

<u>Size (feet²)</u>	<u>Price (\$1000)</u>	<u>Bedrooms</u>	<u>Linear regression</u>
2184	400	3	
1916	232	2	
1537	315	3	
852	178	3	
		price	

size →

Training Set

(How do you represent h ?)

Learning Algo.

$$h(x) = \theta_0 + \theta_1 x_1$$

size

"h func"

→ estimated price 2 parallel features

"hypothesis func"

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$$h(x) = \sum_{j=0}^n \theta_j x_j$$

$$x_1 = \text{size } 3^{1/2} \text{ feature}$$

where $\sum x_0 = 1$ for n parameters
or $n+1$ features

$$x = \text{"Inputs"} / \text{features} \quad \theta = \text{parameters} \quad h(x) = \sum_{j=0}^n \theta_j x_j$$

$y = \text{"Output"}$ $M = \# \text{ training examples}$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

Target variable
(# rows in table above)
 (x, y) = training example

$n = \# \text{ features}$

for n features there are $n+1$ parameters (θ)

$(x^{(i)}, y^{(i)})$ = i^{th} training example

$$n_1^{(1)} = 2104 \quad \left| \begin{array}{l} \text{choose } \theta \text{ s.t. } h(\mathbf{x}) \approx y \\ \text{for training examples} \\ \text{Then } \theta(\mathbf{x}) = h(\mathbf{x}) \end{array} \right.$$
$$n_1^{(2)} = 1416$$

Minimize $\frac{1}{2} \sum_{i=1}^n (h_\theta(\mathbf{x}^{(i)}) - y^{(i)})^2$

and θ

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

Choose θ s.t. $J(\theta)$ is minimized

Gradient descent

Start with θ (say $\theta = \vec{0}$)

Keep changing θ to resolve $J(\theta)$

$$\theta_j^0 := \theta_j - \alpha \left(\frac{\partial}{\partial \theta_j} J(\theta) \right) \quad (\theta = \theta_1, \theta_2, \dots)$$

learning rate

assignment operator $\rightarrow a := a + 1$
 \hookrightarrow increment value
of a by 1

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_\theta(\mathbf{x}) - y)^2$$

$$= \lambda \frac{1}{2} (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_{\theta}(x) - y)$$

$$= (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (\theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n - y)$$

for $j=1 \Rightarrow (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_1} (\theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n - y)$

$$= (h_{\theta}(x) - y) \cdot x_1$$

$$\therefore = (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (\theta_0 x_0 + \theta_1 x_1 + \dots + \theta_j x_j + \dots + \theta_n x_n)$$

$$= (h_{\theta}(x) - y) \cdot x_j$$

$$\theta_j := \theta_j - \alpha (h_{\theta}(x) - y) \cdot x_j$$

(only one training example)

Repeat until convergence.

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

for all training examples

$$\left[\frac{\partial J(\theta)}{\partial \theta_j} \right]$$

(for $j=0, 1, \dots, n$) a.k.a. Batch gradient descent

Repeat {

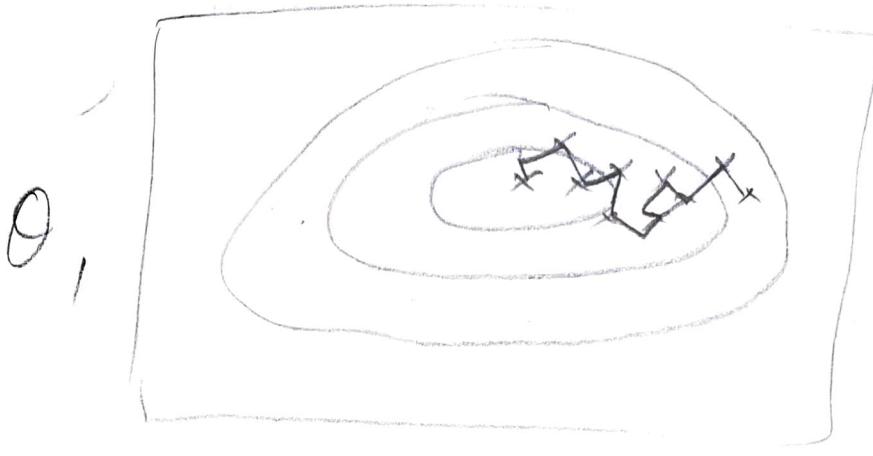
for every j

For $i = 1 \text{ to } m$ }

$$\theta_j := \theta_j - \alpha (\ln(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

}

} \Rightarrow stochastic gradient descent



$$\theta_0$$

$$\theta = (\dots)$$

18008969999

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \end{bmatrix}$$

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

$$A \in \mathbb{R}^{2 \times 2} \quad A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

$$f(A) : \mathbb{R}^{2 \times 2} \rightarrow \mathbb{R}$$

$$f(A) = A_{11} + A_{12}^2$$

$$f\left(\begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix}\right) = 5 + 6^2 = 41$$

$$\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \frac{\partial f}{\partial A_{12}} \\ \frac{\partial f}{\partial A_{21}} & \frac{\partial f}{\partial A_{22}} \end{bmatrix}$$

$$\nabla_A f(A) = \begin{bmatrix} 1 & 2A_{12} \\ 0 & 0 \end{bmatrix}$$

$$\nabla_\theta J(\theta) \stackrel{\text{set}}{=} \vec{0}$$

If A is a square $(A \in \mathbb{R}^{n \times n})$

$$\underbrace{\text{trace of } A}_{\text{sum of diagonal entries}} \quad \boxed{\text{tr } A = \text{tr } A^T}$$

$$f(AB) = \text{tr } AB$$

$$\nabla_A f(A) = B^T$$

$$\boxed{\text{tr } AB = \text{tr } BA}$$

$$\boxed{\text{tr } ABC = \text{tr } CAB}$$

$$\nabla_A \text{tr } A A^T C = CA + C^T A$$

$$\left\{ \frac{d}{da} a^2 c = 2ac \right\} \text{ similar}$$

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

$$X\theta = \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \\ \vdots \\ -(x^{(m)})^T \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$$

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

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad J(\theta) = \frac{1}{2} (X\theta - y)^T (X\theta - y)$$

$$X\theta - y = \begin{bmatrix} h_\theta(x^{(1)}) - y^{(1)} \\ h_\theta(x^{(2)}) - y^{(2)} \\ \vdots \\ h_\theta(x^{(m)}) - y^{(m)} \end{bmatrix}$$

$$\boxed{z^T z = \sum_i z_i^2}$$

$$(X\theta - y)^T (X\theta - y)$$

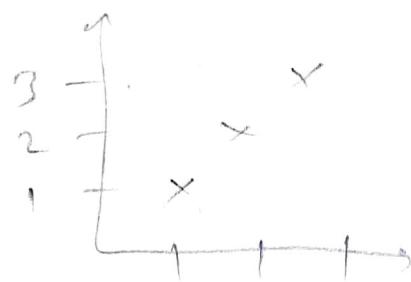
$$\begin{aligned}
 \nabla_{\theta} J(\theta) &= \nabla_{\theta} \frac{1}{2} (x^T \theta - y)^T (x^T \theta - y) \\
 &= \frac{1}{2} \nabla_{\theta} (\theta^T x^T - y^T)(x^T \theta - y) \\
 &= \frac{1}{2} \nabla_{\theta} (\theta^T x^T x \theta - y^T x^T \theta \\
 &\quad - y^T x \theta + y^T y) \\
 &\leftarrow \sum \{(ax-b)(ax-b) = a^2x^2 - axb -\right. \\
 &\quad \left. - bax + b^2\}\} \\
 &= \frac{1}{2} [x^T x \theta + x^T x \theta - x^T y \\
 &\quad - x^T y] \\
 &= x^T x \theta - x^T y \stackrel{\text{set}}{=} \vec{0}
 \end{aligned}$$

~~$x^T x \theta = x^T y$~~

NORMAL EQU.

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

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



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

i.e. $\theta_0 = 0 \Rightarrow h_{\theta}(x) = \theta_1 x$

$$J(\theta) = \frac{1}{2 \times 3} (0-1)^2 + (0-2)^2 + (0-3)^2 \\ = \frac{1}{6} (1+4+9) = \frac{7}{3}$$

$$J(\theta_i) = \frac{1}{2m} \sum_{i=1}^m (\theta_i x_i - y^{(i)})^2$$

gradient descent

$$\min J(\theta_0, \theta_1)$$

θ_0, θ_1

- ① Start at $\theta_0 = 0, \theta_1 = 0$ (it can be any value though)
- ② Keep changing θ_0 & θ_1 to reduce $J(\theta_0, \theta_1)$
- ③ until we reach minimum

$$\theta_j := \theta_j + \sqrt{\theta_0 \theta_1} \quad \theta_0 = 1, \theta_1 = 2$$

$$\theta_0 := \theta_0 + \sqrt{\theta_0 \theta_1} = 1 + \sqrt{1 \times 2} = 1 + \sqrt{2}$$

$$\theta_1 := \theta_0 + \sqrt{\theta_0 \theta_1} = 2 + \sqrt{1 \times 2} = 2 + \sqrt{2}$$

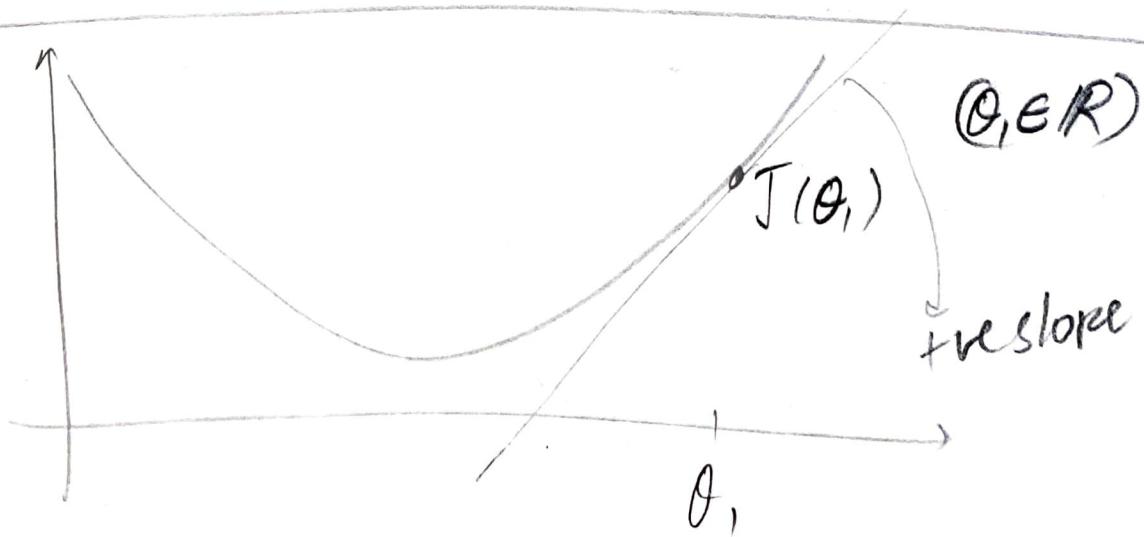
Correct Simultaneously update:

$$\text{temp } 0 = \theta_0 - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0}$$

$$\text{temp } 1 = \theta_1 - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1}$$

$$\theta_0 = \text{temp } 0$$

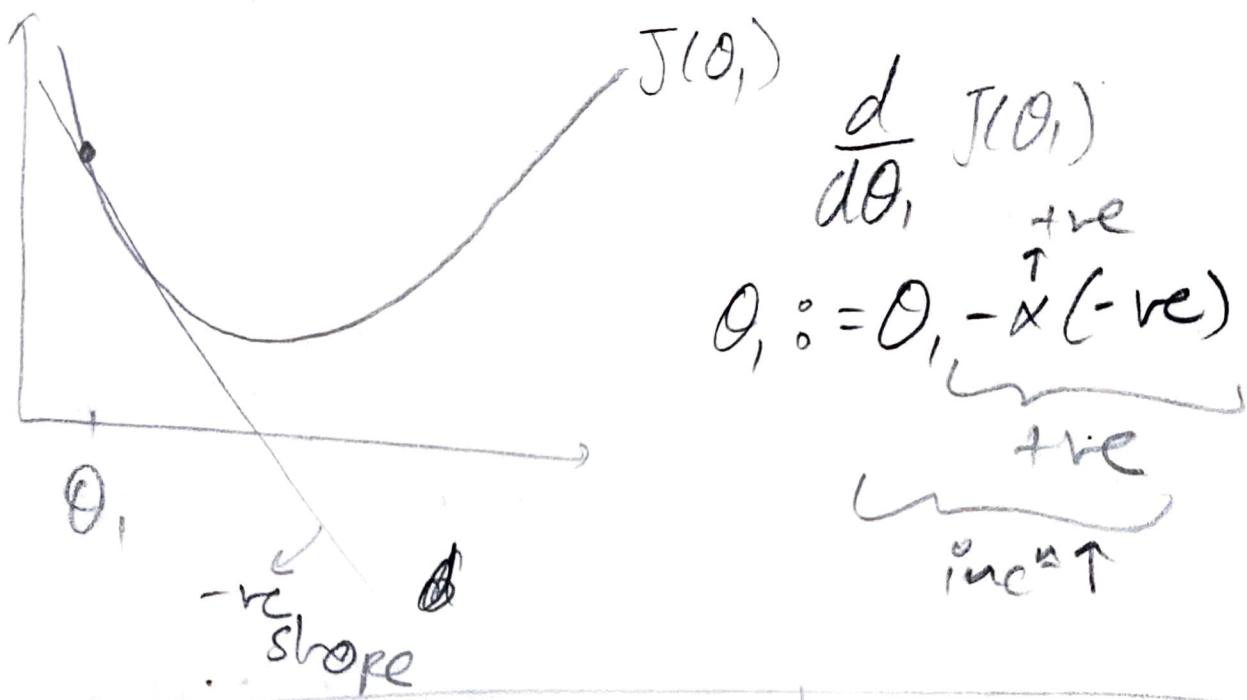
$$\theta_1 = \text{temp } 1$$



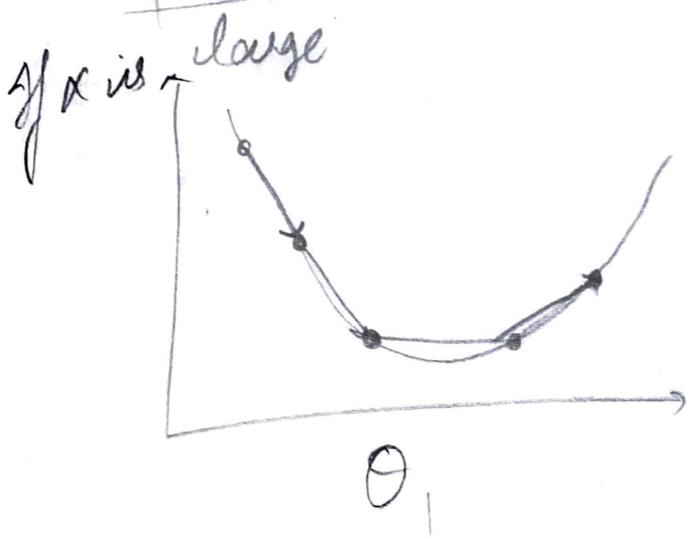
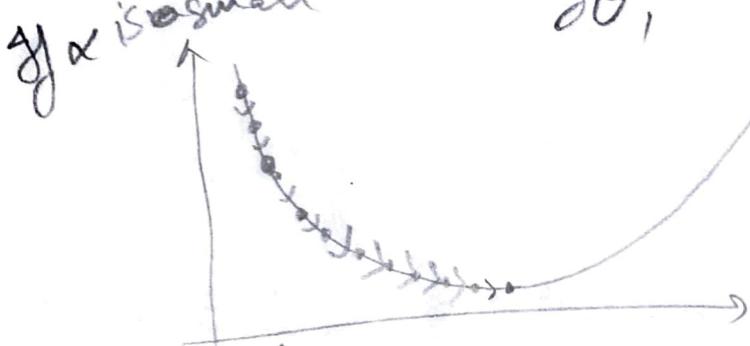
$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} (J(\theta_1))$$

$$\theta_1 := \theta_1 - \alpha (+ve)$$

↓
 ↗ $\theta_1 (+ve)$ ↗ $-ve$
 decⁿ ↓



$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$



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

$$\frac{1}{m} \cdot \cancel{\frac{\partial}{\partial \theta_j} (h_\theta(x^{(i)}) - y^{(i)})} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot \cancel{\frac{\partial}{\partial \theta_j} h_\theta(x^{(i)})} = \sum_{i=1}^m (\theta_0 + \theta_1 x_1^{(i)} - y^{(i)})^2$$

$$\text{for } j=0, \theta_0 \Rightarrow \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m h_\theta(x^{(i)}) - y^{(i)}$$

