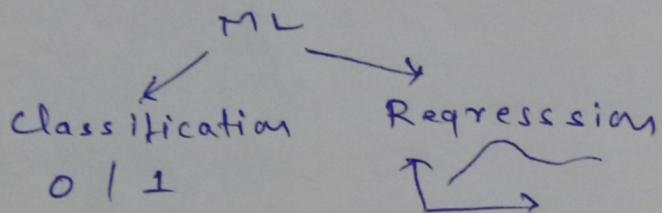
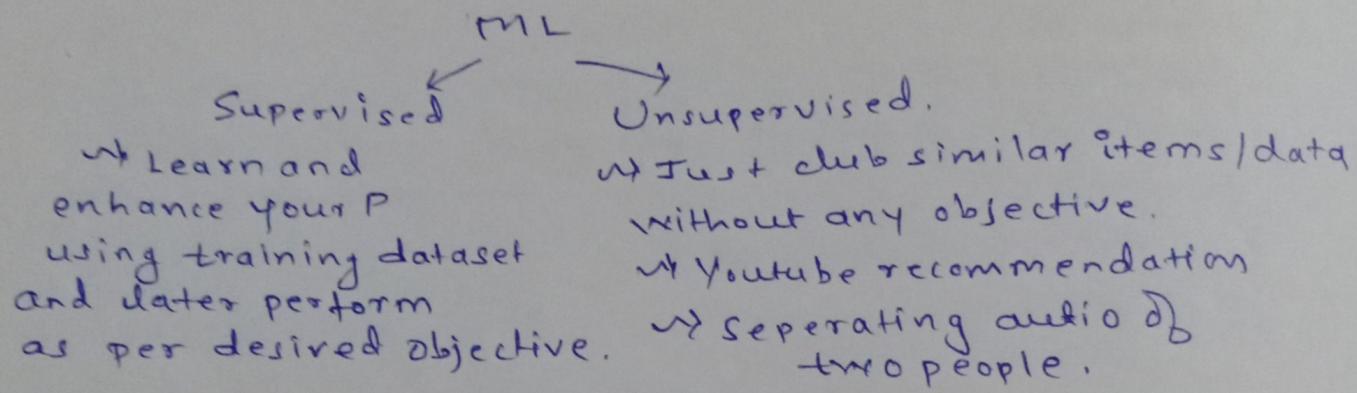


* Machine Learning *

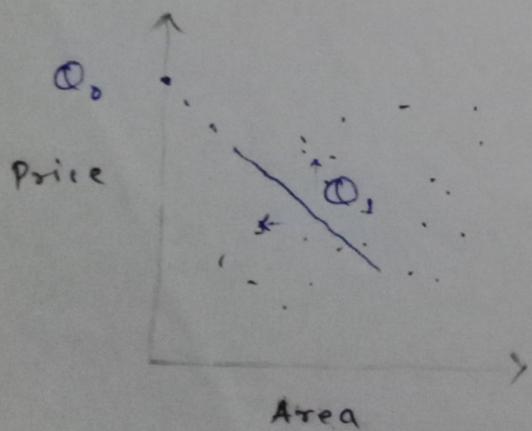
→ Doing tasks without being programmed explicitly.

→ Experience E is obtained by doing task T, where performance is measured by P.

Measure of Performance P in doing a task T, increases with task experience E



* LINEAR REGRESSION *



1) Assuming linear relationship
This we will learn later

Find:

