}(x, y)$ - plot same sized vectors against each other

→ plot (x, y, 'r--o') → red colour
→ -- → line style
→ o → marker style

→ hold on → plot on same curve

→ hold off → prot on diff. alleles.

→ plot(v1) → plotting single vector plots values against index of array.

→ `plot(sample, v1, 'ro-', "linewidth", 4)`

→ title ("sample mass") → give title to previous plot

→ $y_{\text{label}}("mass(g)") \rightarrow y_{\text{ans}} \text{ label.}$

→ legend ("Exp A", "Exp B") → add legend over graph

→ Tables (add details)

□ → Exp A
★ → Exp B

→ Tables → load datafile
= elements

$d = \text{element_density}$ \rightarrow assigns density column to variable d