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

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

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

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

$$= \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \frac{\partial}{\partial \theta_j} (\theta_0 + \theta_1 x_1^{(i)} - y^{(i)})$$

$$\text{for } j=0, \frac{\partial}{\partial \theta_0} J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \frac{\partial}{\partial \theta_0} (\theta_0 + \theta_1 x_1^{(i)} - y^{(i)})$$

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

$$\text{for } j=1, \frac{\partial}{\partial \theta_1} J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

Gradient descent for linear regression
repeat until convergence Σ

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

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

(Update θ_0 & θ_1 simultaneously)

$$h_\theta(x^{(0)}) = \theta_0 + \theta_1 x^{(0)}$$

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

x	y
3	2
1	2
0	1
4	3

$$J(0, 1) = ? = \frac{1}{8} (1 + 1 + 1 + 1)$$
$$= \frac{1}{2} [532]^2$$

$$h_\theta(6) = \theta_0 + \theta_1 n$$

$$2 \quad 6 \quad \begin{bmatrix} -4 \\ 0 \end{bmatrix} = -2 + \frac{1}{2} \times 6 \begin{bmatrix} 2 & 2 & 2 \\ 2 & 2 & 2 \end{bmatrix}$$
$$-3 \quad -3 \quad \begin{bmatrix} 1 \\ 1/2 \end{bmatrix} = 1 \quad 2 \cancel{\begin{pmatrix} 4-4 \\ 4-4 \\ 4-4 \end{pmatrix}} + 2$$

$$11 \quad 16 \quad \begin{bmatrix} 4 \\ 3 \\ 12 \\ 6 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 2 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} 2 \\ 12 \\ 31 \end{bmatrix} - \begin{bmatrix} 3 \\ 12 \\ 11 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 2 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 3 \\ 4 & 0 \\ 2 & 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 5 \end{bmatrix} = \begin{array}{l} ② \times 1 \\ 3 \times 2 \end{array} \begin{bmatrix} 526 \\ 1052 \\ 2104 \end{bmatrix} \begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix} \begin{bmatrix} -40 \\ 0.25 \end{bmatrix} \begin{bmatrix} 486 \\ 314 \\ 384.5 \\ 173 \end{bmatrix}$$

$$\begin{bmatrix} 16 \\ 4 \\ 7 \end{bmatrix} \begin{bmatrix} 526 \\ -40 \\ 486 \end{bmatrix} \begin{bmatrix} 267 \\ 383.5 \end{bmatrix} \begin{bmatrix} 708 \\ 354 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 1 & 5 & 173 \\ 0 & 3 & 0 & 4 & 1 \\ -1 & -2 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \\ 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 15 \\ 13 \\ -7 \end{bmatrix}$$

$$1+6+2+5 \\ 0+9+0+4 \\ -1-6 \quad 0 \quad 0$$

$$y = ax + b$$

~~$y = xp$~~

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

$$X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} P = \begin{bmatrix} a \\ b \end{bmatrix}$$

~~$y_0(n) = -40 + 0.25n$~~

~~$y = 40 + \frac{n}{4}$~~

~~$\frac{n}{4} = y - 40$~~

~~$X = \begin{bmatrix} x \end{bmatrix} \quad B = \begin{bmatrix} 40 \end{bmatrix}$~~

$$\begin{bmatrix} 1 & 3 & 2 \\ 4 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 11 \\ 9 \\ 14 \end{bmatrix}$$

$$1+0+10=11 \quad 3+3+4$$

$$4+0+5=9 \quad 12+0+2$$

$$\begin{bmatrix} 1 & 3 \\ 2 & 5 \end{bmatrix} \begin{bmatrix} 0 & 7 \\ 3 & 2 \end{bmatrix} = \begin{bmatrix} 9 & 17 \\ 15 & 12 \end{bmatrix}$$

$0+9$ ~~2+1~~ $2+10$

$0+15$ $1+6$

$$\rightarrow \begin{bmatrix} 1 & 3 \\ 2 & 5 \end{bmatrix} \begin{bmatrix} 0 \\ 3 \end{bmatrix} = \begin{bmatrix} 9 \\ 15 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 3 \\ 2 & 5 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 7 \\ 12 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 3 \\ 2 & 5 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 2 & 3 \end{bmatrix} = \begin{bmatrix} ? & 9 \\ a & b \\ c & d \end{bmatrix}$$

$$\begin{bmatrix} 1 & 3 \\ 2 & 5 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 7 \\ 10 \\ 10 \end{bmatrix}$$

$$\left. \begin{array}{l} a=10 \\ b=12 \\ c=10 \\ d=15 \end{array} \right\}$$

$$\begin{bmatrix} 1 & 3 \\ 2 & 5 \\ 6 & 5 \end{bmatrix} \begin{bmatrix} 0 \\ 3 \end{bmatrix} = \begin{bmatrix} 9 \\ 12 \\ 15 \end{bmatrix}$$

$2+8$ $6+10$

$0+9$ $0+12$

House Sizes	
2104	
1416	
1534	
852	

- 3 competing hypothesis
- $$h_0(n) = -40 + 0.25n$$
- $$h_0(x) = 200 + 0.1x$$
- $$h_0(n) = -150 + 0.4n$$

$$\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix} \times \begin{bmatrix} -40 & 200 & -150 \\ 0.25 & 0.1 & 0.4 \end{bmatrix}$$

$$= \begin{bmatrix} 486 & 410 & 692 \\ 314 & 342 & 416 \\ 344 & 353 & 464 \\ 173 & 285 & 191 \end{bmatrix}$$

identity matrix

informally,

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

$$A \cdot I = I \cdot A = A$$

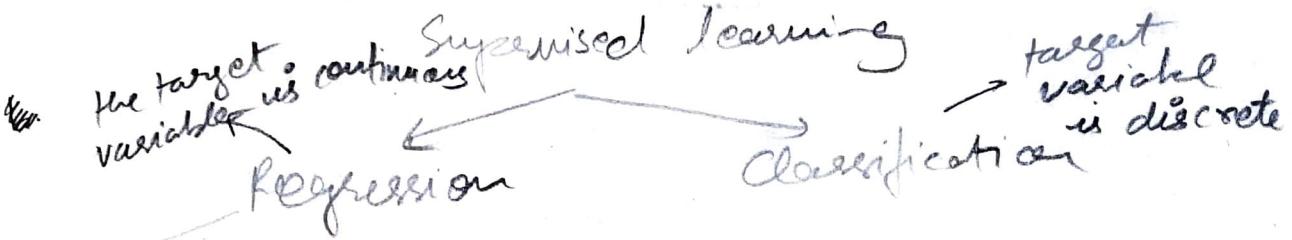
$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$
nxn nxn nxn nxn

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 3 \\ 2 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \\ 2 \end{bmatrix}$$

$$1+0+0=1$$

$$0+3+0=3$$

$$0+0+2=2$$



① Linear Regression

WORKFLOW



$h(n)$ function

$n = \# \text{ features}$ ($n=2$)

$m = \# \text{ training examples}$

$$\text{Ex} = (m=4)$$

In general,

$$h(n) = \sum_{j=0}^n \theta_j x_j$$

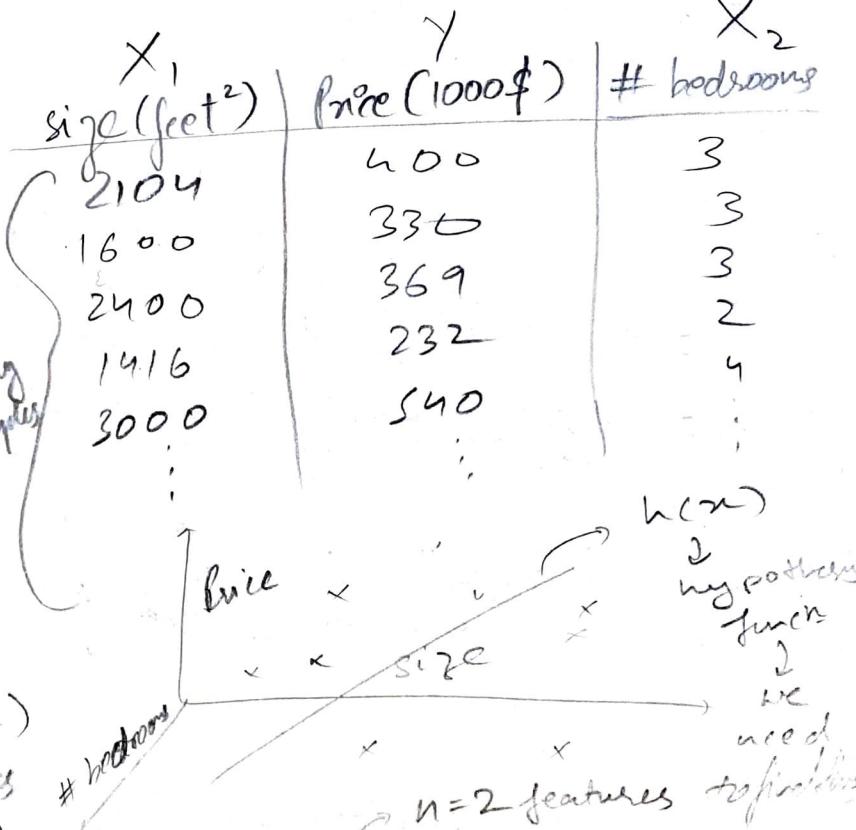
where, $\theta_0 = 1$ [for n features]

and for n features

there are $n+1$ parameters

$$\text{Also, } \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \text{ & } x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

for n features
there are $n+1$
parameters (θ)



$$h(n) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$\theta_0, \theta_1, \theta_2 \rightarrow 3 \text{ parameters (or, weights)}$

$x_1, x_2, \text{size} \rightarrow \text{inputs (or, features)}$

$x_1 = \text{size (in feet}^2\text{)}$

$x_2 = \# \text{ bedrooms}$

$$h(\theta) = \sum_{j=0}^m \theta_j x_j = \theta^T x \quad \text{where } \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}$$

$(x^{(i)}, y^{(i)})$ is training example

 $x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_m \end{bmatrix}$

$$h_\theta(x) = h(x)$$

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

~~J'(theta) = 0~~

(cost func.)

Task \Rightarrow choose θ s.t. $J(\theta)$ is minimum.

$$\frac{d\sigma^2}{d\sigma} = 2n$$

① LMS Algo (Gradient descent)

Starts with θ (say $\theta = \vec{\theta}_0$)

$$2n = 0$$

$$\sigma = 0$$

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

assignment operator ~~using learning rate~~
~~(size of step towards a minimum)~~

$$\theta - 0.01 \times 20$$

$$\theta - 0.020$$

$$1.98 \theta - 0.001$$

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta_j} &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_\theta(x) - y)^2 \\ &= \frac{1}{2} \cdot x (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left(\sum_{i=0}^m \theta_i x_i - y \right) \\ &= (h_\theta(x) - y) \cdot (x_j) \end{aligned}$$

$$\boxed{\theta_j := \theta_j - \alpha (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}}$$

Repeat this until convergence

(for $j = 0, 1, \dots, n$)

$$\theta_j := \theta_j - \alpha \sum_{i=0}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

{

Maths

Stochastic gradient descent.

Loop \mathcal{E}

for $i = 1$ to $i = m$ \mathcal{E}

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

(for $j = 0 \rightarrow n$)

\mathcal{E} Useful for large datasets
(i.e. large value of m)

Matrix derivatives

$f: \mathbb{R}^{mn} \rightarrow \mathbb{R}$
gradient of $f(A)$ w.r.t. A

$$\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \frac{\partial f}{\partial A_{12}} & \cdots & \frac{\partial f}{\partial A_{1n}} \\ \frac{\partial f}{\partial A_{21}} & \frac{\partial f}{\partial A_{22}} & \cdots & \frac{\partial f}{\partial A_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial A_{m1}} & \frac{\partial f}{\partial A_{m2}} & \cdots & \frac{\partial f}{\partial A_{mn}} \end{bmatrix}_{mn \times n}$$

\downarrow $\rightarrow (i, j)^{\text{th}}$ element
derivative of f w.r.t. A

If A is a 2×2 matrix $\Rightarrow A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$

and func $f: \mathbb{R}^{2 \times 2} \rightarrow \mathbb{R}$

$$f(A) = \frac{3}{2} A_{11} + 5A_{12}^2 + A_{21}A_{22}$$

$$\nabla_A f(A) = \begin{bmatrix} \frac{3}{2} & 10A_{12} \\ A_{22} & A_{21} \end{bmatrix}$$

Trace of a matrix

$A \rightarrow$ ~~mn~~ matrix fully

$$\text{tr } A = \sum_{i=0}^n A_{ii} = A_{11} + A_{22} + \cdots + A_{nn}$$
 (only for square matrix)

* If $a \in \mathbb{R}$ i.e. a is 1×1 matrix
 then, $\text{tr}(a) = a$ prop ①

If A & B are two matrices s.t. $\text{tr } BA = \text{tr } AB$

$$\text{tr } ABC = \text{tr } CAB = \text{tr } BCA$$

$$\text{tr } ABCD = \text{tr } DABC = \text{tr } CDAB = \text{tr } BCDA$$

$$\text{tr } A = \text{tr } A^T$$

$$\text{tr } (A + B) = \text{tr } A + \text{tr } B$$

$$\text{tr } aA = a \text{tr } A$$

$$\nabla_A \text{tr } AB = B^T$$

$$\nabla_A \text{tr } f(A) = (\nabla_A f(A))^T$$

$$\nabla_A \text{tr } ABATC = CAB + C^T AB^T$$

$$\nabla_A |\text{tr } A| = |\text{tr } A| (A^{-1})^T$$

$$x = \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \\ \vdots \\ -(x^{(m)})^T \end{bmatrix} \quad \vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

