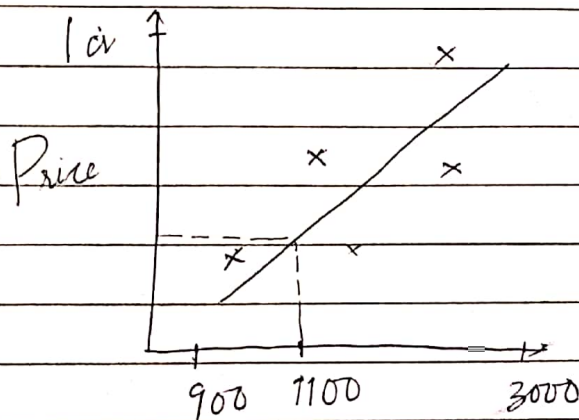


Linear Regression

Problem - House Price Prediction
 $n=1 \rightarrow$ input features

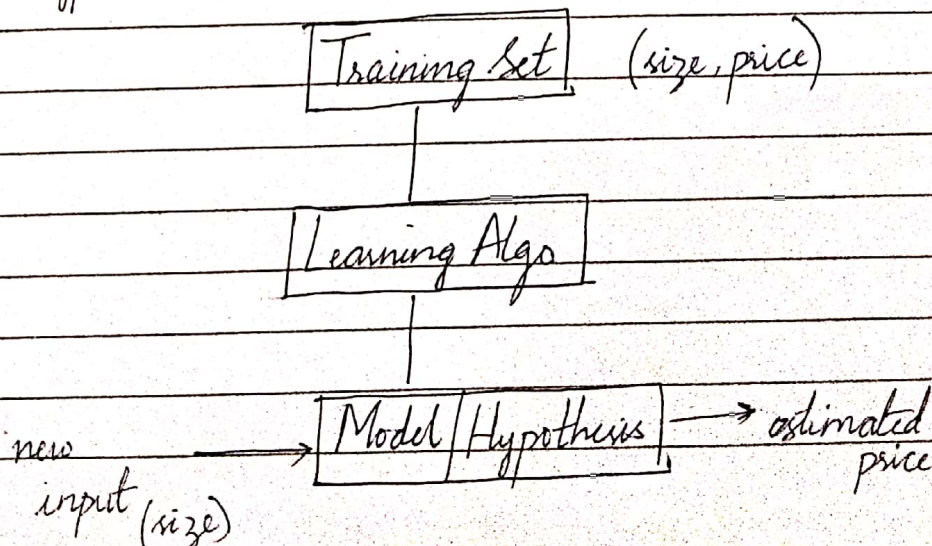
Input	Size (Sq.ft)	Price (lakhs)	Output / Target Variable
$m=5$ training examples	900	12	
	1200	20	
	1500	30	
	1800	40	
	2400	50	



Trying to find the best fit line.

Notations : Training set (x, y)
input/features (x)
output/target variable (y)
 $m = \#$ training examples

Hypothesis / Model



Representation of our Hypothesis

$$y = mx + c \quad (\text{equation of line})$$

\uparrow slope \uparrow intercept

$$h(x) = \overset{\text{bias}}{w_0} + w_1 x_1 \rightarrow \text{for single variable}$$

$$h(x) = w_0 x_0 + w_1 x_1$$

$w =$ Parameters /
Weight
(initialized to
0 / random
value)

$$h(x_1, x_2) = w_0 + w_1 x_1 + w_2 x_2$$

\downarrow in case of two input features (size, BHK)

Generalized \rightarrow
$$h(x) = \sum_{i=0}^n w_i x_i \quad x_0 = 1$$

\uparrow vectorized form of input features

$$w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} \quad X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

Job of learning algorithm is to find best value of
Parameters / Weights

\rightarrow find/choose w such that $h(x) \approx y$ for the given training set

$$\rightarrow \text{minimize } J(w) = \frac{1}{2} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$$

