

120

~~WOW THE WORLD THIS FOR DAY~~

m-estimate = $\frac{nc}{n+m}$ soch aur $\frac{nc+m}{n+m}$ probability of m-estimate
 but agar koi arist nahi karta, toh zero
 ho jaoge ga

Laplaces smoothing: khetay uniform probability hai (when we have no proper knowledge of underlying distribution)

$$P(\cdot) = \frac{nc+1}{n+m} \quad \text{to avoid } \div \text{ by zero}$$

ab continuous values koh liye kday
 Karen:

→ use Gaussian Naive Bayes:

Univariate gaussian $P(X|Y; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\frac{(x-\mu)^2}{\sigma^2}}$

weights

nakalay keliye

$$P(X|Y; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\frac{(x-\mu)^2}{\sigma^2}}$$

if conditional independence

nhi karta toh wo maine nahi.



Probability of Play Tennis

$\langle \text{Sunny}, \text{cool}, \text{high}, \text{strong} \rangle$
 $n=4, t=2$

$n \times t = 8$, 2 priors (Prob of yes
 & Prob of no)

Labels can
 be on chart

$$P(\text{Yes}) = 9/14$$

$$P(\text{No}) = 5/14$$

$P(\text{Sunny})$

$$\left\{ \begin{array}{l} P(\text{Yes}), P(\text{Sunny}|\text{Yes}), P(\text{cool}|\text{Yes}) \\ 0.005 \quad P(\text{strong}|\text{Yes}), P(\text{strong}|\text{Yes}) \\ P(\text{No}), P(\text{Sunny}|\text{No}) \\ 0.0206 \end{array} \right.$$

if category $P(\text{Sunny}|\text{Yes}) = \frac{P(\text{Sunny} \& \text{Yes})}{\text{Yes}}$

$$\text{toh } P(\text{No}) \leftarrow \frac{0.0206}{0.0206 + 0.005} = \frac{2}{9}$$

but probability
 manage
 to divide
 both
 by sum
 of

$$P(\text{No}) \cdot P(\text{Sunny}|\text{No})$$

$$5/14$$

$$P(a_1, \dots, a_n) = \prod_{j=1}^n \prod_{i=1}^n P(a_i | v_j) \cdot P(v_j)$$

$$V_{MAP} = \arg \max_{v_j \in V} P(v_j | a_1, \dots, a_n) \cdot P(v_j)$$

↳ classes kept

$$P(V_{MAP}) = \arg \max_{v_j \in V} P(a_1, \dots, a_n | v_j) \cdot P(v_j)$$

count each find

$P(2^n \rightarrow \text{for binary classification})$
 k binary possibilities

$$P(a_1, a_2, \dots, a_n, v_j)$$

balanced half

$$= P(a_1 | a_2, \dots, a_n, v_j) \cdot P(a_2, \dots, a_n, v_j)$$

$$= P(a_1 | a_2, \dots, a_n, v_j) \cdot P(a_2 | a_3, \dots, a_n, v_j)$$

$$= P(a_1 | a_2, \dots, a_n, v_j) \cdot P(a_2 | a_3, \dots, a_n, v_j) \cdot P(a_3 | a_4, \dots, a_n, v_j) \cdot \dots \cdot P(a_n | v_j) \cdot P(v_j)$$

V_{MAP} conditional independence (as the class label)

$$P(a_1 | v_j) = P(a_1 | v_j, a_2, \dots, a_n)$$

V_{MAP} vector of probabilities $\boxed{n \times \text{no. of classes}}$ $\rightarrow \text{no. of attributes} \times \text{targets}$

$$V_{MAP} = \arg \max_{v_j} P(a_1 | v_j) \cdot \dots \cdot P(a_n | v_j) \cdot P(v_j)$$