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

$$x\theta - \vec{y} = \begin{bmatrix} (x^{(1)})^T \theta \\ \vdots \\ (x^{(m)})^T \theta \end{bmatrix} - \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$= \begin{bmatrix} h_\theta(x^{(1)}) - y^{(1)} \\ h_\theta(x^{(m)}) - y^{(m)} \end{bmatrix}$$

$$\text{由 } z^T z = \sum_i z_i^2$$

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

$$= J(\theta)$$

$$\therefore \nabla_{A^T} \text{tr} A B A^T C = B^T A A^T C + B A^T C$$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \frac{1}{2} (x\theta - \vec{y})^T (x\theta - \vec{y})$$

$$= \frac{1}{2} \nabla_{\theta} (\theta^T x^T - \vec{y}^T) (x\theta - \vec{y})$$

$$= \frac{1}{2} \nabla_{\theta} (\theta^T x^T x \theta - \theta^T x^T y - \vec{y}^T x \theta + \vec{y}^T \vec{y})$$

$$= \frac{1}{2} \nabla_{\theta} \text{tr} (\theta^T x^T x \theta - \theta^T x^T y - \vec{y}^T x \theta + \vec{y}^T \vec{y})$$

$$= \frac{1}{2} \nabla_{\theta} (\text{tr } \theta^T x^T x \theta - 2 \text{tr } \vec{y}^T x \theta)$$

$$= \frac{1}{2} (x^T x \theta + x^T x \theta - 2 x^T \vec{y})$$

$$= x^T x \theta - x^T \vec{y}$$

$$\nabla_{\theta} J(\theta) \stackrel{\text{set}}{=} \vec{0} \text{ for } \theta \text{ s.t. } J(\theta) \text{ min min}^u$$

$$\nabla_{\theta} J(\theta) = x^T x \theta - x^T \vec{y} = \vec{0}$$

$$x^T x \theta - x^T \vec{y} = \vec{0}$$

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

$$3(3^{-1}) = 1 \quad 12(12^{-1}) = 1 \quad 0(0^{-1}) = 1 \times \text{WRONG} \quad (\text{undefined})$$

only for non-zero square matrices

$\underset{\substack{\downarrow \\ \text{(inverse)}}}{A A^{-1}} = I \quad \bullet = A^{-1} A$

(A matrix) of A matrix

$A = \begin{bmatrix} 3 & 4 \\ 2 & 16 \end{bmatrix}_{2 \times 2} \quad A^{-1} = \begin{bmatrix} 0.04 & -0.1 \\ 0.05 & 0.075 \end{bmatrix}_{2 \times 2} \quad AA^{-1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}_{2 \times 2}$

$$A^T = \begin{bmatrix} 1 & 2 & 0 \\ 2 & 3 & 5 \\ 0 & 1 & 9 \end{bmatrix}_{m \times n}^{2 \times 3} \xrightarrow{\text{Transpose}} A^T = \begin{bmatrix} 1 & 2 \\ 3 & 5 \\ 0 & 9 \end{bmatrix}_{n \times m}^{3 \times 2} = B$$

$$A_{12} = 2 \quad \rightarrow \quad A_{ij}^o = B_{ji}^o \quad B_{32} = 9$$

$$B_{21} = 2$$

$$A_{23} = 9$$

$$\begin{bmatrix} 1 & 2 \\ 0 & 3 \\ 1 & 4 \end{bmatrix}^T = \begin{bmatrix} 1 & 2 \\ 0 & 3 \\ 1 & 4 \end{bmatrix}$$

$$A = \begin{bmatrix} 1 & 0 & 1 \\ 2 & -2 & -1 \\ 3 & 0 & 0 \end{bmatrix}$$

$$A^{-1} = \frac{1}{|A|} \text{adj}^o(A)$$

$$\text{adj}^o(A) = \begin{bmatrix} +0 & -3 & +6 \\ -0 & +3 & -0 \\ +2 & -3 & +2 \end{bmatrix}$$

$$\text{adj}^o(A) = \begin{bmatrix} 0 & -3 & 6 \\ 0 & -3 & 0 \\ 2 & 3 & -2 \end{bmatrix}^T = \begin{bmatrix} 0 & 0 & 2 \\ -3 & -3 & 3 \\ 6 & 0 & -2 \end{bmatrix}$$

$$|A| = 1(0 - 0) - 0() + 1(0 + 6) = 6$$

$$A^{-1} = \begin{bmatrix} 0 & 0 & 1/3 \\ -1/2 & -1/2 & 1/2 \\ 1 & 0 & -1/3 \end{bmatrix}$$

~~$$A = \begin{bmatrix} 0 & -1/2 & 1 \\ 0 & 1/2 & 0 \\ 1/3 & 1/2 & 1/3 \end{bmatrix}$$~~

in reverse

Linear Regression via multivariable

size (feet ²) x_1	# bedrooms x_2	# floors x_3	Age x_4	Price (\$1000)
2104	5	1	45	460
1416	3	2	40	232
1537	3	2	30	315
852	2	1	36	178
:	:	:		

this data (training set) has 4 features

$$n = \# \text{ features}$$

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

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

$$x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix} \in \mathbb{R}^n \quad x_3^{(2)} = 2 \quad x_1^{(n)} = 852$$

$$h_0(x) = \theta_0 + \underbrace{\theta_1 x_1}_{\text{size}} + \underbrace{\theta_2 x_2}_{\text{#bedrooms}} + \underbrace{\theta_3 x_3}_{\text{#floors}} + \underbrace{\theta_4 x_4}_{\text{Age}}$$

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

$$h_0(x) = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n \quad (\text{where, } x_0^{(i)} = 1)$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \Rightarrow \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \Rightarrow \mathbb{R}^{n+1}$$

$\boxed{x_0 = 1}$

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

$$\boldsymbol{\theta}^T = [\theta_0 \ \theta_1 \ \dots \ \theta_n]_{1 \times (n+1)}$$

row matrix vector

$$\boldsymbol{\theta}^T \mathbf{x} = [\theta_0 \ \theta_1 \ \dots \ \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$$= \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

\Rightarrow a.k.a multivariate linear regression.

$$\therefore h_{\boldsymbol{\theta}}(\mathbf{x}) = \sum_{i=0}^n \theta_i x_i = \boldsymbol{\theta}^T \mathbf{x}$$

where, $x_0 = 1$

$$\begin{aligned} J(\theta_0, \theta_1, \dots, \theta_n) &= J(\boldsymbol{\theta}) = \frac{1}{2m} \sum_{i=1}^m (h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{2m} \sum_{i=1}^m \left(\sum_{j=0}^n \theta_j x_j^{(i)} - y^{(i)} \right)^2 \\ &\text{where, } x_0 = 1 \\ &= \frac{1}{2m} \sum_{i=1}^m (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)})^2 \end{aligned}$$

Gradient descent?

Repeat until convergence

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) \quad (\text{for every } j = 0, 1, \dots, n)$$

{ }

{ simultaneous update }

$$\begin{aligned}
 \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \left(\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \right) \\
 &= \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \frac{\partial}{\partial \theta_j} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \\
 &= \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \times \frac{\partial}{\partial \theta_j} \left(\sum_{j=1}^n (\theta_j x_j^{(i)}) - y^{(i)} \right) \\
 &= \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} = x_j^{(i)}
 \end{aligned}$$

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

gradient descent for $n \geq 1$:

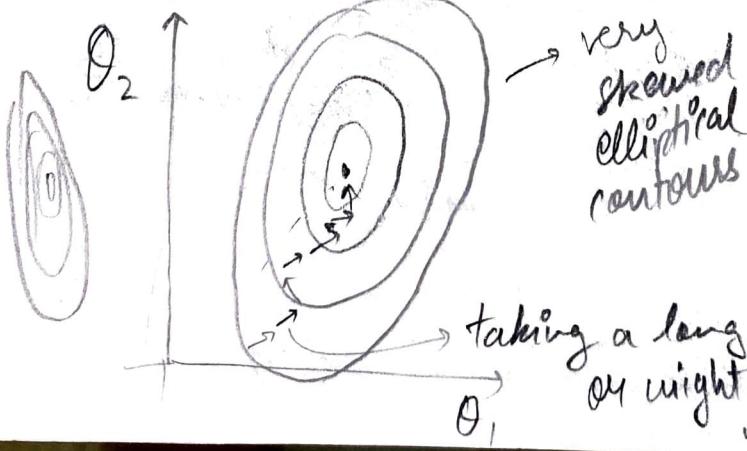
Repeat until convergence ε

$$\begin{aligned}
 \theta_j &:= \theta_j - \alpha \left(\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \right) \cdot x_j^{(i)} \\
 \} \quad &(\text{for every } j = 0, 1, 2, \dots, n) \\
 &[\text{simultaneous update}]
 \end{aligned}$$

Feature Scaling \Rightarrow Idea: Make sure features are on a similar scale.

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

x_2 = # bedrooms ($1-5$)

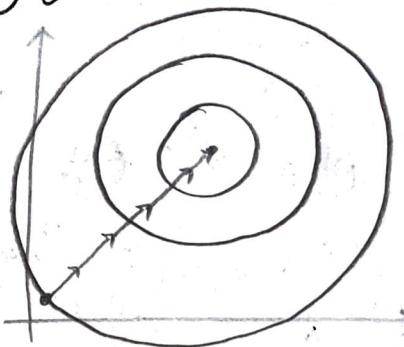


but scaling the feature properly by multiplying or dividing the inputs by an appropriate.

$$x_1 = \text{size (feet}^2) ; x_2 = \frac{\# \text{bedrooms}}{5} \quad 0 \leq x_1 \leq 1 \quad 0 \leq x_2 \leq 1$$

2000

θ_2



$J(\theta)$

→ contour of funcⁿ
is less skewed
∴ reaching to
minima is easier

θ_1

Get every feature into approx. range. $-1 \leq x_i \leq 1$ where, $x_0 = 1$

$$\text{eg.) } 0 \leq x_1 \leq 3 \checkmark ; -2 \leq x_2 \leq 0.5 \checkmark ; -100 \leq x_3 \leq 100 \checkmark \\ ; -0.0001 \leq x_4 \leq 0.0001 X$$

Mean Normalization

replace x_i with $x_i - \bar{x}_i$ to make features have approximately zero mean (Don't apply $x_0 = 1$)

avg. size.

avg # bedroom

$$\text{eg.) } x_1 = \frac{\text{size} - 1000}{2000} ; x_2 = \frac{\# \text{bedrooms} - 2}{5}$$

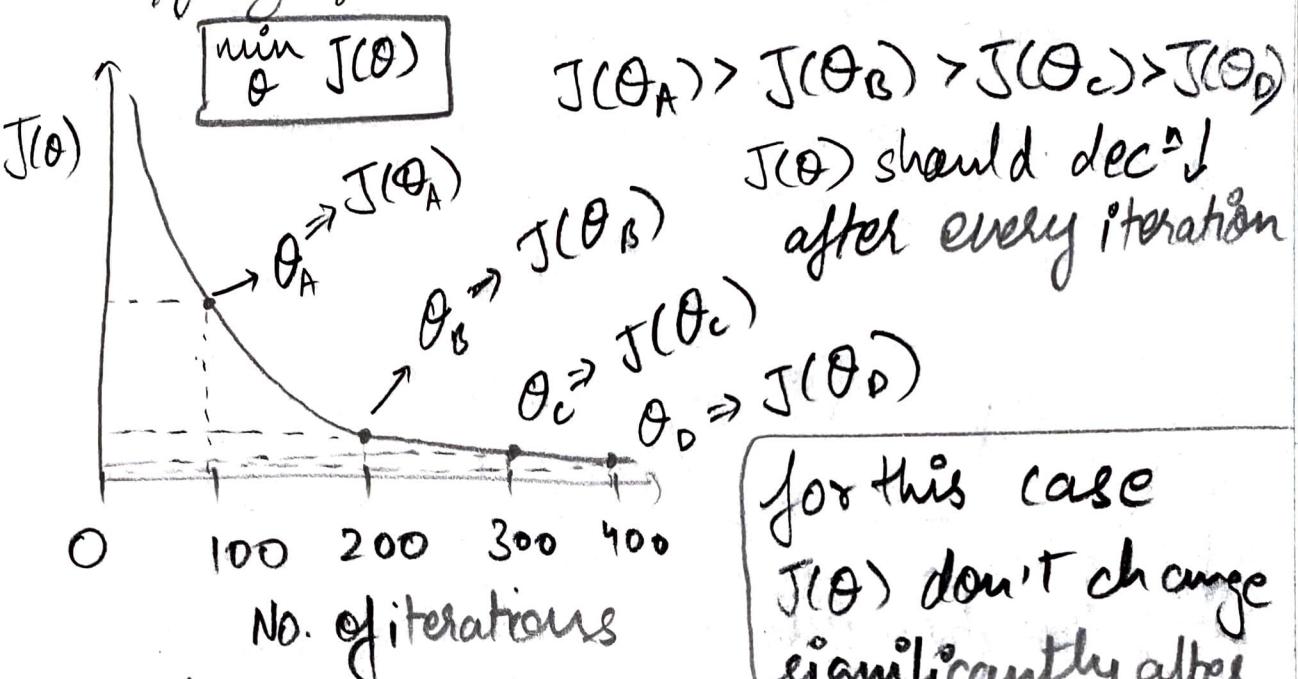
$$\text{generally, } -0.5 \leq x_1 \leq 0.5 ; -0.5 \leq x_2 \leq 0.5$$

$$x_1 \rightarrow \frac{x_1 - \bar{x}_1}{s_1} \quad \begin{matrix} \text{avg. value of } x_1 \text{ in} \\ \text{training set} \end{matrix}$$

$$\begin{matrix} \text{for } x_2 \\ x_2 \rightarrow \frac{x_2 - \bar{x}_2}{s_2} \end{matrix} \quad \begin{matrix} \text{range (max-min)} \\ \text{std deviation of } x_2 \end{matrix}$$

$\tilde{x}_i = \frac{\text{Age} - \text{avg Age}}{\text{range of Age}}$ | feature scaling
 ↓ can make
 input features 20 avg Age gradient descent
 range of Age much faster

Debunking Gradient descent.