best $h: A \rightarrow P$
such that the squared error is min average

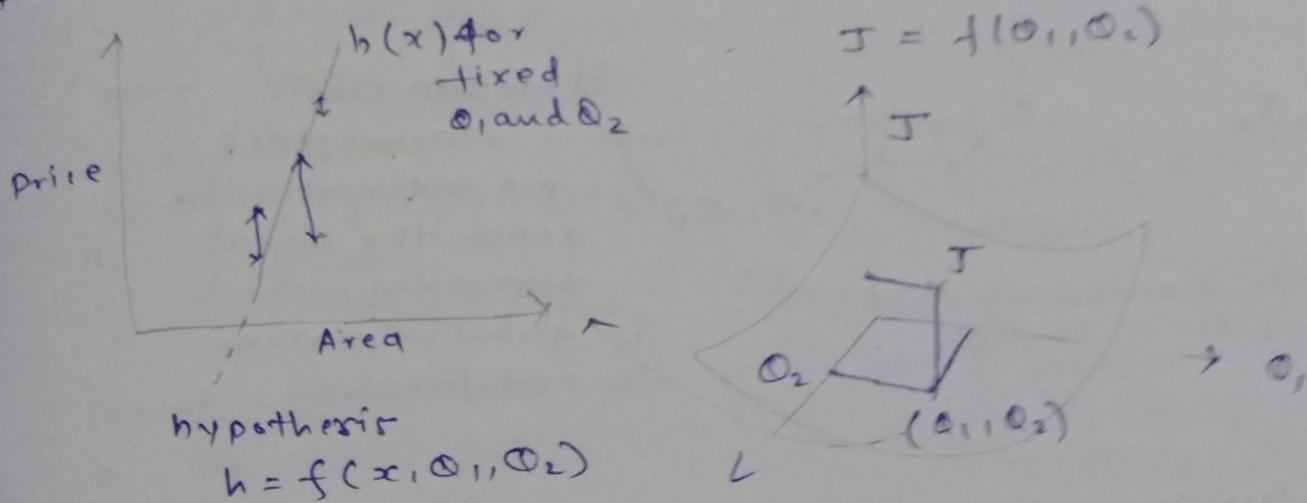
2) Error-squared function: $J(\Theta_0, \Theta_1)$

$$= \frac{1}{2m} \sum_{i=0}^m (h(x) - y)^2$$

* Gradient Descent *

→ It's about finding local minima.

* Understand this condition



We need such θ_0 and θ_1 that J is min.

* Understand this equation.

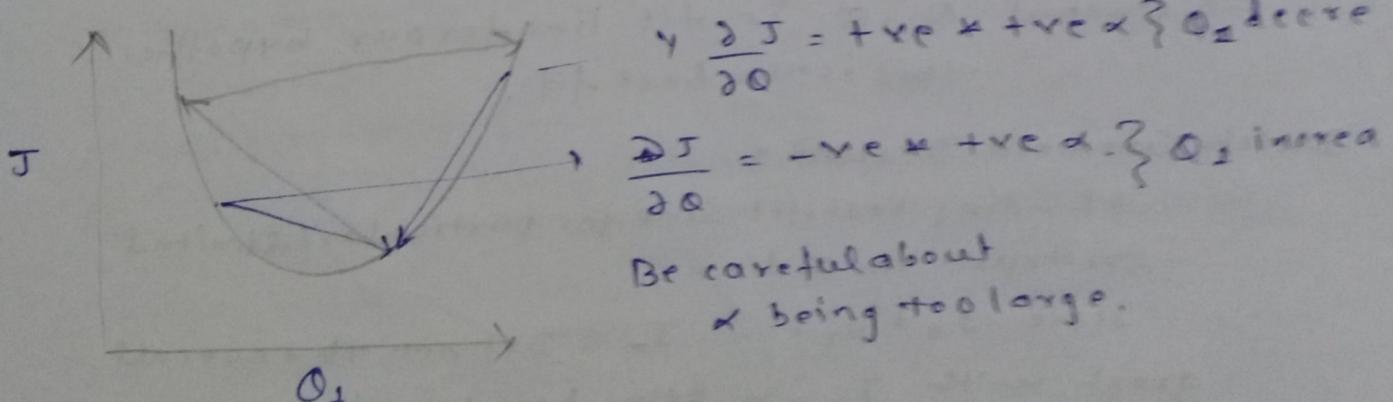
→ learning rate.

$$\left\{ \begin{array}{l} \theta_0 = \theta_0 - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} \\ \theta_1 = \theta_1 - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} \end{array} \right\} \text{Do until convergence}$$

This is incorrect
do simultaneous update.

* Understand the learning rate.

