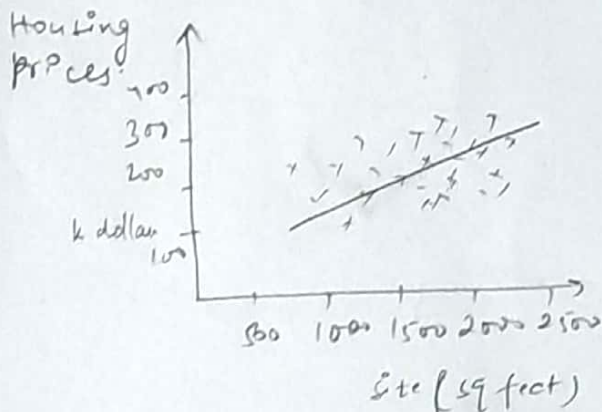


⇒ Model Representation

Supervised learning - Given the 'right answer' for each example in data

↳ Regression problem → Predict real value output

↳ Classification → Discrete-valued output



→ supervised learning, we have training set.

Size (x)	Price (y)
2109	460
1416	232
1534	315
	178

} m=4

m → No. of training examples

x's → "input" variable / features

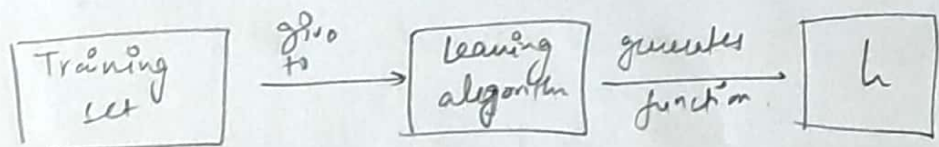
y's → "output" variable / "target" variable.

(x, y) → one training example

(x⁽ⁱ⁾, y⁽ⁱ⁾) → ith training example.

for ex $x^{(2)} = 1416$
 $y^{(2)} = 232$ dk.

working



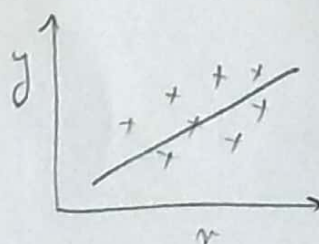
Size of house (x) → h (hypothesis) → estimate price (estimated value of y).

↑
prediction function generated by algorithm

So, h maps x's to y's

→ Representing h .

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



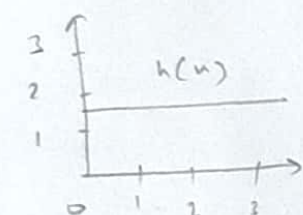
$$h(x) = \theta_0 + \theta_1 x$$

Linear regression with one variable
 Univariate linear regression.

⇒ Cost function → helps to determine best possible straight line for data.

Hypothesis: $h_0(x) = \theta_0 + \theta_1 x$

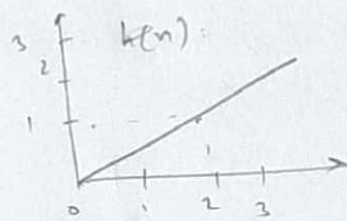
θ_i 's → Parameters, How to choose θ_i 's!



$$\theta_0 = 1.5$$

$$\theta_1 = 0$$

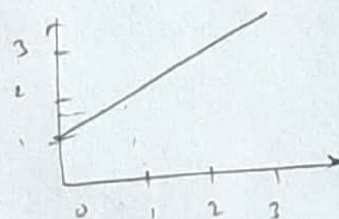
$$h(x) = 1.5$$



$$\theta_0 = 0$$

$$\theta_1 = 0.5$$

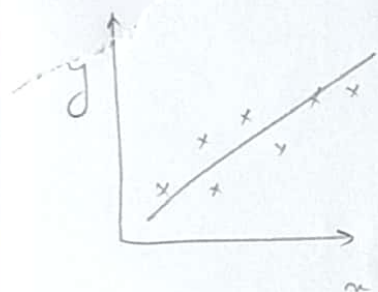
$$h(x) = 0.5x$$



$$\theta_0 = 1$$

$$\theta_1 = 0.5$$

$$h(x) = 1 + 0.5x$$



we need to find θ_0, θ_1 so that $h_0(x)$ is close to y for our training examples (x, y) (determine best fit)

So, we need to solve a minimisation problem

minimise θ_0, θ_1

such that $(h_0(x) - y)^2$ is small.

we need to sum over the training set.

So minimise $\frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$ w.r.t θ_0, θ_1 .

why? $\frac{1}{2m}$ is average.