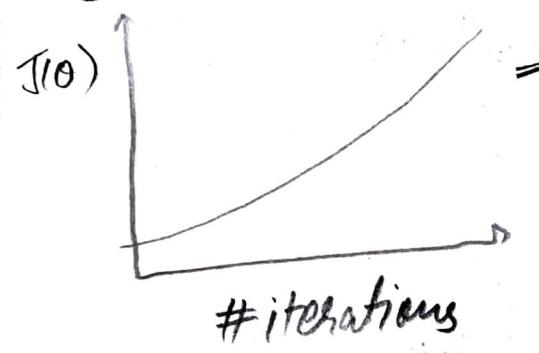
No. of iterations for gradient descent to get $\min_{\theta} J(\theta)$ can vary significantly

for this case
 $J(\theta)$ don't change significantly after 300 iterations

→ Eg of a automatic convergence test:
 Declare convergence if $J(\theta)$ is less than 10^{-3} in one iterations

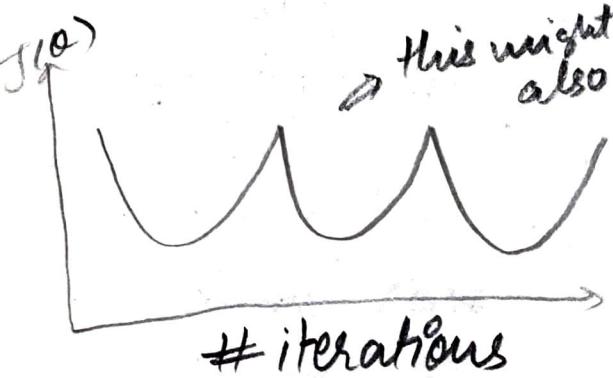
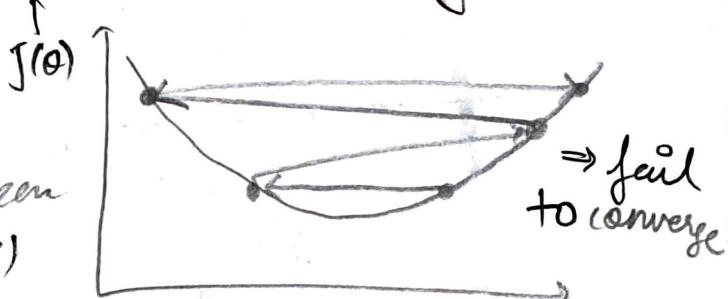
→ this threshold is not definitive and plotting $J(\theta)$ vs #iteration graph is better to declare convergence.

J(θ) v/s # iteration looks like



→ this means α (i.e. learning rate is too large) and decⁿ the value of α

this might also happen
the code have bug(s)



→ this might also happen when α is large $\theta \rightarrow$
if α is large

- for sufficiently small α , $J(\theta)$ should decrease on every iteration (theorem)
- but if α is too small then, can be slow to converge

Summary,

- α is too small \Rightarrow slow convergence
 - α is too large \Rightarrow $J(\theta)$ may not converge with every iteration as it may not decrease or convergence may happen but slowly.
- $x^3 \rightarrow 0.003 \rightarrow x^3 \rightarrow 0.03 \rightarrow \dots$ try to choose like this in practice
($\alpha = \dots, 0.001, 0.01, 0.1, \dots, 1 - \text{so on}$)

Housing Prices Prediction

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

by linear regression



new feature frontage

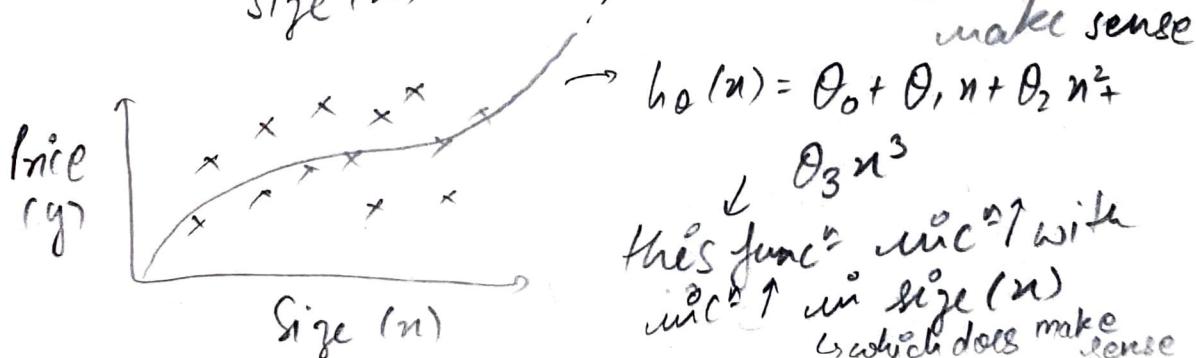
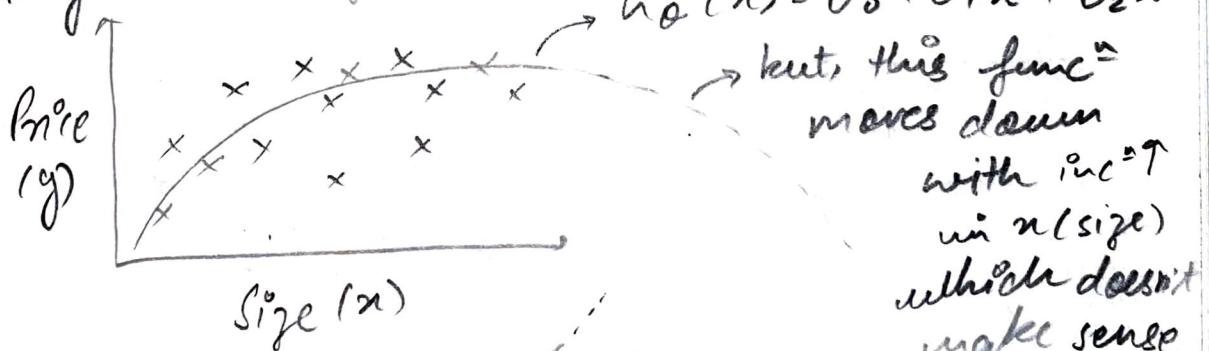
$$\text{Area} = x = \text{frontage} \times \text{depth}$$

$$h_0(x) = \theta_0 + \theta_1 x = \theta_0 + \theta_1 \times \text{Area}$$

→ this new feature formed using two given features of training data set

∴ only two parameters instead of three

Polynomial Regression



$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

$$= \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3$$

$x_1 = (\text{size})$; $x_2 = (\text{size})^2$; $x_3 = (\text{size})^3$

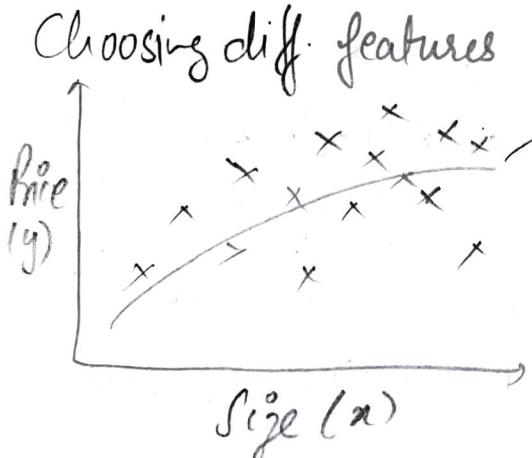
③

↳ 3 different features
using single given feature
(size $\rightarrow x$)

POLYNOMIAL REGRESSION

If, size: $1 \rightarrow 10^3$ size $^3 = 1 \rightarrow 10^9$
then, size $^2 = 1 \rightarrow 10^6$

↳ quite large range
use feature scaling to improve gradient descent.



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

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

↳ using different set of features for same set of training data

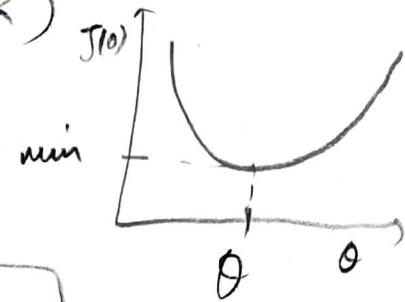
NORMAL EQUATION: Method to solve for θ analytically

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

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

$$\frac{d}{d\theta} J(\theta) = 2a\theta + b \stackrel{\text{set}}{=} 0$$

$$\text{Solve for, } \theta = -\frac{b}{2a}$$



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

$$\frac{\partial}{\partial \theta_j} J(\theta) \stackrel{\text{set}}{=} 0 \quad (\text{for every } j)$$

then solve for $\theta_0, \theta_1, \theta_2, \dots, \theta_n$

Eg) $m=4$

No	size (feet ²)	# bedrooms	# floors	Age	Price
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}_{m \times (n+1)} \quad y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}_m$$

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

normal eqⁿ

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

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

(design matrix)

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

Eg) If $x^{(i)} = \begin{bmatrix} 1 \\ x_i^{(i)} \end{bmatrix}$ here, $n=1$

$$X = \begin{bmatrix} 1 & x_1^{(1)} \\ 1 & x_1^{(2)} \\ \vdots & \vdots \\ 1 & x_1^{(m)} \end{bmatrix}_{m \times 2} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}_m$$

$$\theta = (X^T X)^{-1} X^T y \quad \{(X^T X)^{-1}\} \text{ is inverse of } X^T X$$

Let $A = X^T X$, $A^{-1} = (X^T X)^{-1}$

Octave: piuv $\underbrace{(X^T X)^{-1} X^T y}_\theta$ not need
for feature
scalin

$$\theta = (X^T X)^{-1} X^T y \Rightarrow \text{s.t. } \min_{\theta} J(\theta)$$

let m training examples, n features

Gradient Descent

- ⇒ Need to choose ' α '
- ⇒ Needs many iterations
- ⇒ works well even when n is large

If $n \approx 10^6$ or larger,

** If $X \in \mathbb{R}^{m \times n+1}$ then $X^T \in \mathbb{R}^{n+1 \times m}$

Set, $X^T X = A \in \mathbb{R}^{n+1 \times n+1}$ $O(n^3) \Rightarrow$ time complexity of inverting a matrix
 \downarrow
 $\begin{matrix} n \times m \\ m \times n \end{matrix} \xrightarrow{\text{same}}$ Square matrix \therefore , for large value of n inverting a matrix is computationally expensive.

What if $(X^T X)$ is non-invertible? (i.e. singular matrix/degenerate matrix)

Octave: $\text{pinv}(X^T X)^{-1} X^T y$

pseudo inverse \rightarrow can solve

If $X^T X$ is non-invertible
→ Redundant features exist (i.e. linearly dependent)

e.g.) $x_1 = \text{size(feet}^2\text{)}; x_2 = \text{size(m}^2\text{)}$

Normal Equation

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

$n = 100, 1000, 10,000$
 $\downarrow \quad \downarrow \quad \downarrow$
fast reasonable \Rightarrow might be slow

$x_1 = kx_2$ (linearly dependent)

$$\hookrightarrow x_1 = (3 \cdot 28)^2 x_2$$

\rightarrow If $x_1 = kx_2$ then $X^T X$ is non-invertible

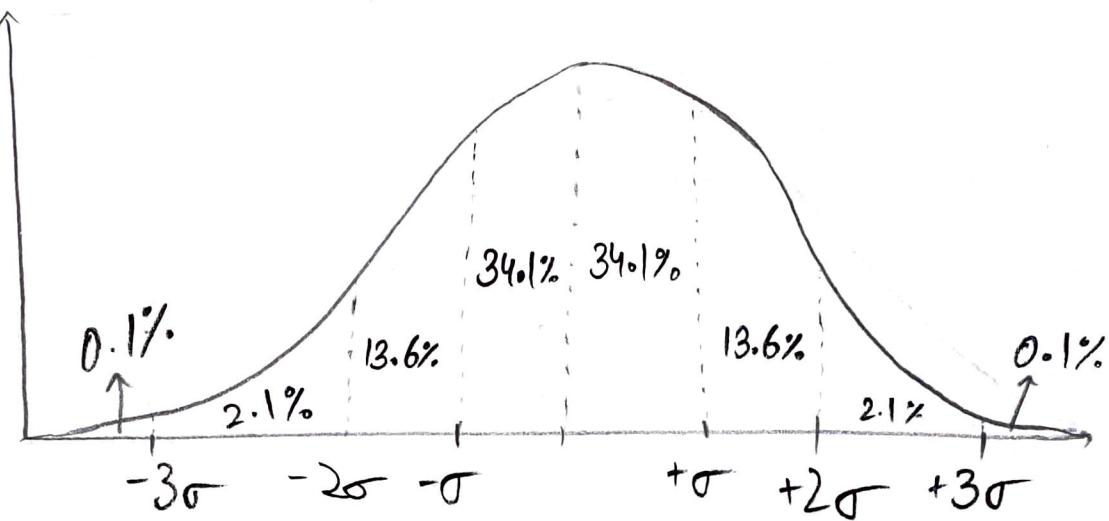
\Rightarrow Too many features (Eg $m \leq n$) \Rightarrow $\begin{matrix} m=10 \\ n=101 \end{matrix}$

\hookrightarrow Delete some features

sometimes work
but, not always \Rightarrow there is too little data.

$$0 \in R^{101}$$

Standard deviation



$$\sigma = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$$

$$\text{Data} = 1, 2, 3, 4, 5$$

$$\bar{x} = \frac{1+2+3+4+5}{5} = \boxed{3}$$

$$\sum (x - \bar{x})^2$$

$$x=1 \Rightarrow (1-3)^2 = 4 \quad \left| \begin{array}{l} x=3 \Rightarrow (3-3)^2 = 0 \\ x=4 \Rightarrow (4-3)^2 = 1 \end{array} \right| \quad \left| \begin{array}{l} x=5 \Rightarrow (5-3)^2 = 4 \end{array} \right.$$

$$\sum (x - \bar{x})^2 = 4 + 1 + 0 + 1 + 4 = 10$$

$$\sigma = \sqrt{\frac{10}{5-1}} = \sqrt{\frac{10}{4}} = \sqrt{\frac{5}{2}} = 1.58$$

$$\sigma = \sqrt{\frac{(3.6)^2 + (1.6)^2}{4} \times 2 + 0}$$

$$= \sqrt{\frac{(12.96 + 2.56) \times 2 + 0}{4}} = 2.78$$

$$\bar{x} = \frac{0+2+4+5+7}{5}$$

$$= \frac{18}{5} = 3.6$$

Vectorization

$$\textcircled{1} \quad h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j = \theta^T x$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \rightarrow \begin{array}{l} \theta_0 \rightarrow \text{theta}(1) \\ \theta_1 \rightarrow \text{theta}(2) \\ \theta_2 \rightarrow \text{theta}(3) \end{array}$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \rightarrow \begin{array}{l} x_0 \rightarrow x(1) \\ x_1 \rightarrow x(2) \\ x_2 \rightarrow x(3) \end{array}$$

Unvectorized implementation (MATLAB)

prediction = 0.0;

for $j = 1 : n + 1$

prediction = prediction + theta(j) * x(j)

end;

Vectorized Implementation (MATLAB)

prediction = (theta') * x;

Unvectorize Implementation (C++)

double prediction = 0.0;

for (int j = 0; j <= n; j++)

prediction += theta[j] * x[j];

Vectorized Implementation (C++)

double prediction = theta.transpose() * x;

② Vectorization of gradient descent

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad (\text{for all } j=0, 1, \dots, n)$$

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

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

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

(for $j=0, 1, 2$)