NB v_j

$$V_{MAP} = \arg \max_{v_j} P(v_j) \cdot \prod_{i=1}^n P(a_i | v_j)$$

Bayesian rule lagao ga \rightarrow gen model location
this

$$P(Y|X) = \frac{P(X|Y) \cdot P(Y)}{P(X)}$$

discriminative likelihood
model frame
form

Naive Bayes Classifier Posterior Prob
 \rightarrow generative model

$$V_{\text{map}}^{\text{label}} = \underset{v_j \in V}{\operatorname{argmax}} P(v_j | a_1, \dots, a_n)$$

(6th chapter)

before bayes rule

a_1	a_2	\dots	a_n	v_j
-------	-------	---------	-------	-------

a_0 yahan par
nhi add hoga

j is class keliye max no. aayega.
wo class hai

\rightarrow after bayes rule

$$P(V_j | a_1, a_n) \underset{v_j \in V}{\operatorname{argmax}} \frac{P(a_1, \dots, a_n | v_j), P(v_j)}{P(a_1, \dots, a_n)}$$

RP
 $V_{\text{map}} \rightarrow$ class
hai
tak
confidence
label

normalize
jaise
keliye

(\hookrightarrow dont need

Date: 1/12/20

MON TUE WED THU FRI SAT
000000

divide data into sets
(k-fold cross-validation)

avg performance
kappa validation error increase
avg no. of epochs, wss pattern
Karin.

Naive Bayes Classifier

goal of ML supervised learning?

decision boundary determines KNN

$$f: \mathbb{X} \rightarrow \mathbb{Y}$$
$$h(\mathbf{x}) = \begin{cases} 1 & \text{if } \dots \\ 0 & \text{else} \end{cases}$$

Discriminative Models

- Models that find boundaries
- trying to find class label of given data.

$$P(\mathbf{y}|\mathbf{x})$$

(directly learns
Vikay)

Indirectly \Rightarrow generative models

\hookrightarrow they learn

$$P(\mathbf{x}|\mathbf{y}) \cdot P(\mathbf{y})$$

~~AP~~
~~discrete~~
~~model~~

~~BT AD~~ \rightarrow Re-estimate beta, prior prob?

Generative model $P(\mathbf{x}|\mathbf{y}) \cdot P(\mathbf{y})$

E |
no of epochs

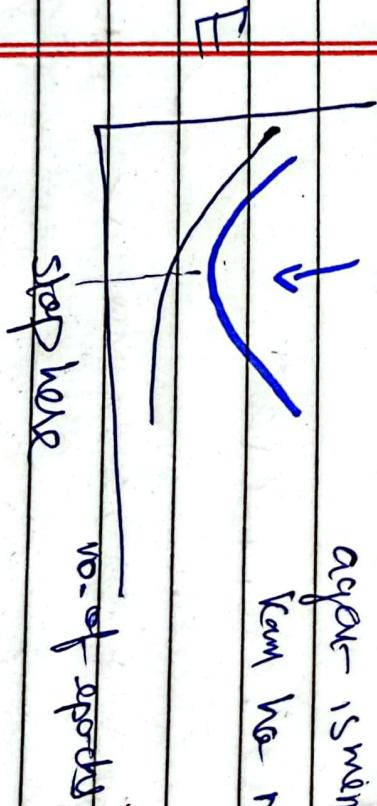
model hi iss koi bhi
training mein sahi ha chalein
 \rightarrow High bias

testing error high
 \rightarrow High variance
(may infer more have done
anything)

find a sweet spot bw high bias
& high variance

validation error

error is men koi
kam ha raha error



stop here

no. of epochs

data is an asset

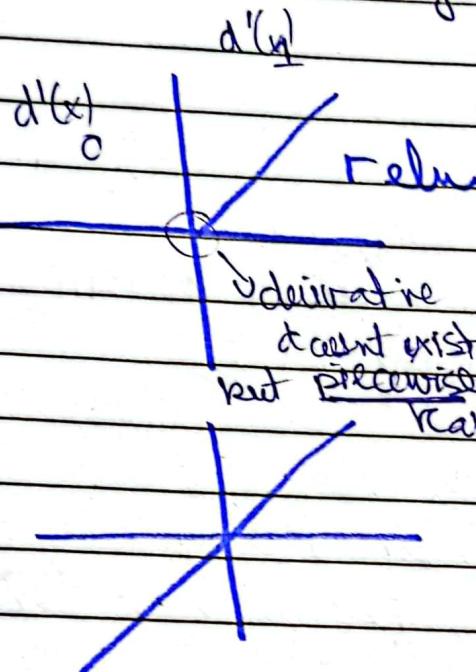
agar koi data hon to:

Overs - validation

iterations ziyada tak minimize
with time

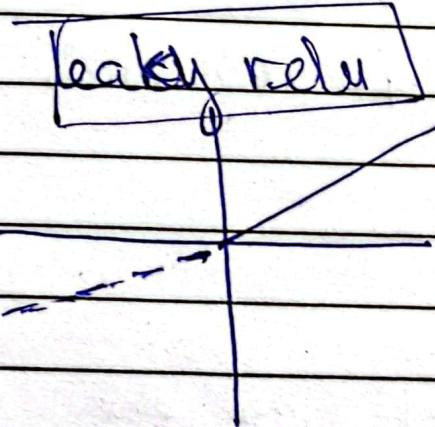
C)

plot heuristic for IC to see where
error bda (askay heuristic)



derivative
doesn't exist
but piecewise
function use
really hoga, uska
derivative hoga

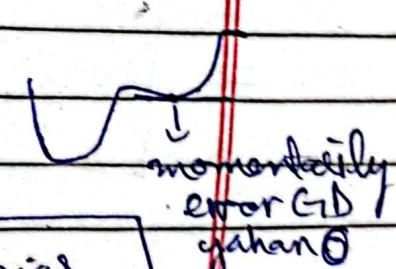
bottleneck wala
hali



EBPA ka kaha nikna :

stopping condition:

- error gradient kam ho jayein
- epoch r
- weight change ziyadi [high-bias]
nhi ho rhi
- * - validation keh zaribe
control karan



model itna complex
nhi keh it can't understand
your data

O mein karen

(v) weights ki initialization

rules for initialization

 $f(w^T x)$ $w^T x$

keep weights in range

keh +, - bhi

ho sakte hain range

ho

eg:

[-1.5, 1.5]

avoid 0

ng bhi zigzag
dps likhein

agel sochit small

hao (near 0)

toh not activation

main time lagayga

$$\text{std of weight } \sigma_w = \frac{1}{\sqrt{T_m}}$$

↳ no. of connections
to a node

scikit-learn → documentation check

(vi) Learning constant

learning constant ke variable
bana do

change lc :

outer mein bda,
inner main chota

Date: 1/12/20

MON TUE WED THU FRI SAT

Classes kitadaad = no. of nodes in
OR layer

binary classification 0 1 0
0 0 1

one hot encoding

binary encoding

Input layers = no. of features
nodes

hidden layer nodes = $D/\alpha \text{ or } 1/\beta \text{ or } 1/10$

ki range
(fixed formula ah)

hui zi yada whi rakhin
everything ho jayege

generalization kitar ha jati :
generalized performance judge
kija rahi

(iv) Scaling of Inputs

we scaling with GDB.

scaling kora karo, numerical scaling
text ke koi numericals leni padhegi

3 - scale scaling
min-max scaling (eg: normalized)

(required : de-select them) weight amness

1 - individual is scaled magi

approx previous weight change by
current weight change divided
by learning rate

$$\Delta w_t \leftarrow \Delta w_t + \beta \Delta w_{t-1}$$

↓
momentum
factor
(is a fraction)

$$[0,1]$$

hyperparameter
0.95, 0.9

(iii) Antisymmetric O/P function
in hidden nodes $\rightarrow f(-x) = -f(x)$

don't use sigmoid.