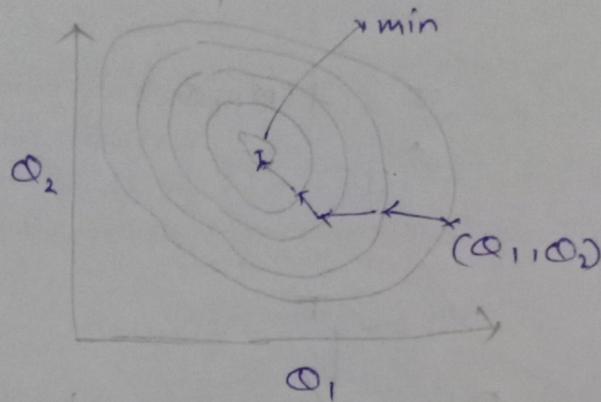
$$J = f(\theta_0) \text{ suppose } \Rightarrow h = \theta_0 x$$



Back to normal:

$$J = f(\theta_0, \theta_1)$$

Contour plot



If we run until convergence we automatically reach the local minimum which is global minimum for cost function

Conclusion

$$h: A \rightarrow P \Rightarrow \theta_0 + \theta_1 x$$

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

\Rightarrow find $\theta_0, \theta_1 \Rightarrow$ s.t $J = \min J(\theta_0, \theta_1)$

Gradient descent

Repeat until convergence {

$$\text{temp1} := \theta_1 - \alpha \frac{\partial J}{\partial \theta_1}$$

$$\text{temp2} := \theta_2 - \alpha \frac{\partial J}{\partial \theta_2}$$

$$\theta_1 = \text{temp1}$$

$$\theta_2 = \text{temp2}$$

\Rightarrow Output θ_0, θ_1 is fit line for our hypothesis
~ this is the best fit.

Batch Gradient Descent:

if finding best fit line for particular/limited spectrum on x-line.

Here $J_1 + J_2 + \dots + J_n \}$ defines the overall graph with J_1, \dots, J_n being best fits for their respective domains

$$J_1 = \frac{1}{2N_1} \sum_{i=1}^{N_1} (h - y)^2$$

$$J_2 = \frac{1}{2N_2} \sum_{i=1}^{N_2} (h - y)^2$$

Best of the procedure is the same

* Linear Algebra *

Properties of Matrix Multiplication:

1. $AXB \neq BXA$ } Not always but may be.
2. $(AXB)XC = AX(BXC)$

$$3. \underset{m \times m}{I} \times \underset{m \times n}{A} = \underset{m \times n}{A} \times \underset{n \times n}{I} = A.$$

For square matrices inverse doesn't exist:

$$4. \text{ non } A(A^{-1}) = A^{-1}A = I \quad \left. \begin{array}{l} \text{where} \\ A_{\text{LHS}}^{-1} = A_{\text{RHS}}^{-1} \end{array} \right\} \text{Not for all } A.$$

Note: what actually
is an inverse
of matrix.

Transpose of a matrix:

$$5. B_{ji} = A^T = A_{ij}$$

* Multidimensional linear regression *

Earlier : h : Area \rightarrow Price

Now : h : Bedrooms Area + Floors + Age \rightarrow Price.

Nomenclature :

$X^{(i)}$ \rightarrow i^{th} training set element
 $X^{(i)}_n \rightarrow$ n^{th} attribute of that training set.)

$$h = \underbrace{Q_0}_{1} x_0 + Q_1 x_1 + Q_2 x_2 + \dots + Q_n x_n. \quad \left. \begin{array}{l} Q = \begin{bmatrix} Q_0 \\ Q_1 \\ \vdots \\ Q_n \end{bmatrix} \end{array} \right\}$$

$$\therefore h = Q^T \cdot X$$

Intuition :

$Q_0 \rightarrow$ base price

$Q_1 \rightarrow$ prop. w.r.t no. of bedrooms

$Q_2 \rightarrow$ prop. w.r.t. no. of floors

$\vdots \vdots$

$$X = \begin{bmatrix} x_0 & 1 \\ x_1 & \\ \vdots & \\ x_n & \end{bmatrix}$$

* Multivariate cost function *

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

↓ Training set
specially this
don't confuse with
 $\theta_0 x_0 + \theta_1 x_1 + \dots$

Training set

Parameter.