Vectorized Implementation:

$$\theta := \theta - \alpha \delta \quad \delta \in \mathbb{R}^{n+1}$$

where, $\delta = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$

$$\delta = \begin{bmatrix} \delta_0 \\ \delta_1 \\ \delta_2 \end{bmatrix} \quad \delta_0 = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$(h_\theta(x^{(1)}) - y^{(1)}) x^{(1)} + (h_\theta(x^{(2)}) - y^{(2)}) x^{(2)} + (h_\theta(x^{(3)}) - y^{(3)}) x^{(3)}$$

$$X = \begin{bmatrix} x_0^{(1)} \\ x_1^{(1)} \\ x_2^{(1)} \\ \vdots \\ x_0^{(3)} \\ x_1^{(3)} \\ x_2^{(3)} \end{bmatrix}$$

$$\forall j \quad u(j) = 2v(j) + 5w(j) \quad (\text{for all } j)$$

$$u = 2V + 5W$$

vectorization

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

(ODE \Rightarrow) $J = \text{cost function } J(X, y, \theta)$

$$m = \text{length}(y)$$

$$\text{Predictions} = X^* \theta \Rightarrow \mathbb{R}^{n+1 \times 1}$$

\downarrow

$$\mathbb{R}^{m \times n+1}$$

$\underbrace{\qquad\qquad\qquad}_{\mathbb{R}^{m \times 1}}$

\downarrow

SqrError = $(\text{Predictions} - y)^2$

$$\text{Predictions} = \begin{bmatrix} h_{\theta}(x^{(1)}) \\ h_{\theta}(x^{(2)}) \\ \vdots \\ h_{\theta}(x^{(m)}) \end{bmatrix} \in \mathbb{R}^{m \times 1}$$

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

$$\text{Prediction} - y = \begin{bmatrix} h_{\theta}(x^{(1)}) - y^{(1)} \\ h_{\theta}(x^{(2)}) - y^{(2)} \\ \vdots \\ h_{\theta}(x^{(m)}) - y^{(m)} \end{bmatrix} \in \mathbb{R}^{m \times 1}$$

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

1 iteration of gradient descent

$$X = \begin{bmatrix} x_0' & x_1' \\ x_0^2 & x_1^2 \\ x_0^3 & x_1^3 \\ \vdots & \vdots \\ x_0^m & x_1^m \end{bmatrix}_{mx2} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}_{2 \times 1} \quad y = \begin{bmatrix} y_1' \\ y_2' \\ y_3' \\ \vdots \\ y^m \end{bmatrix}_{mx1}$$

prediction - y = $\begin{bmatrix} h_\theta(x^{(1)}) - y^{(1)} \\ h_\theta(x^{(2)}) - y^{(2)} \\ h_\theta(x^{(3)}) - y^{(3)} \\ \vdots \\ h_\theta(x^{(m)}) - y^{(m)} \end{bmatrix}_{mx1}$

X^* theta $\begin{bmatrix} x_0' & x_1' \\ x_0^2 & x_1^2 \\ x_0^3 & x_1^3 \\ \vdots & \vdots \\ x_0^m & x_1^m \end{bmatrix}_{2 \times m} \times \begin{bmatrix} h_\theta(x^{(1)}) - y^{(1)} \\ h_\theta(x^{(2)}) - y^{(2)} \\ h_\theta(x^{(3)}) - y^{(3)} \\ \vdots \\ h_\theta(x^{(m)}) - y^{(m)} \end{bmatrix}_{mx1}$

$$X^T(\text{prediction} - y) = \begin{bmatrix} x_0' x_0^2 & \dots & x_0^m \\ x_1' x_1^2 & \dots & x_1^m \end{bmatrix}_{2 \times m} \times \begin{bmatrix} h_\theta(x^{(1)}) - y^{(1)} \\ h_\theta(x^{(2)}) - y^{(2)} \\ h_\theta(x^{(3)}) - y^{(3)} \\ \vdots \\ h_\theta(x^{(m)}) - y^{(m)} \end{bmatrix}_{mx1}$$

$$= \begin{bmatrix} h_\theta(x^{(1)}) - y^{(1)} \cdot x_0^{(1)} + h_\theta(x^{(2)}) - y^{(2)} \cdot x_0^{(2)} + \dots + h_\theta(x^{(m)}) - y^{(m)} \cdot x_0^{(m)} \\ h_\theta(x^{(1)}) - y^{(1)} \cdot x_1^{(1)} + h_\theta(x^{(2)}) - y^{(2)} \cdot x_1^{(2)} + \dots + h_\theta(x^{(m)}) - y^{(m)} \cdot x_1^{(m)} \end{bmatrix}_{2 \times 1}$$

$$\theta := \theta - \alpha^* (X^T(\text{predictions} - y))$$

Feature Normalization

$$x_i \leftarrow \frac{x_i - \bar{u}_i}{\sigma_i} \quad \text{mean } \bar{u} = [u_1, u_2, u_3]_{1 \times 3}$$

$$X = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & x_3^{(1)} \\ x_1^{(2)} & x_2^{(2)} & x_3^{(2)} \\ \vdots & \vdots & \vdots \\ x_1^{(m)} & x_2^{(m)} & x_3^{(m)} \end{bmatrix}_{m \times 3}$$

$$X - \bar{u} = \begin{bmatrix} x_1^{(1)} - u_1 & x_2^{(1)} - u_2 & x_3^{(1)} - u_3 \\ x_1^{(2)} - u_1 & x_2^{(2)} - u_2 & x_3^{(2)} - u_3 \\ \vdots & \vdots & \vdots \\ x_1^{(m)} - u_1 & x_2^{(m)} - u_2 & x_3^{(m)} - u_3 \end{bmatrix}$$

$$\sigma = [\sigma_1, \sigma_2, \sigma_3]$$

$$(X - \bar{u}) / \sigma = \begin{bmatrix} \frac{x_1^{(1)} - u_1}{\sigma_1} & \frac{x_2^{(1)} - u_2}{\sigma_2} & \frac{x_3^{(1)} - u_3}{\sigma_3} \\ \vdots & \vdots & \vdots \\ \frac{x_1^{(m)} - u_1}{\sigma_1} & \frac{x_2^{(m)} - u_2}{\sigma_2} & \frac{x_3^{(m)} - u_3}{\sigma_3} \end{bmatrix}$$

Classification \Rightarrow

- \rightarrow Email \Rightarrow Spam / Not spam
- \rightarrow Online \Rightarrow Fraudulent (Y/N)
- Trans fraction
- Tumor \Rightarrow Malignant / Benign?

Two class classification

$y \in \{0, 1\} \rightarrow 0 \Rightarrow \text{"Negative Class" (e.g., benign)} \quad 1 \Rightarrow \text{"Positive Class" (e.g., malignant)}$

$y \in \{0, 1, 2, 3, 4, 5\} \Rightarrow$ Multi-class classification

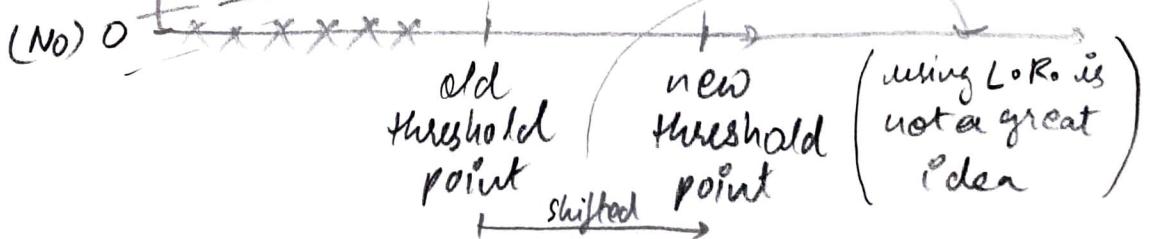


Threshold classifier output $h_0(x)$ at 0.5:

If $h_0(x) \geq 0.5$, predict "y=1"

If $h_0(x) < 0.5$, predict "y=0"

But, if



for Classification $\Rightarrow y = 0 \text{ or } 1$

$h_{\theta}(x)$ can be > 1 or $< 0 \Rightarrow$ linear regression is not a great learning Algo for classification problem

\Rightarrow Logistic Regression $\Rightarrow 0 \leq h_{\theta}(x) \leq 1$

\hookrightarrow this is a better learning algo. for classification problems

Logistic Regression Model

Want $0 \leq h_{\theta}(x) \leq 1$

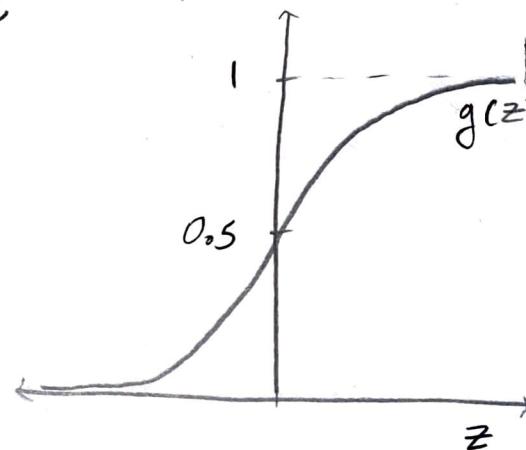
$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

* we will use diff. algo to fit parameters θ

$$g(z) = \frac{1}{1 + e^{-z}}$$

$z \in \mathbb{R}$

\downarrow
 $g: (0,1)$ Sigmoid function
 \Downarrow OR OR, Logistic function



$$0 < g(z) < 1$$

Interpretation of Hypothesis Output

$h_{\theta}(x)$ = estimated probability
that $y=1$ on output x

Eg) If $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorsize} \end{bmatrix}$

$h_{\theta}(x) = 0.7 \rightarrow$ probability of $y=1$ for given x
is 0.7

↳ Tell patient that 70% chance of
tumor being malignant

$h_{\theta}(x) = P(y=1|x; \theta) \Rightarrow$ "Probability that
 $y=1$, given x , para-
meterized by θ "

$$\therefore P(y=1|x; \theta) + P(y=0|x; \theta) = 1$$

$$P(y=0|x; \theta) = 1 - P(y=1|x; \theta) = 0.3$$

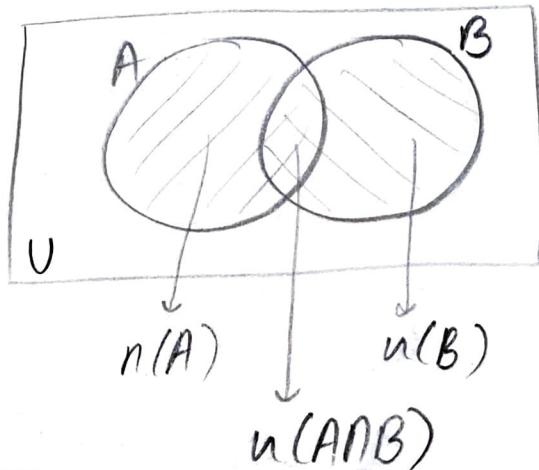
Tell the patient that 30% chance
of tumor

OR

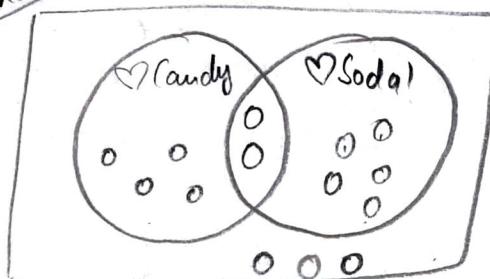
Probability of $y=0$ for given x
is 0.3

Conditional Probability

$$P(A|B) = \frac{n(A \cap B)}{n(B)} = \frac{n(A \cap B)/n(S)}{n(B)/n(S)} = \frac{P(A \cap B)}{P(S)}$$



Example



In Stat land people like to eat candy & soda

	Loves Candy	Doesn't love Candy	
Loves Soda	2	5	Total loves soda = 7
Doesn't love Soda	4	3	
Total Loves Candy = 4+2=6	8		# Doesn't love candy