+1 we this benefits:

antic
otic
+1
we this
 \rightarrow symmetric about origin
 \rightarrow linear range b/w

-1
 \rightarrow want to avoid
Saturated ranges

$$f(x) = e^x - e^{-x} \rightarrow \text{we divide} +1$$
$$\frac{e^x - e^{-x}}{e^x + e^{-x}} \quad a-1$$

$$f'(x) = 1 - (f(x))^2$$

of
- same positive in sigmoid

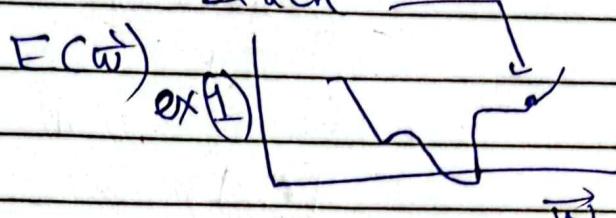
how many nodes in output in MLP
 \rightarrow 0
 \rightarrow 0
 \rightarrow classification seen as 0

batch men hum phasen gray

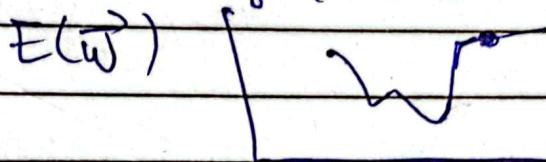
Heuristics

(i) SGD

to sakte aik error surface, we get stuck



but ho sakte ex 2 keh hisab seh
you are here



complicated networke alor so SGD is good in this situation

for batch

but agar error func use Karna, toh skip GD & use:

Newton's method

training example
We network error wrt one effor

(ii) momentum

- adam? solver
- line search
 - ↳ improve method:
conjugate
gradient descent

error surface :

error vs plot kastay wrt weight.

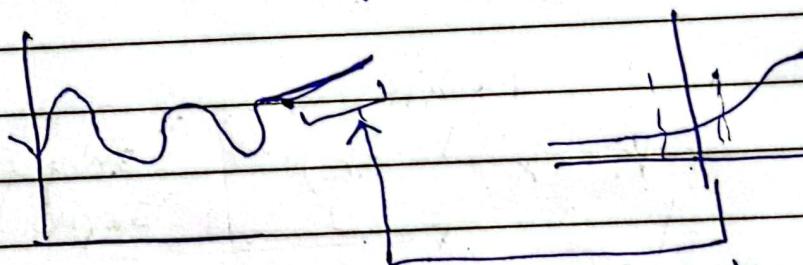
hyper surface

3D. ← more than 3D ka sochna

think
of it
as 3Dwith
billions of parameters

Why we take different error surface :

→ error surface flat nhi



non-convex optimization
model acha as it works,
and is non-linear and
has a lot of free parameters

yahan
ka smooth
generate
kar rakh

error function derived in class
wrt weight that we minimized

$$E(w) = \frac{1}{2} \sum_k (t_k - o_k)^2$$

network error represent kar rakh ley;

$$E(w) = \frac{1}{2N} \sum_{i=1}^N \sum_k (t_{ik} - o_{ik})^2$$

Name Bayes Classifier

review it today
error back propagation:

Heuristics to improve the performance of EBPA:

(It is basically a discriminative classifier)
(for stochastic implementation)

$\downarrow x_i$ Single training set
 i Kth node Venkyai hai

$$\Delta w_{kj} = \eta \cdot \delta_k \cdot o_{kj}$$

δ_k (delta k)

w_{kj}

o_{kj}

d_k

o_k

δ_k

o_k

Date: 1/20

MON TUE WED THU FRI SAT
○ ○ ○ ○ ○

$$Q_0 \quad w_{31} = 0.1 \quad w_{b3} = -0.5$$

$$x_2 \rightarrow 1 \rightarrow 2 \rightarrow 3 \quad w_{21} = 0.2 \quad w_{32} = 1$$

$$x_2 \rightarrow 2 \rightarrow 4 \rightarrow 5 \quad w_{42} = 0.3 \quad w_{54} = 0.5 \quad b_5 = 1$$

$$w_{b4} = 0.3 \quad w_{b5} = -0.25$$

$$\begin{aligned} \psi &= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix} & w_{b4} &= 0.3 \\ (\text{i}) \Delta w_{53} &=? \end{aligned}$$

$$(\text{ii}) S_4 = ?$$

cost creates non-convex function inefficiency (it doesn't align with the probability)