$$1 \quad x_0^1, x_1^1, x_2^1, x_3^1$$

$$2 \quad x_0^2, x_1^2, x_2^2, x_3^2$$

$$\begin{aligned}\therefore J(\theta) &= \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{2m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2\end{aligned}$$

* Gradient Descent function *

Earlier:

$$\theta_k = \theta_k - \frac{\alpha}{2} \frac{\partial J(\theta_0, \theta_1, \theta_2, \dots)}{\partial \theta_k}$$

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

$$= \theta_k - \frac{\alpha}{m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) x_k^{(i)}$$

* How to make Gradient Descent work faster in action *

Suppose:

(a) Parameter 1 $\{ \text{size} \} \in (0, 200,000)$

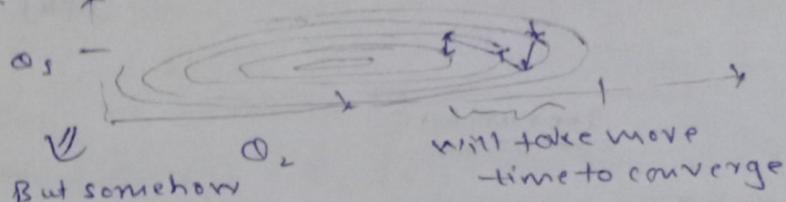
(b) Parameter 2 $\{ \text{bed} \} \in (0, 2)$

$$\text{Price} = \text{C}_0 + \text{C}_1 S + \text{C}_2 B$$

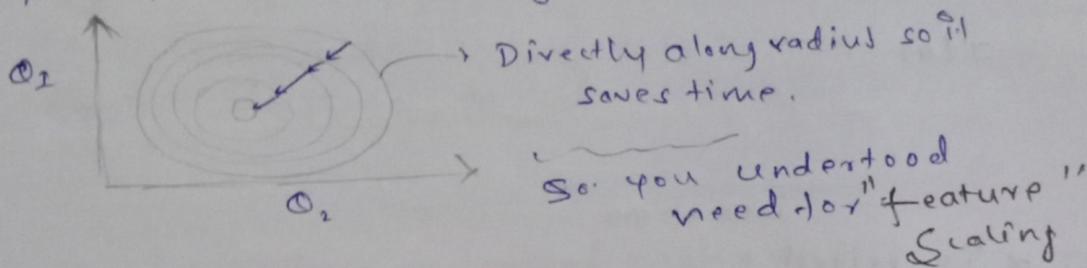
assume
0 ↓ >>> >>>

LLL 0, w + C_2

(c) \therefore on \mathbb{R}^2 axis



But if C_1 and C_2 are of same order



Now let's work with normalisation:

Before that:

$$\text{for Price} = (C_1 S) + (C_2 B)$$

\downarrow why these two terms should

nearly be of same order?

Because their existence matter

that is why we have included them

Easier way to digest scaling:

$$\text{Price} = (\underbrace{C_1 \times 2 \times 10^5}_{\text{C}_1}) \left(\frac{S}{2 \times 10^5} \right) + \dots$$

* Why normalisation?

From where {which point} to begin gradient descent calculation...
 ↓ let's say $(0, 0, 0, 0, \dots)$

$$Data = \frac{\text{Data} - \text{Mean}}{\text{range} (\text{max} - \text{min})}$$

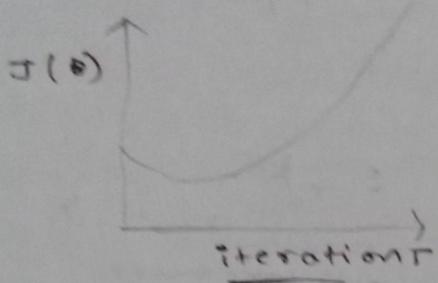
We can have std deviation also.

* Learning Rate *

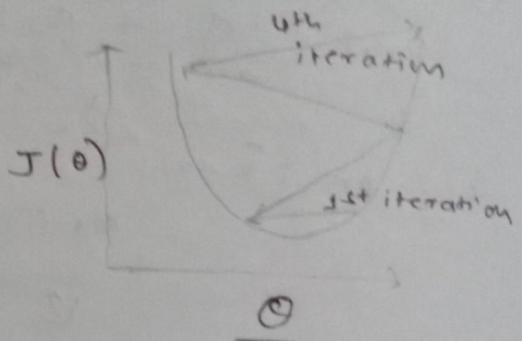
what should be learning rate α ?