$$P(\text{Loves Candy}) = \frac{6}{14} = \frac{3}{7}$$

$$P(\text{Loves Candy} | \text{Loves Soda})$$

$$P(\text{Loves Soda}) = \frac{7}{14} = \frac{1}{2}$$

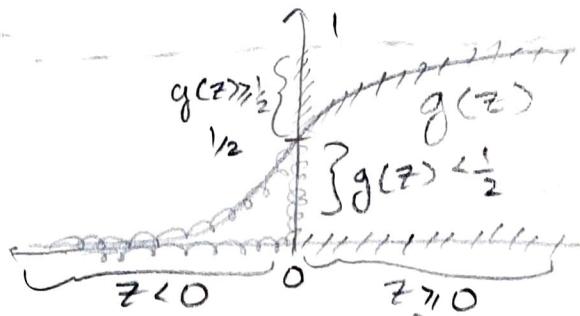
$$P(\text{Loves Soda}) = \frac{2}{7}$$

$$P(\text{Love Soda} | \text{Doesn't love candy}) = \frac{5}{8}$$

Decision Boundary

$$h_0(x) = g(\theta^T x)$$

$$g(z) = \frac{1}{1+e^{-z}}$$



$$h_0(x) = g(\theta^T x) = P(y=1|x; \theta)$$

Suppose,

Predict "y=1" if $h_0(x) \geq 0.5$

Predict "y=0" if $h_0(x) < 0.5$

∴ $g(z) \geq 0.5$ when $z \geq 0$

∴ $h_0(x) = g(\theta^T x) \geq 0.5$ whenever $\theta^T x \geq 0$

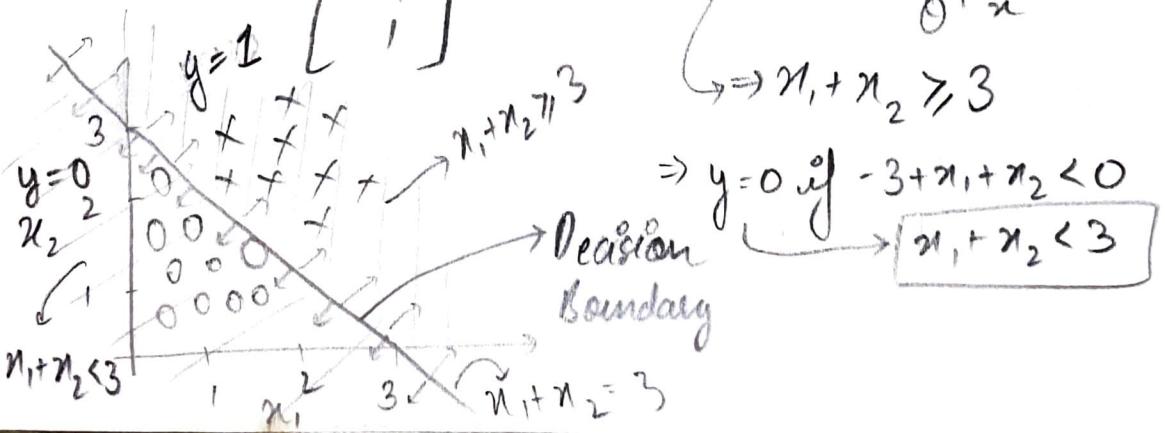
By, ∴ $g(z) < 0.5$ when $z < 0$

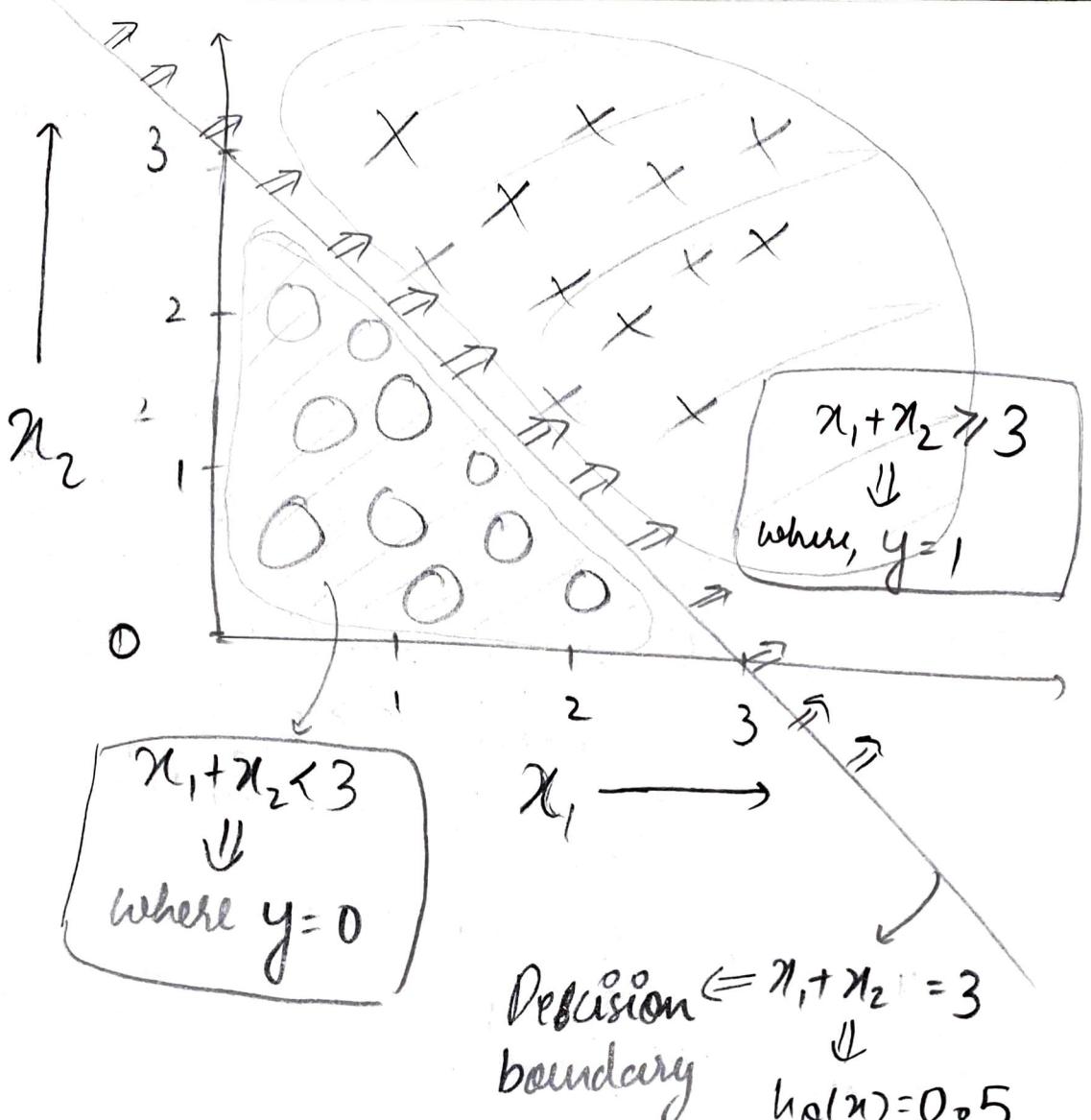
∴ $h_0(x) = g(\theta^T x) < 0.5$ whenever $\theta^T x < 0$

$$\text{let } h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$\text{let } \theta_0 = -3; \theta_1 = 1, \theta_2 = 1$$

$$\therefore \theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \quad \text{Predict } y=1 \text{ if } -3 + x_1 + x_2 \geq 0$$



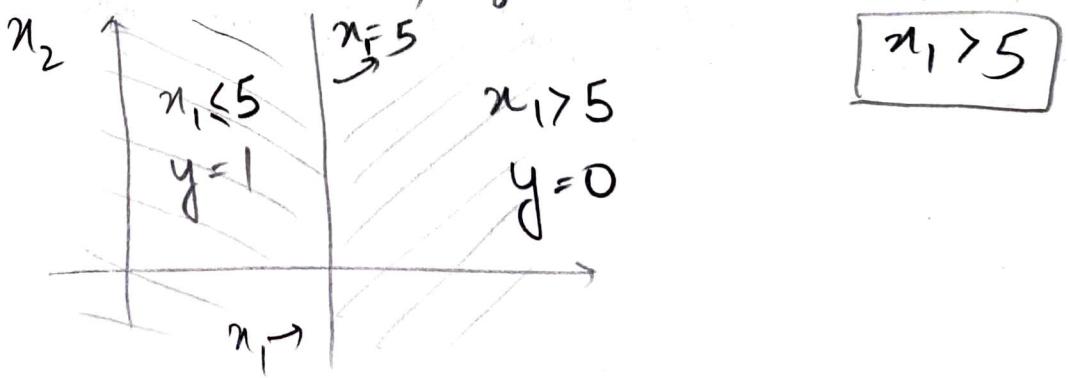


$$g(\theta = \begin{bmatrix} 5 \\ -1 \\ 0 \end{bmatrix}) \Rightarrow h_0(x) = g(5 - x_1)$$

$$h_0(x) \geq 0.5 \Rightarrow \theta^T x = 5 - x_1 \geq 0$$

$$x_1 \leq 5$$

Thus, $h_0(x) < 0.5 \Rightarrow 5 - x_1 < 0$



Non-linear decision boundaries

let $h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$



$$\text{if } \theta_0 = -1, \theta_1 = 0, \theta_2 = 0, \theta_3 = 1, \theta_4 = 1$$

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

Predict:

$$\textcircled{1} \quad y=1 \Rightarrow -1 + x_1^2 + x_2^2 \geq 0$$

$$x_1^2 + x_2^2 \geq 1$$

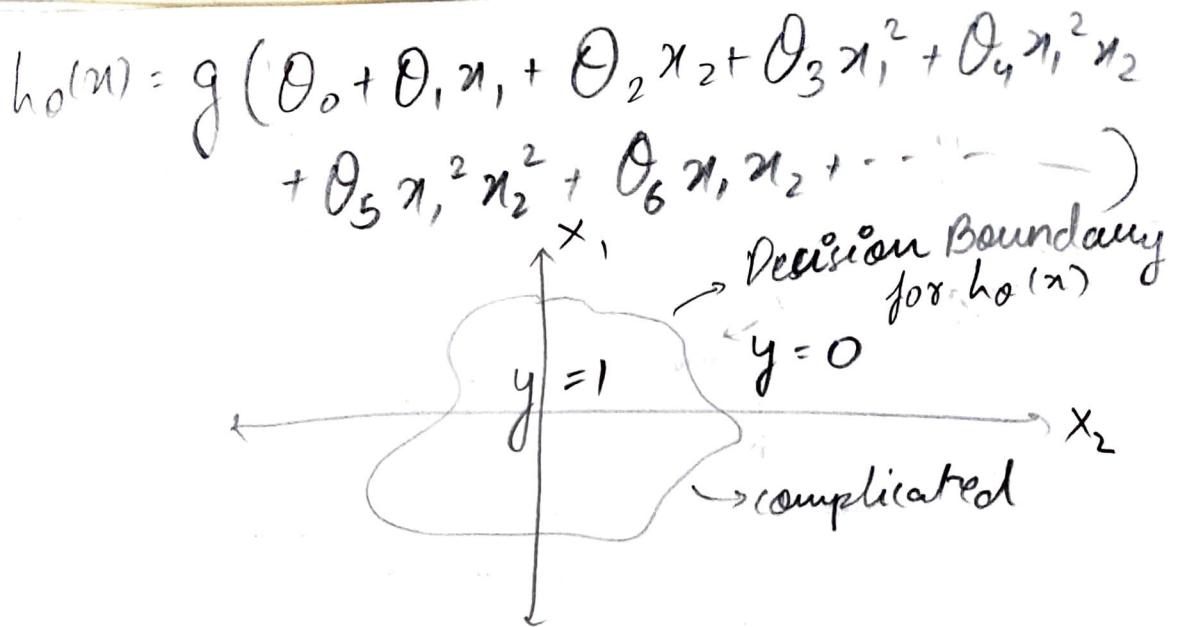
Eqn of circle ($x_1^2 + x_2^2 = 1$)

$$x_1^2 + x_2^2 = 1 \quad \textcircled{2} \quad y=0 \Rightarrow -1 + x_1^2 + x_2^2 < 0$$

↳ Eqn of circle ($(0,0); r=1$)

$$x_1^2 + x_2^2 < 1$$

* Decision Boundary is the property of parameters
 θ in $h_\theta(x)$ and not of training set



Training Set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

examples: $x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}_{n+1}$ $x_0 = 1, y \in \{0, 1\}$

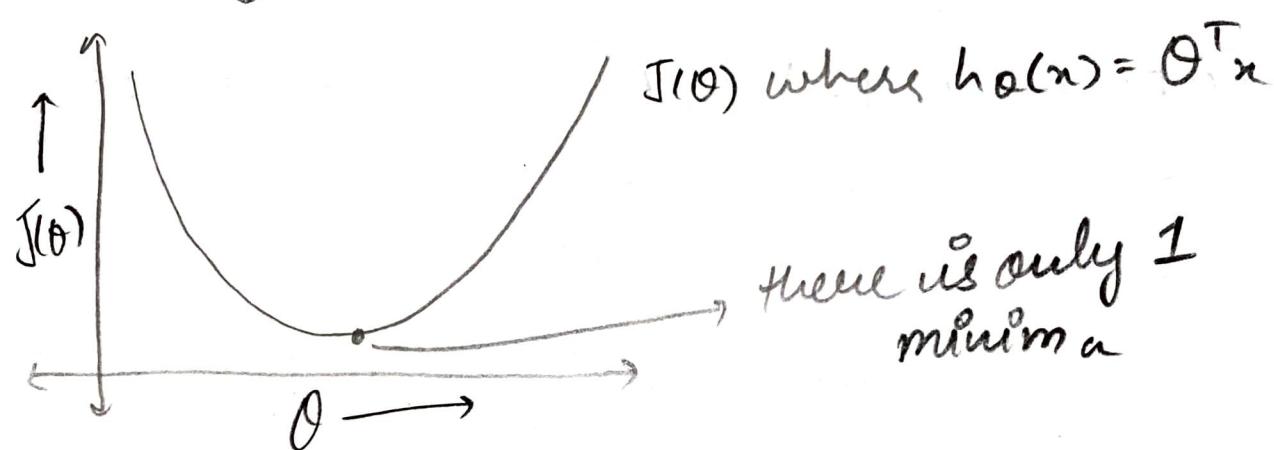
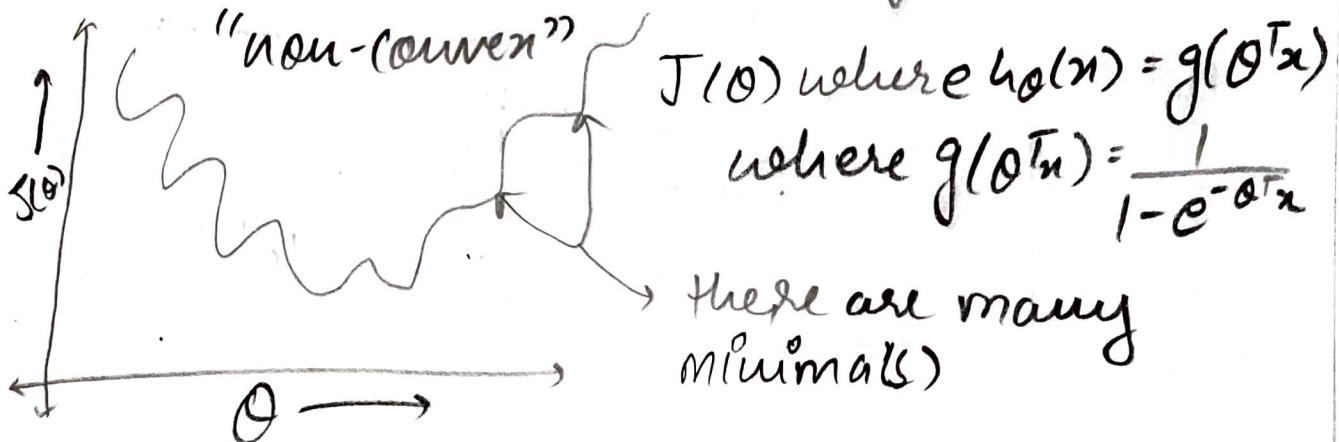
$$h_0(x) = \frac{1}{1 + e^{-\theta^T x}}$$

(cost funcⁿ in Linear Regression) $\Rightarrow J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_0(x^{(i)}) - y^{(i)})^2$

let, Cost($h_0(x^{(i)}), y^{(i)}$) = $\frac{1}{2} (h_0(x^{(i)}) - y^{(i)})^2$
OR

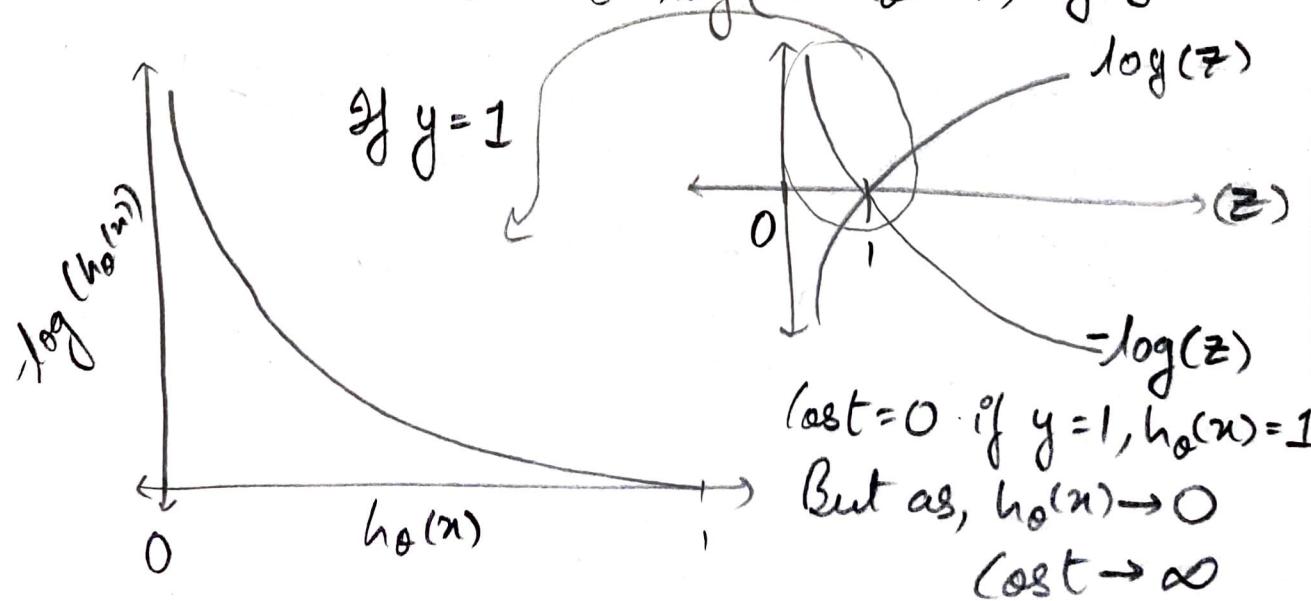
Cost($h_0(x), y$) = $\frac{1}{2} (h_0(x) - y)^2$
[for linear Regression]