Cost Function

loss \rightarrow single data point / example

cost \rightarrow over all the data points / batch

\rightarrow purpose: Identify quality of your result

Loss function

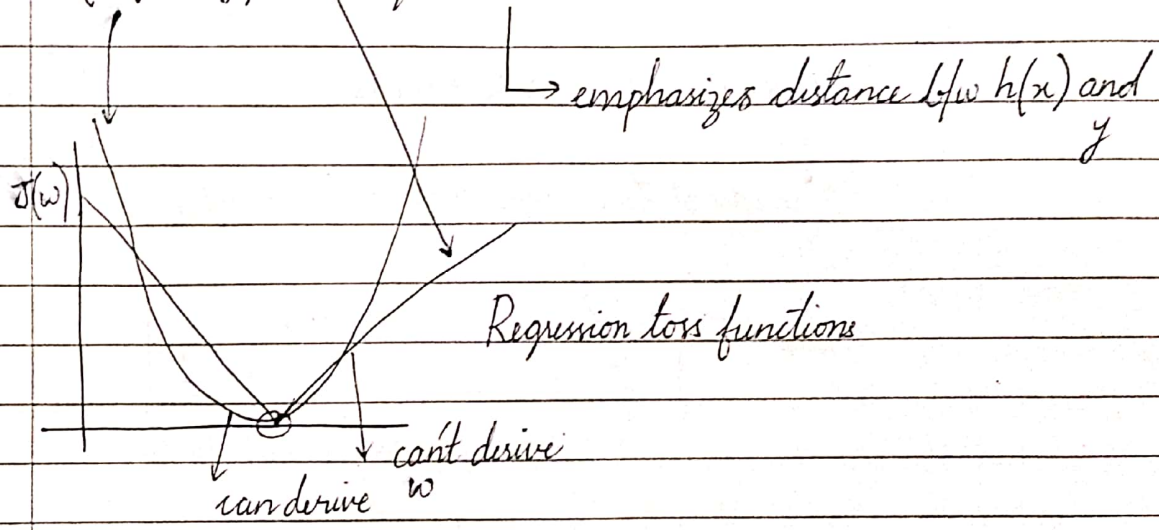
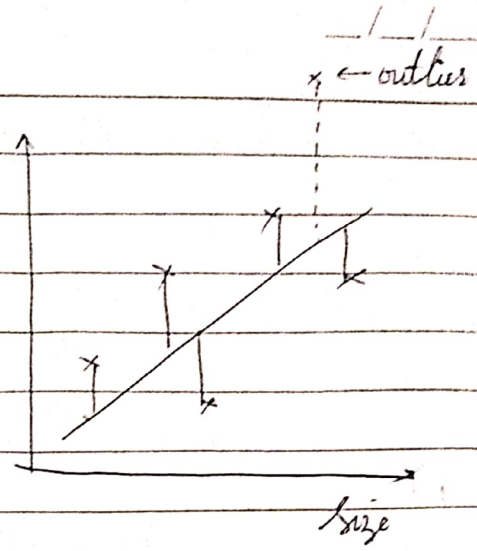
$h(x) \rightarrow$ predicted value
 $y \rightarrow$ ground truth / original value

$h(x) - y$ - margin

$|h(x) - y|$ - absolute loss

$(h(x) - y)^2$ - squared loss

Price



Loss and Accuracy (difference)

y	predicted o/p ($h(x)$)		if $h(x) \geq 0.5$ = 1
0	0.2	0	
1	0.8	1	
1	0.3	0	else = 0
0	0.1	0	
1	0.9	1	

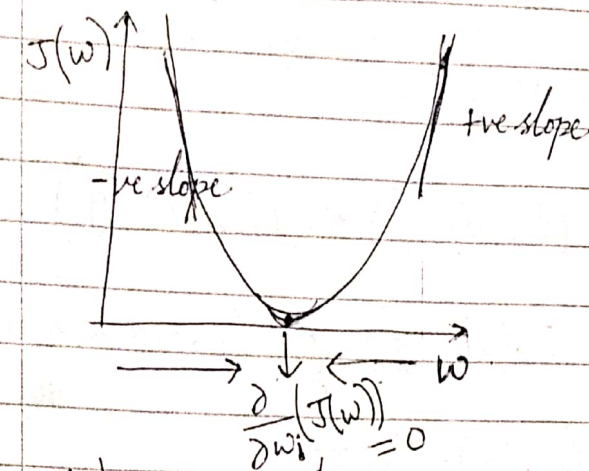
$4/5 \times 100 \rightarrow$ Accuracy (Classification)
 $(0.2 - 0)^2 + (0.8 - 1)^2 + \dots \rightarrow$ loss (Regression)
 (consistency)

Learning Algo: Gradient Descent.

$$y = (3x_1 + 2x_2)^2$$

$$\frac{\partial y}{\partial x_1} = 2(3x_1 + 2x_2)^2 \cdot 3$$

$$\frac{\partial y}{\partial x_2} = 2(3x_1 + 2x_2)^2 \cdot 2$$

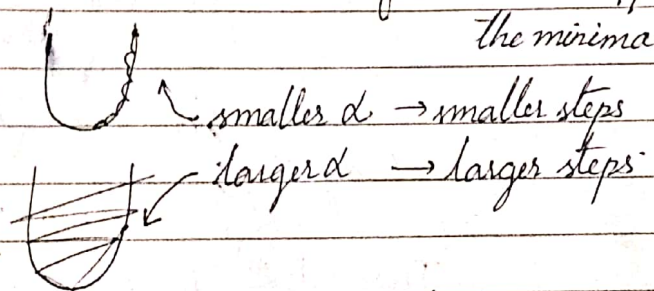


$$w_i = w_i - \alpha \frac{\partial}{\partial w_i} (J(w))$$

partial derivative of $J(w)$
slope at that pt

learning rate.

how fast shall we approach
the minima.