① How to make sure that gradient descent is working properly?

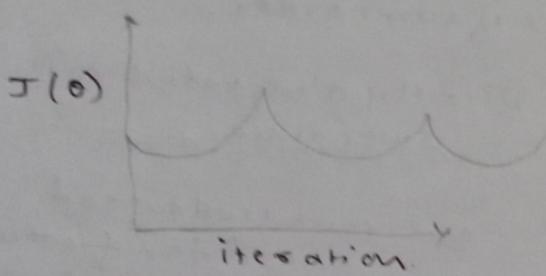
(i) $\alpha > \gamma$



Because \Rightarrow

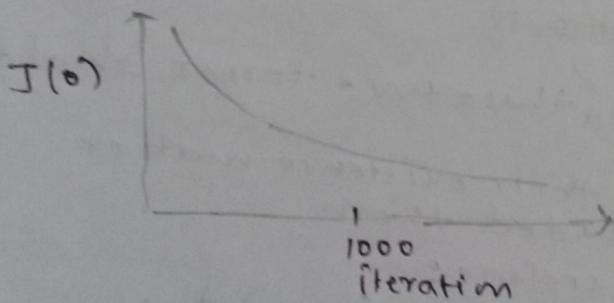


(ii) $\alpha < \gamma$

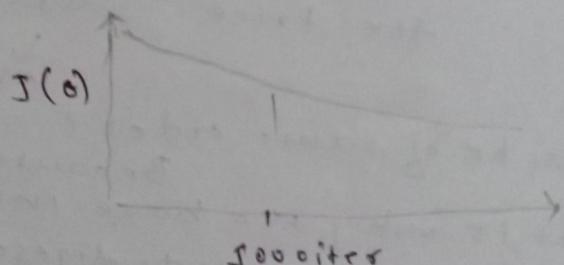


(iii) $\alpha = \gamma$ perfect :

"For $\alpha \in \{ \text{sufficiently small } \lambda \}$, $J(\theta)$ should decrease with every iteration"



perfect α
is largest correct
 α .



α very small
very

Mathematicians were pissed off and hence they couldn't derive a method for perfect α hence they said,

Try:

$$0.001 \cdot \gamma \cdot 0.01 \cdot \gamma \cdot 0.1 \dots$$

$\frac{0.001}{0.003} \quad \frac{0.01}{0.03}$

And check
 $J(\theta)$'s iterations
graph.

Pick the largest possible value.

* What feature to take *

$h: A \rightarrow P$ or $h: \text{width, length} \rightarrow P$

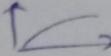
* polynomial regression *

Let hypothesis be:

$$h: \text{size} \rightarrow P = \theta_0 + \theta_1 (\text{size})^1 + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3$$

Then behave with $(\text{size})^2$ as a separate entity,
and normalise my feature scale or independently.

Ex. let size $\in (0, 1000)$

and perfect fit be parabola 

$$\therefore h \approx \theta_0 + \theta_1 \text{size} + \theta_2 \sqrt{\text{size}}$$

$$\therefore x_1 \leftarrow \frac{\text{size} - 500}{1000}$$

$x_2 \leftarrow \frac{\sqrt{\text{size}} - 10}{32}$

This is how.

* How we do it in a vectoised manner *

$$X = \begin{bmatrix} x_0 & x_1 & \dots & x_n \end{bmatrix}^T \quad \left. \begin{array}{l} \text{number of features} \\ \text{No. of training set} \end{array} \right\} m \times n$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}^T \quad \left. \begin{array}{l} m \\ \text{No. of training set} \end{array} \right\} n$$

For linear regression or polynomial regression
mapped to linear regression

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

use for $n \approx 10,000$

No feature scaling required.
No iterations straight answer

* Normal equation

* for Linear Regression

Intuition

$$v = y - \theta^T x$$

We want to minimise

∇v that's it.