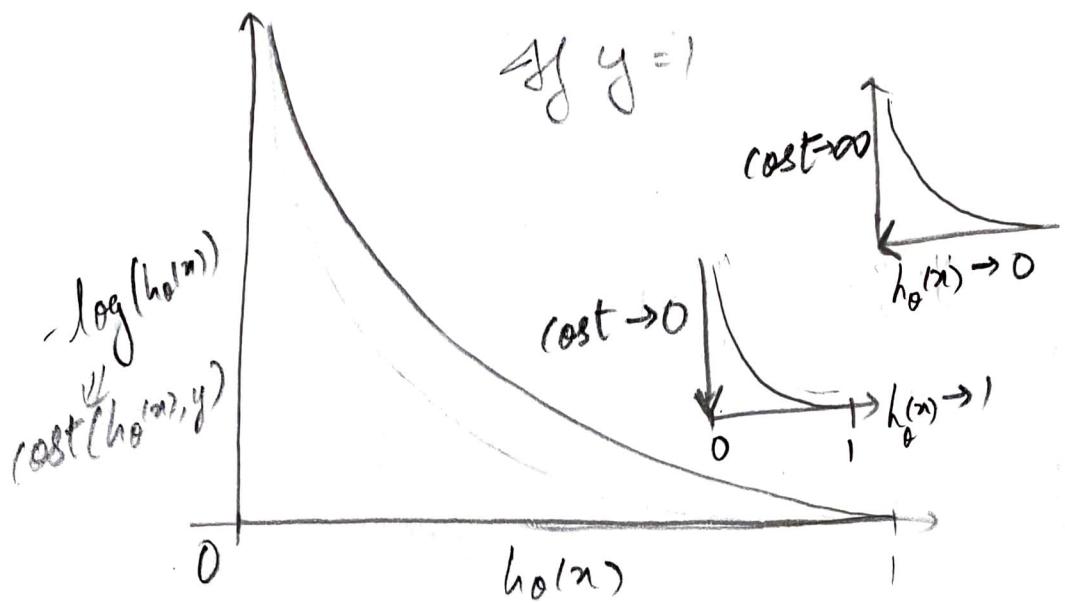
But for Logistic Regression the above cost funcⁿ won't work as it will return a non-convex funcⁿ



Logistic Regression Cost Function

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$



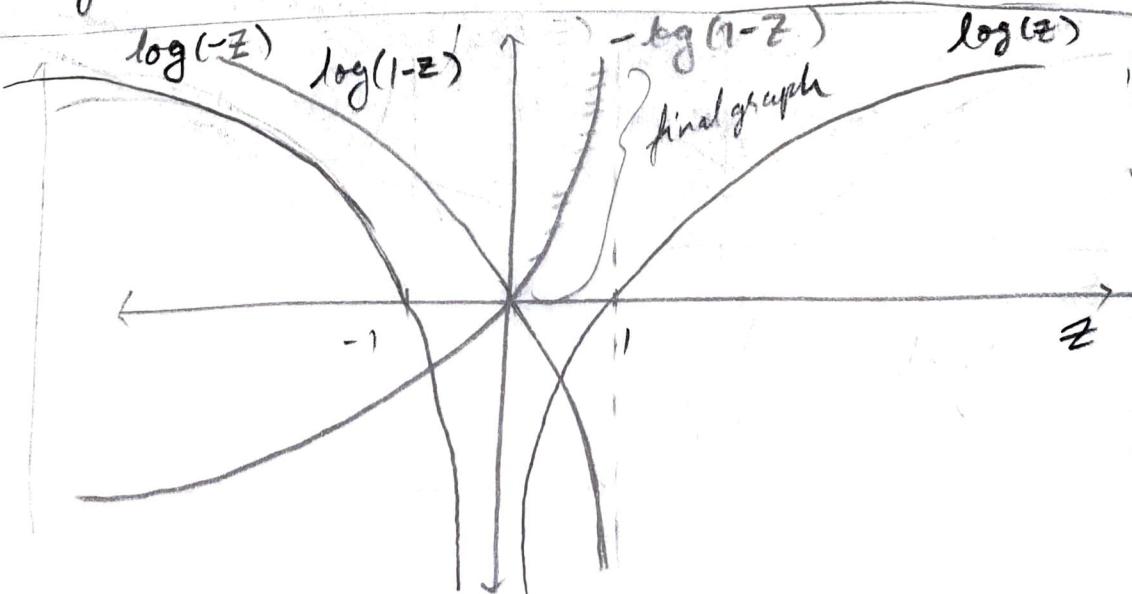


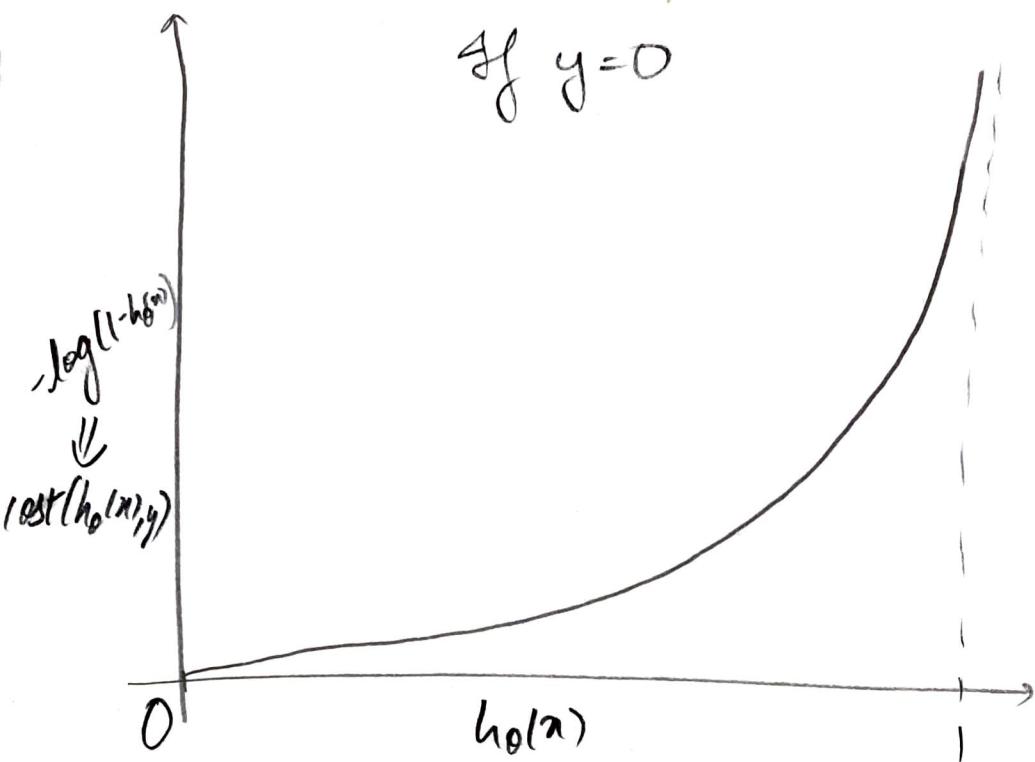
$\text{cost} = 0$ if $y = 1, h_0(x) = 1$

But, as $h_0(x) \rightarrow 0$

$\text{cost} \rightarrow \infty$

Captures intuition that if $h_0(x) = 0$,
(i.e. predict $P(y=1|x; \theta) = 0$), but $y=1$,
we'll penalize learning algorithm
by a very large cost





$\text{Cost} = 0 \text{ if } y=0, h_\theta(x)=0$

But as, $h_\theta(x) \rightarrow 1$

$$\log y \rightarrow \infty$$

(captures intuition that if $h_\theta(x) = 0$,
 (i.e.) predict $P(y=0|x; \theta) = 1$), but
 $y=0$, then we'll penalize learning
 algorithm by a very large cost

logistic Regression cost function:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$\text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y=1 \\ -\log(1-h_\theta(x)) & \text{if } y=0 \end{cases}$$

Note $\Rightarrow y=0$ or 1 always

$$\text{Cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x))$$

If $y=1 \Rightarrow \text{Cost}(h_\theta(x), y) = \boxed{-\log(h_\theta(x))} - 0$

$\hookrightarrow \boxed{1-y=0}$

If $y=0 \Rightarrow \text{Cost}(h_\theta(x), y) = 0 - (1-0) \log(1-h_\theta(x))$

$\hookrightarrow \boxed{y=0}$

$$\text{Cost}(h_\theta(x), y) = \boxed{-\log(1-h_\theta(x))}$$

$$\therefore J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$= -\frac{1}{m} \left[\sum_{i=1}^m [y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)}))] \right]$$

To fit Parameters θ $\min_{\theta} J(\theta)$

To make a prediction: Output $h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$
given new x

$$P(y=1|x; \theta) \leftarrow$$

Gradient descent for Logistic Regression.

We need $\min_{\theta} J(\theta)$:

Repeat until convergence {

$$\theta_j^{\circ} = \theta_j^{\circ} - \alpha \frac{\partial}{\partial \theta_j} (J(\theta)) \quad (\text{for } j=0, 1, 2, \dots, m)$$

} (Simultaneously update all θ_j°)

$$\begin{aligned}\frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \left(-\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) - (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)})) \right) \\ h_{\theta}(x) &= \sum_{j=0}^n \theta_j x_j \\ &= -\frac{1}{m} \left(\left(y^{(i)} \frac{1}{h_{\theta}(x^{(i)})} x_j^{\circ} - (1-y^{(i)}) \frac{1}{1-h_{\theta}(x^{(i)})} (-x_j^{\circ}) \right) \right) \\ &= \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{\circ}\end{aligned}$$

}, Repeat {

$$\theta_j^{\circ} = \theta_j^{\circ} - \frac{\alpha}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{\circ} \quad (\text{for } j=0, 1, 2, \dots, n)$$

} (Simultaneous update all θ_j°)

* Looks identical to Linear Regression even though hypothesis func $\Rightarrow h_{\theta}(x) = \theta^T x$ & $h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \quad n = \begin{bmatrix} n_0 \\ n_1 \\ \vdots \\ n_n \end{bmatrix}$$

Linear Regression

Logistic Regression

Vectorization of Gradient Descent for Logistic Regression. ($\theta := \theta - \alpha \delta$)

$$X = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ x_0^{(2)} & x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_0^{(m)} & x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \Rightarrow \mathbb{R}^{n+1 \times 1}$$

$$h_{\theta}(x) = \begin{bmatrix} h_{\theta}(x^{(1)}) \\ h_{\theta}(x^{(2)}) \\ \vdots \\ h_{\theta}(x^{(m)}) \end{bmatrix} \Rightarrow \mathbb{R}^m \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} \Rightarrow \mathbb{R}^m$$

$$\bar{J}(\theta) = \frac{1}{m} (-y^T \log(h_{\theta}(x)) - (1-y)^T \log(1-h_{\theta}(x)))$$

$$h_{\theta}(x) = X\theta := \begin{bmatrix} \theta_0 x_0^{(1)} & \theta_1 x_1^{(1)} & \dots & \theta_n x_n^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ \theta_0 x_0^{(m)} & \theta_1 x_1^{(m)} & \dots & \theta_n x_n^{(m)} \end{bmatrix}$$

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

Code \Rightarrow

$$\text{theta} = \text{theta} - (\alpha/m) * X^T * (X * \text{theta} - y)$$

Mathematical \Rightarrow

$$\boxed{\theta := \theta - \frac{\alpha}{m} X^T (g(X\theta) - y)}$$

Advanced Optimisation for Logistic Regression

-Plan.

Given θ , we have code that can compute
 $\rightarrow J(\theta)$

$\rightarrow \frac{\partial}{\partial \theta_j} J(\theta) \text{ (for } j=0, 1, 2, \dots, n)$

Optimisation Algorithms:

- \rightarrow Gradient descent (One we are using)
 - \rightarrow Conjugate gradient
 - \rightarrow BFGS
 - \rightarrow L-BFGS
- $\left. \begin{array}{l} \text{Advantages:} \\ \rightarrow \text{No need to manually} \\ \text{pick } \alpha \\ \rightarrow \text{Often faster than gradient} \\ \text{descent} \end{array} \right\}$

Disadvantage: More complex

$$\text{Ex: } \theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \& J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$$

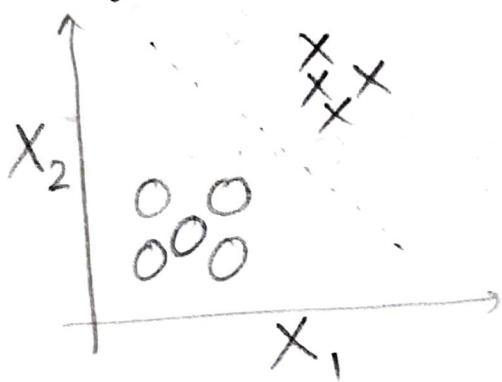
$$\text{for } \min_{\theta} J(\theta) \Rightarrow \theta_1 = 5, \theta_2 = 5$$

$$\frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5) \quad \left\{ \begin{array}{l} \frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5) \end{array} \right.$$

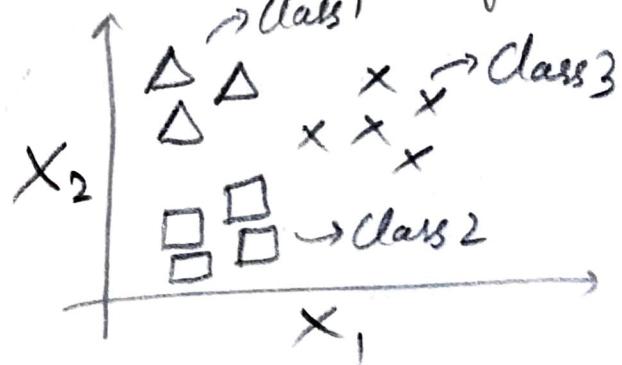
Logistic Regression for Multi-class Classification

- ① Email foldering \Rightarrow Work, Friends, Family, / tagging
- $y=1$ $y=2$ $y=3$
 $y=4$ Hobby
- ② Medical diagrams \Rightarrow Not ill, cold, flu
- $y=1$ $y=2$ $y=3$
- ③ Weather \Rightarrow Sunny, Cloudy, Rain, Snow
- $y=1$ $y=2$ $y=3$ $y=4$

Binary Classification



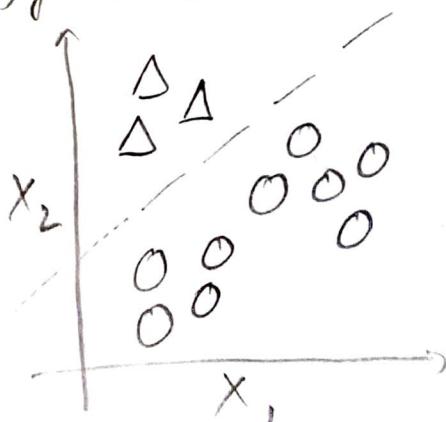
Multi-class Classification



One vs All (or, one v/s rest)

① for class 1

② for class 2

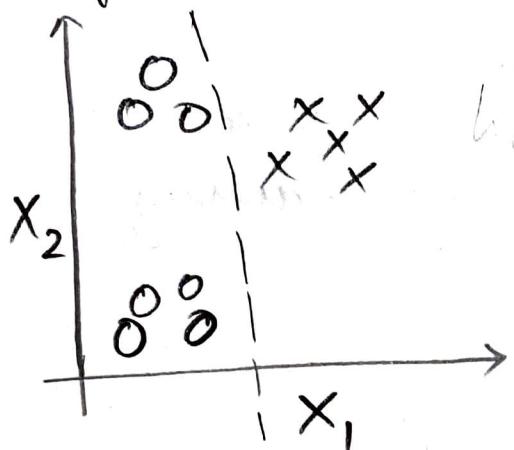


$$h_0^{(1)}(x) = P(y=1|x; \theta)$$

$$h_0^{(2)}(x)$$

$$P(y=2|x; \theta)$$

③ for Class 3



$$h_0^{(3)}(x) = P(y=3|x; \theta)$$

$$h_0^{(i)}(x) = P(y=i|x; \theta)$$

$$(i=1, 2, 3)$$

One v/s all

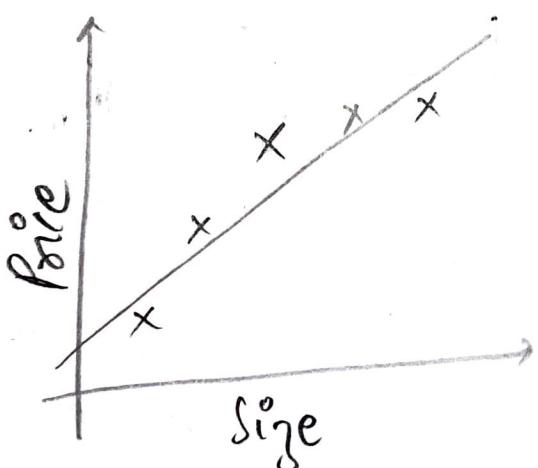
Train a logistic regression
 $h_0^{(i)}(x)$ for each class i to
predict the probability
that $y=i$.