1

DWIT

Decision

2. For each hidden node:

$$\Delta w_{j,i} = \eta s_j \cdot x_{ji}$$

known
calculated
hidden node for
delta weight find
column?

$$s_j - o_j(1-o_j) \sum_{k \in \text{output}} w_{kj} \cdot s_k$$

$$E_d(\vec{w}) = \frac{1}{2} \sum_k (t_k - o_k)^2$$

dth example
of pattern

$$E_d(\vec{w}) = \frac{1}{2} \sum_d \sum_k (t_k - o_k)^2$$

for batch mode: $w_{j,i} + \sum_d \Delta w_{j,i}$

why does stochastic converge faster:

local learning

ex example ex all weight local
minimum mean phas jaise kisi
decreasing mean na phasay

(stochastic kisiyaab hai
practically)

1. $0 \rightarrow 0 \rightarrow 0 \xrightarrow{\text{output}} \text{Lauferer}$

~~options~~ 2.0 → 0 → 0

→ forward propagate the signals

1. For each o/p node:

$$\text{error term } \hat{e}_k = (f_k - o_k) \cdot o_k(1-o_k)$$

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

error
derivative
term

Dukjipriya Singh
Lecturer

or Symmetric
Kth node fully

jth
Connection

kehliger
durch

create non-linear cost function

Date: 1/20 - inefficiency 000000

EBPA ($n_m, n_{in}, n_{op}, n_{nd}, D$)

→ Create a n/w with n_m, n_{op}, n_{nd}
Initially the n/w with small weights
[-0.05, +0.05]

→ keep weight small

w_{ji} with self loop

$\sum_i w_{ji}$ with some cross activation

net_j

but why small so?

initializing?

with 0.0

value

but why 0.0

error change
will note toh

is large O seh

say initializing keh
abrupt bhai gradient
change bhi aaye.

$$o(\text{adj}) = \frac{1}{1 + e^{-\text{net}_j}}$$

more based
fully connected
sigmoidal

Date: / /20

MON TUE WED THU FRI SAT
00000000

$$P(y=c | z) = \frac{e^{z_c}}{\sum_{j=1}^c e^{z_j}}$$

out prob are probabilities

$$L(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{c=1}^C y_i^{(c)} \ln(P(y=c | x^{(i)}))$$

predicted prob
hence
the log likelihood

$$\frac{\partial L(\theta)}{\partial \theta_c} = \frac{1}{m} \sum_{i=1}^m \sum_{c=1}^C y_i^{(c)} - P(y=c | x^{(i)})$$

$$\frac{\partial \ell_i = \ell_i - d \frac{\partial L}{\partial \theta_c}}{\sum_{j=1}^c e^{z_j}}$$

Date: / /20

MON TUE WED THU FRI SAT
00000000

$$\begin{array}{c|c} x & y \\ \hline & \vdots \\ x_1 & \dots x_d & y \end{array}$$

$$\mathbb{R} \rightarrow \mathbb{R}$$

$$D : \mathbb{R}^{d+1} \rightarrow \mathbb{R}$$

$$n > D \quad \text{summed around } y = x^T g$$

$$X \in \mathbb{R}^{n \times (d+1)} \quad (d+1) \times 1$$

$$\begin{aligned} \sigma_j &= y_j - x^T j \\ &= y_j - x^T (x^T x)^{-1} x^T y \\ &= x^T (y_j - x^T y) = 0 \end{aligned}$$

$$\begin{aligned} x^T y_j - x^T y &= 0 \\ x^T y &= x^T y_j \\ x^T y &= x^T x^T y \end{aligned}$$