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

we need to minimise $J(\theta_0, \theta_1)$ } → cost function

$\left(\frac{1}{2}\right)$ is put just as a convenience for computation of gradient descent.
 → squared error function.

Cost function Intuition

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

θ_0, θ_1
parameters

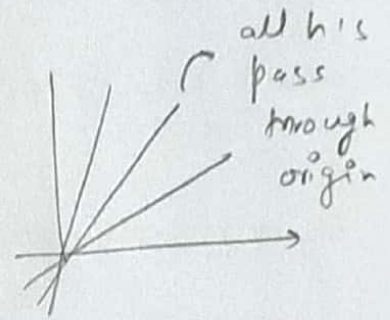
cost function.
 $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$

goal \rightarrow minimize $(J(\theta_0, \theta_1))$

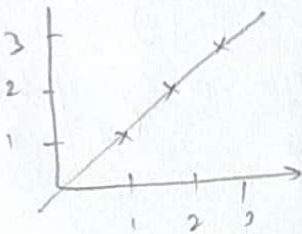
let's take $h_0(x) = \theta_1 x$ i.e. $\theta_0 = 0$

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

minimize wrt θ_1



$h_0(x) \rightarrow$ for fixed θ_1 , this is a fn of x



Training set $(1,1), (2,2), (3,3)$

$\theta_1 = 1$

let's calculate $J(\theta_1)$ for $\theta_1 = 1$

$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$

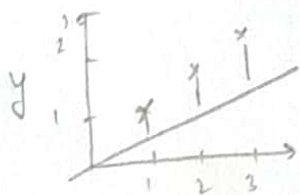
$= \frac{1}{2m} \times 0 = 0$

$J(1) = 0$

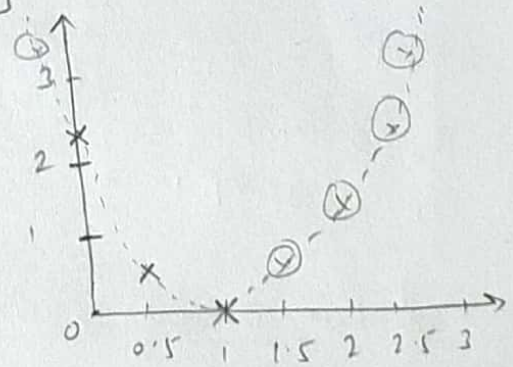
So $J(1) = 0$

let's plot it over J and θ curve J

let's take $\theta_1 = 0.5$



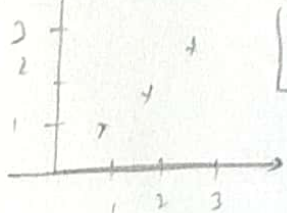
$J(0.5) = \frac{1}{2m} ((0.5-1)^2 + (1-2)^2 + (1.5-3)^2)$
 $= 0.58 = J(0.5)$



$J(0) = \frac{1}{6} (1^2 + 2^2 + 3^2)$

$\frac{14}{2} = \frac{7}{1}$

$\theta_1 = 0$



$J(0) = \frac{7}{2} = 2.35$

So, each value of θ_1 corresponds to a different straight line $h_0(x)$ and thus a different J is obtained

Our objective is to minimize $J(\theta_1)$ and find values of θ_1 so that J is minimum.

We can see that $\theta_1 = 1$ gives best possible str. line fit.

⇒ Cost function intuition - 2

$$\theta_0 = 1 + \sqrt{1.2} \quad 1 + \sqrt{2}$$

$$\theta_1 = 2 + \sqrt{2}$$

⇒ Gradient Descent

↳ Here, we will use GD to minimize an arbitrary function J .

Q. we have $J(\theta_0, \theta_1)$

want $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

Outline → start with some θ_0, θ_1

→ keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$ until we hopefully reach at minimum.

say $(\theta_0 = 0, \theta_1 = 0)$.

GD algorithm

Assignment operator $|$ Truth table
 $a := a + 1$
 $a = b$ ✓
 $a = a + 1$ ✗

repeat until convergence. {

$$\theta_j := \theta_j - \alpha \frac{d}{d\theta_j} J(\theta_0, \theta_1) \quad (\text{for } j=0 \text{ and } j=1)$$

* simultaneously update θ_0 and θ_1 .