On a new input x , to make a prediction, pick the
class i that maximizes

$$\max_i h_0^{(i)}(x)$$

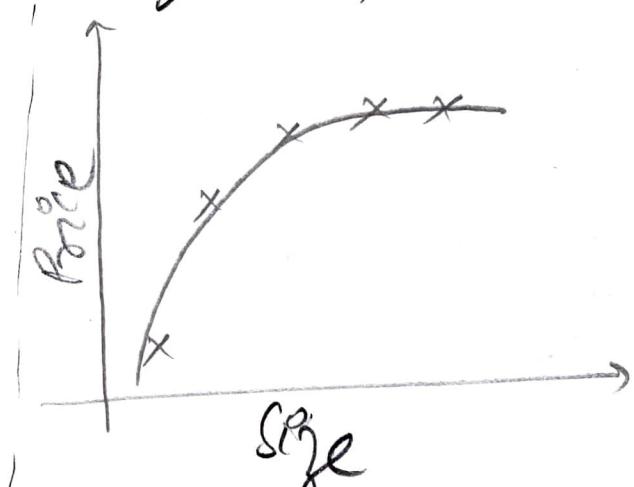
⇒ The problem of overfitting

Eg) Linear Regression (housing Prices)



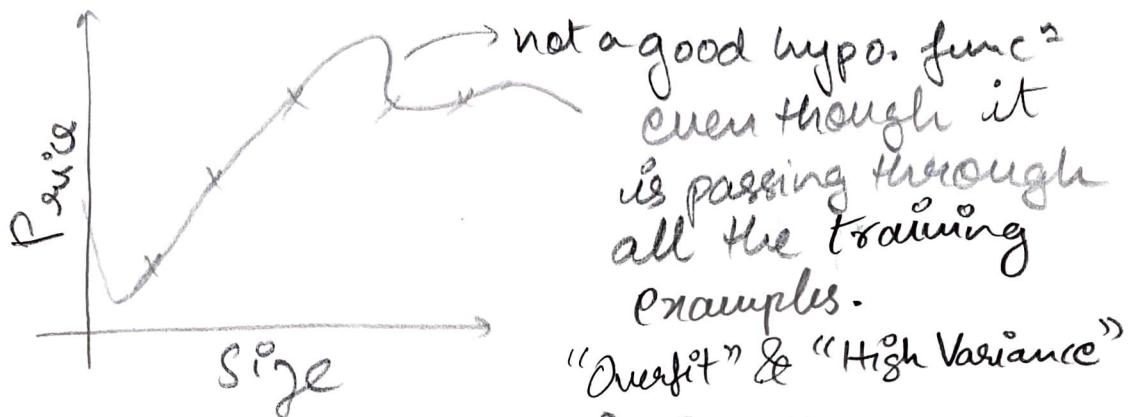
$$\theta_0 + \theta_1 x_1$$

Underfit, High bias



$$\theta_0 + \theta_1 x_1 + \theta_2 x_1^2$$

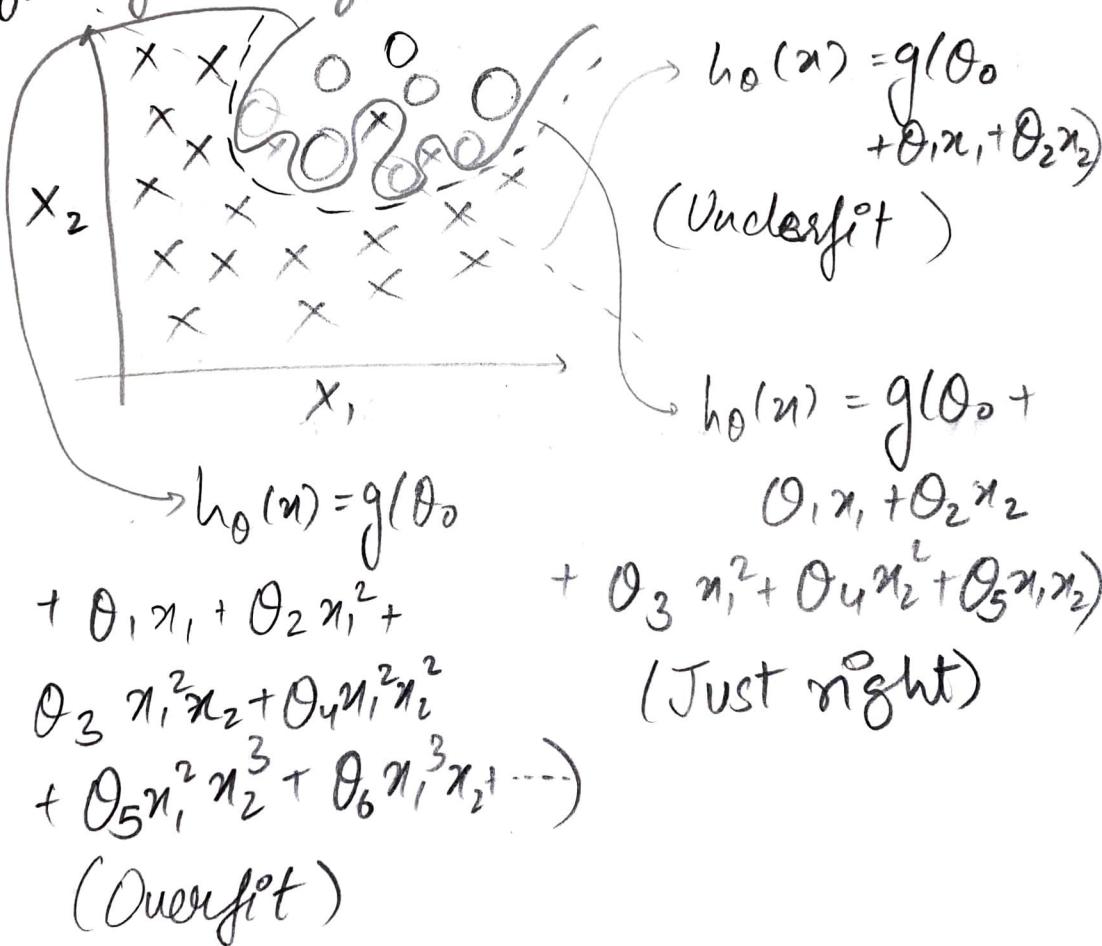
Just Right



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Overfitting \Rightarrow If we have too many features, the learned hypothesis may fit the training

Eg) Logistic Regression



Addressing Overfitting:

$x_1 = \text{size}$

$x_2 = \# \text{ bedrooms}$

$x_3 = \# \text{ floors}$

:

$x_{100} = \text{Area}$



too many features
with less training data
will make hypo funcⁿ overfitting

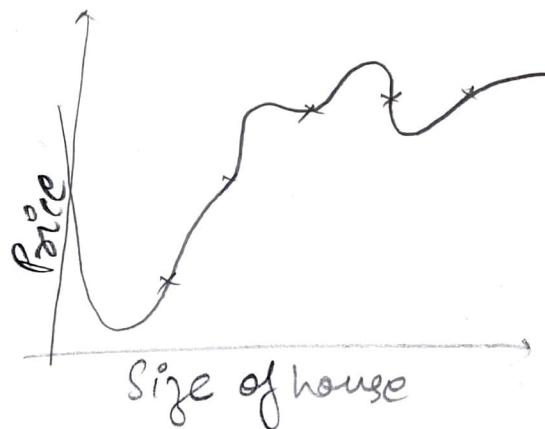
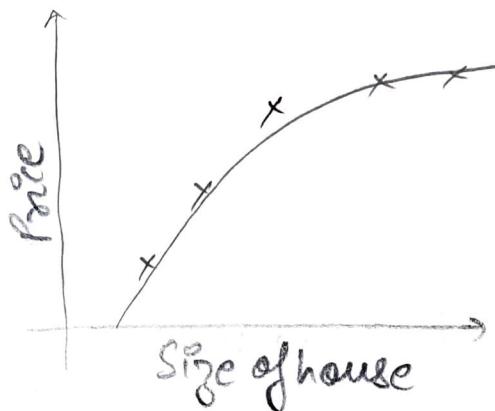
① Reducing # features.

- ↳ Manually select which features to keep
- ↳ Model selection algorithm

② Regularizations

- ↳ keep all the features, but reduce magnitude / values of parameters θ_j
- ↳ Works well when we have a lot of features, each of which contributes a bit to predicting y

Intuition



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

Suppose we penalize and make $\theta_3 \approx 0$ & $\theta_4 \approx 0$.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000\theta_3^2 + 1000\theta_4^2$$

$\hookrightarrow \theta_3 \approx 0 \text{ & } \theta_4 \approx 0 \Rightarrow \theta_0 + \theta_1 x + \theta_2 x^2 + 0 + 0$

Regularization

Small values for parameters θ_1, θ_2 -

θ_m

→ simpler hypothesis

→ less prone to overfitting

e.g.) Housing :

→ features $\rightarrow x_1, x_2, \dots, x_{100}$

→ Parameters $\rightarrow \theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

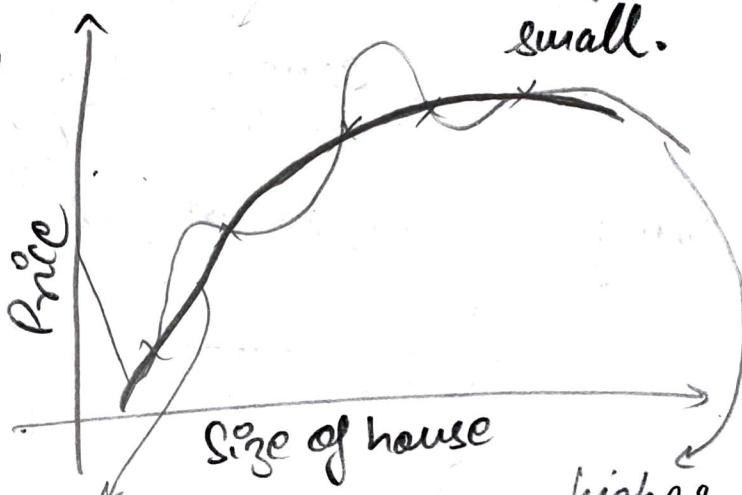
$$\min_{\theta} J(\theta)$$

Regularization parameter

this term helps fitting training data well.

Regularization term

this term helps in keeping the parameters small.



higher order hypo. funcⁿ

without regularization
with regularization
(smoother and better)

higher order hypo. funcⁿ

If λ is very large,

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

for very large $\lambda \Rightarrow \theta_1 \approx 0, \theta_2 \approx 0 \dots, \theta_n \approx 0$

$\therefore h_\theta(x) = \theta_0 \Rightarrow$ constant line // x

\downarrow to x axis
Underfit or 'high bias'

Regularization linear regression.

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

$$\min_{\theta} J(\theta)$$

Grad. Desc.

Repeat {

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

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

}

$$\theta_j^{(0)} = \theta_j^{(0)} \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$\left[1 - \alpha \frac{\lambda}{m} < 1\right]$ Just little bit less than 1.

$$\theta_j^{(0)} = \theta_j^{(0)} \left(1 - \alpha \frac{\lambda}{m}\right) - \underbrace{\alpha \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}}_{0.99 \approx 1}$$

Original grad-desc. formula without regularization

Normal Eq \approx (with Regularization)

$$X = \begin{bmatrix} (x^{(1)})^\top \\ \vdots \\ (x^{(m)})^\top \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \Rightarrow \mathbb{R}^m$$

$\min_{\theta} J(\theta)$

$$\theta = (X^T X + \lambda \begin{bmatrix} 0 & 0 & \dots \\ 0 & 1 & 0 & \dots \\ \vdots & \vdots & \ddots & \dots \\ 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 1 \end{bmatrix})^{-1} X^T y$$

($n+1$) ($n+1$)

Eg) $n=2 \Rightarrow \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

Non-invertibility (optional / advanced)

Suppose $m \leq n$

\downarrow \rightarrow # features
examples

$$\theta = (X^T X)^{-1} X^T y \quad \begin{matrix} \xrightarrow{\text{Pinv}} \\ \text{pseudo inverse} \end{matrix}$$

If $\lambda > 0$, $\xrightarrow{\text{non-invertible / singular}}$

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

Regularized Logistic Regression

$$J(\theta) = -\left[\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

↓ new cost funcⁿ with regularization term

Grad Des.

Repeat {

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

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

{ $\frac{\partial J(\theta)}{\partial \theta_j}$ $\xrightarrow{\text{new cost func}}$

Advo Optimisation (for new cost func^{*})

function [JVal, gradient] = costFunction(theta)

JVal = [code to compute $\underline{J(\theta)}$]; $\xrightarrow{\text{new cost func}}$

gradient(1) = [code to compute $\frac{\partial}{\partial \theta_0} J(\theta)$];

gradient(2) = [code to compute $\frac{\partial}{\partial \theta_1} J(\theta)$];

⋮

gradient(n+1) = [code to compute $\frac{\partial}{\partial \theta_n} J(\theta)$];

↳ function (@costFunction, ...)

↳ for applying