$$(x^T x)^{-1} \cdot x^T y = 0$$

$$\frac{\partial \sigma_j}{\partial \sigma_j} = 1 \quad \frac{\partial \sigma_j}{\partial \sigma_i} = \frac{\partial \sigma_j}{\partial \sigma_j} + d \left(y_j - h_{\sigma}(x) \right) v_j$$

$$\begin{aligned} \frac{\partial J(\sigma)}{\partial \sigma_j} &= \frac{1}{2} \sum_{i=1}^n \left((h_{\sigma}(x^{(i)}) - y^{(i)})^2 \right) \\ &= y_j - \frac{1}{2} \sum_{i=1}^n v_i (h_{\sigma}(x^{(i)}) - y^{(i)}) \\ &= (h_{\sigma}(k) - y^k) \cdot x_j \end{aligned}$$

d - no of input vars excluding x_0
cost func

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

choose α to minimize $J(\alpha)$

$$\begin{aligned}\alpha_j &= \frac{\partial J(\alpha)}{\partial \alpha_j} \\ &= \alpha_j - \frac{1}{n} \sum_{i=1}^n (h_\alpha(x^{(i)}) - y^{(i)}) x_j\end{aligned}$$

$y - h$ one training ex:

$$\alpha_j = \alpha_j + \alpha (y^{(i)} - h_\alpha(x^{(i)})) x_j^{(i)}$$

\downarrow
 $b - y$

①

$$\alpha_j = \alpha_j + \alpha \sum_{i=1}^n (y^{(i)} - h_\alpha(x^{(i)})) x_j^{(i)}$$

for every

$R \times d \rightarrow R$

$$X \rightarrow n \times (d+1)$$

$$\begin{aligned}\nabla_\alpha J(\alpha) &= X^T X \alpha - X^T \bar{y} \\ 0 &= X^T X \alpha - X^T \bar{y}\end{aligned}$$

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

$$\alpha = (X^T X)^{-1} \cdot X^T \bar{y}$$

$$\frac{x^T x}{x^T x} = \alpha$$

Legend of
variables:
1. X - input
2. α - weight
3. b - bias
4. y - output
5. \bar{y} - average output
6. α_j - weight of j-th feature
7. x_j - j-th feature
8. $x_j^{(i)}$ - j-th feature of i-th training example
9. $y^{(i)}$ - output of i-th training example
10. h_α - hypothesis function

$$! \times ((x_{\theta}(y) - h) =$$

$$\frac{\partial e}{\partial x} (x_{\theta}(y)(h-1) - ((x_{\theta}(y)-1) - f_{\theta}(x)) =$$

$$(x_{\theta}(y) - \frac{f_{\theta}(x_{\theta}(y)-1)}{1}) (h-1) - \frac{f_{\theta}(x)}{1} - h = \frac{\partial e}{\partial \theta}$$

$$((1, x_{\theta}(y)-1)) (y_{\theta}(h-1) + (1) \log(y_{\theta}(h-1))) \sum_{t=1}^T$$

$$(x_{\theta}(y) - 1) \log(y_{\theta}(h-1)) +$$

$$L(\theta) = \log \left(\prod_{i=1}^n y_{\theta}(h_i(x_i)) \right)$$

$$(1 - y_{\theta}(h_i(x_i)))^{-1} \times y_{\theta}(h_i(x_i))$$

$$\begin{aligned} & \frac{(z^b - 1)}{(z^{d+1} - 1)} \cdot \frac{(z)^b}{\frac{(z^{d+1})}{1}} = \\ & \frac{z^b}{(z^{d+1} - 1)} = \\ & \left(\frac{z^b}{1} \right) \frac{z^b}{p} = (z)^b \end{aligned}$$

$$\begin{aligned} & f(x) = x_0 + x_1 z + x_2 z^2 + \dots \\ & f(x) = x_0 + 0 = x_0 \\ & \text{OR } (x)_y = 0 \leftarrow z \\ & T \leftarrow (x)_y : 00 \leftarrow z \end{aligned}$$