} learning rate (defines step size).

correct: simultaneous update

$$\text{temp}_0 := \theta_0 - \alpha \frac{d}{d\theta_0} J(\theta_0, \theta_1)$$

$$\text{temp}_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp}_0$$

$$\theta_1 := \text{temp}_1$$

Incorrect

$$\text{temp}_0 := \theta_0 - \alpha \frac{d}{d\theta_0} J(\theta_0, \theta_1)$$

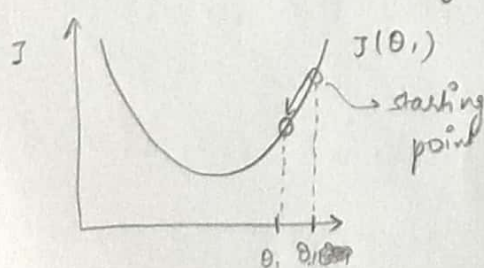
$$\theta_0 := \text{temp}_0$$

$$\text{temp}_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_0, \theta_1)$$

$$\theta_1 := \text{temp}_1$$

⇒ Gradient Descent Intuition

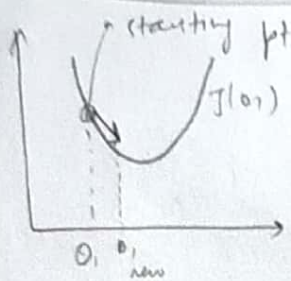
Let's take example of one parameter θ_1 , $J(\theta_1)$. $\theta_1 \in \mathbb{R}$.



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

positive on starting value.

$$\theta_1 := \theta_1 - \alpha (\text{+ve number})$$



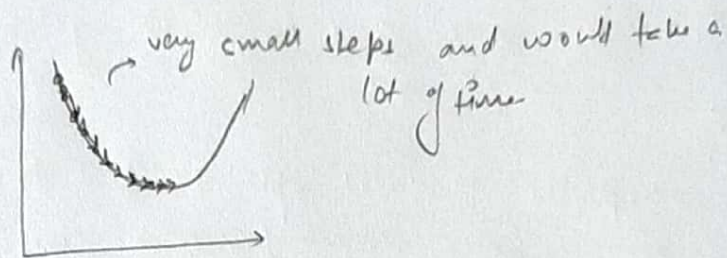
$$\theta_1 := \theta_0 - \alpha \frac{dJ(\theta_0)}{d\theta_0}$$

$\frac{dJ(\theta_0)}{d\theta_0}$ -ve in this case.
total time.

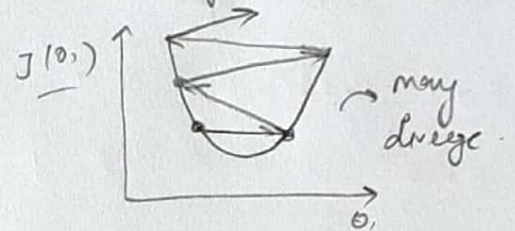
So, θ_1 is greater than previous value.

$$\rightarrow \theta_1 := \theta_0 - \alpha \frac{dJ(\theta_0)}{d\theta_0}$$

If α is too small, gradient descent can be slow.

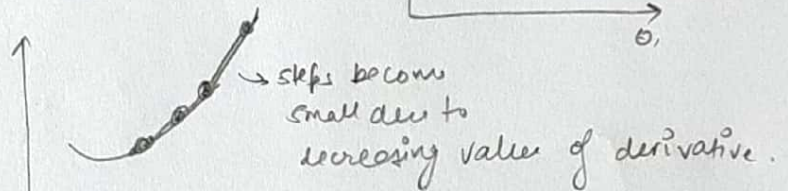


\rightarrow If α too large, GD can overshoot minimum. It may not converge or even diverge.



\rightarrow GD can converge to a local minimum, even with learning rate α fixed.

$$\theta_1 := \theta_0 - \alpha \frac{dJ(\theta_0)}{d\theta_0}$$



\rightarrow So as we approach min, GD automatically takes smaller steps. So we don't need to change α .

\Rightarrow Gradient Descent for Linear Regression

GD algorithm

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{dJ(\theta_0, \theta_1)}{d\theta_j}$$

(for $j=1$ and $j=0$)

}

we need to figure out the partial derivative term now.

$$\frac{dJ(\theta_0, \theta_1)}{d\theta_j} = \frac{d}{d\theta_j} \left[\frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2 \right]$$

$$= \frac{d}{d\theta_j} \left[\frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2 \right]$$

$$\theta_0, j=0; \quad \frac{dJ(\theta_0, \theta_1)}{d\theta_0} = \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})$$

$$\theta_1, j=1; \quad \frac{dJ(\theta_0, \theta_1)}{d\theta_1} = \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) \cdot x^{(i)}$$

Now, we plug these into GD algo

Linear regression model

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

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