To see if gradient descent is working
 $J(0) \rightarrow$ decrease at each iteration
else adjust α

Summary.

$w \leftarrow$ find a best $w \leftarrow [w_0, w_1, w_2, \dots]$

$$\text{minimize } J(w) = \frac{1}{2} \sum_{i=1}^m (h(x)^i - y^{(i)})^2 \leftarrow \text{Squared loss}$$

$$h(x) = w_0 x_0 + w_1 x_1 + w_2 x_2 + \dots + w_m x_m$$

$$w_i = w_i - \alpha \frac{\partial}{\partial w_i} (J(w))$$

Gradient Descent Update Rule.

$$W = [w_0, w_1, w_2, \dots]$$

$$J(w) = \frac{1}{2} (h_w(x) - y)^2$$

$$\frac{\partial J(w)}{\partial w_i} = \frac{2 \times 1}{2} (h_w(x) - y) \frac{\partial}{\partial w_i} (h_w(x) - y)$$

$$= (h_w(x) - y) \cdot \frac{\partial}{\partial w_i} (w_0 x_0 + w_1 x_1 + \dots + w_i x_i + \dots + w_n x_n - y)$$

$$\boxed{\frac{\partial J(w)}{\partial w_i} = (h_w(x) - y) \cdot x_i}$$

$$w_0 \rightarrow w_0 = w_0 - \frac{\partial}{\partial w_0} (J(w))$$

$$w_0 = w_0 - (h_w(x) - y) \cdot x_0$$

Generalizing \rightarrow

$$w_i = w_i - \alpha (h_w(x) - y) x_i$$

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

$$w_i = w_i - \alpha \sum_{i=1}^m (h_w x^{(i)} - y^{(i)}) \cdot x_j$$

size BHK y

$m = \# \text{ training examples - no. of rows}$

$n = \# \text{ features - no. of columns}$

m

ML Life Cycle.

Data \rightarrow Preprocess Data



Split Data (Train Test Validation)



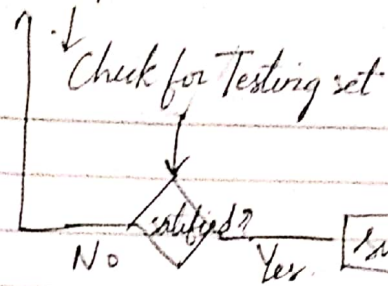
Decide Problem Type (Regression/Classification)



Decide loss function

↓ Decide Hypothesis (Model Structure) \rightarrow LA (gradient descent)

Train Model



high-order polynomial

SVD

Regularization → address high variance problem (overfitting)

L_1 → Lasso Regression → adds 'absolute value of magnitude'

L_2 → Ridge Regression → adds 'squared magnitude'

training accuracy ↑
testing accuracy ↓

$$\text{Loss}(S) = \underbrace{\sum_i \text{Loss}(\hat{y}_i, y_i)}_{\text{minimize}} + \underbrace{\lambda \sum_i |w_i|}_{\text{regularization parameter}}$$

$w \uparrow$ updates weights $|w|$
 $w \downarrow$ regularization parameter (not important) close to 0

L_1 norm $\|w\|_1$

L_2 norm $\|w\|_2$
 $(|w_1| + |w_2| + \dots + |w_n|)$
 $(|w_1|^2 + |w_2|^2 + \dots + |w_n|^2)$

→ Vanilla (Batch) G.D → slower but smooth

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

$$\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i) x_i$$



→ Stochastic G.D.

for i in range(m): → faster but uneven
 $w_i = w_i - \alpha (\hat{y}(i) - y(i)) x_i$ (zig-zag) noise

→ Mini Batch Gradient Descent

for i in range($0, m, 1000$): if mini batch

$$w_i = w_i - \alpha \sum_{j=0}^{p-1} (\hat{y}(i) - y(i)) x_i \quad \text{size} = 1000$$

instead of single sample in SD → we take n data points in each iteration