→ read result 1,0 m/g or
 $h_6(y) = (x)_y$

$$z = x + t$$

~~to calculate product by~~



using ↗

$$\frac{(x_{0+0}) - 1 + t}{t} = \frac{x_{1+0} - 1 + t}{t} = (x)_y$$

function
signature
implementation

$$(1)x \cdot ((1)x^0 - h_6(x)) \sum_{w=1}^{t-1} = \frac{a_w}{(t)_y}$$

$$((1)x^0 - h_6(x)) \sum_{w=1}^{t-1} = \frac{a_w}{(t)_y}$$

$$\frac{a_w}{(t)_y} + b_0 = ?$$

$$\theta_0 = 0.2$$

$$\theta_1 = 0.5$$

$$h_\theta(v) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 v)}}$$

$$h_\theta(0.5) = \frac{1}{1 + e^{-(0.2 + 0.5 \cdot 0.5)}}$$

$$\Rightarrow \frac{1 - e^{-0.45}}{1}$$

$$= \frac{0.5810659}{0.4189341} \approx 1$$

$$h_\theta(0.0) = \frac{1}{1 + e^{-(0.2 + 0.5 \cdot 0)}}$$

$$= \frac{1 + e^{-0.7}}{1 + e^{-0.7}}$$

$$h_\theta(4.5) = \frac{1}{1 + e^{-(0.2 + 0.5 \cdot 1.5)}}$$

$$\Rightarrow \frac{1 - e^{-0.45}}{1}$$

$$= \frac{0.551021}{0.4489789} \approx 1$$

$$= \frac{1 + e^{-0.45}}{1 + e^{-0.45}}$$

$$a) L(\theta) = y \cdot \log(h_\theta(x)) + (1-y) \cdot \log(1-h_\theta(x))$$
$$= 1 \cdot \log(0.616) + (1-1) \cdot \log(1-0.616)$$
$$= -0.7131$$

$$= 0.7131 \cdot \ln 0.3849 - 0.494$$

$$b) L(\theta) = y \cdot \log(h_\theta(x)) + (1-y) \cdot \log(1-h_\theta(x))$$
$$= 0 \cdot \log(0.668) + (1-0) \cdot \log(1-0.668)$$
$$= 0 + \log(0.332)$$

$$= -1.590744$$

$$c) L(\theta) = y \cdot \log(h_\theta(x)) + (1-y) \cdot \log(1-h_\theta(x))$$
$$= 1 \cdot \log(0.750) + (1-1) \cdot \log(1-0.750)$$
$$= -0.415037 + 0$$

$$= -0.415037 + 0$$
$$= -0.415037$$
$$= -0.415$$

Date: 1/120

1 - out put is cleared to
00000000
MON TUE WED THU FRI SAT SUNDAY

→ Back propagate error

- O/P layer

$$\text{error } S_j := \Theta_j(1-\Theta_j)(t_j - \Theta_j)$$



gradient

of O/P function $\Theta_j \rightarrow \text{hypothesis}$

Cohesiveness of whole structure
is two hypothesis

$$[\hat{y}_j - (t_j - \Theta_j)]$$

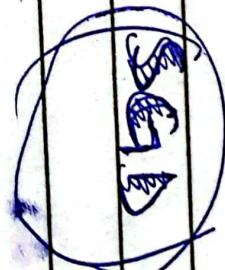
$$A_{w_{ij}} = \eta S_j x_j$$

for out put layer

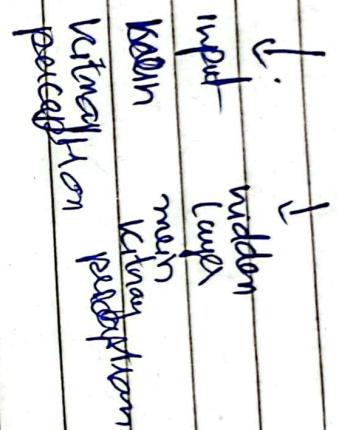
3 : learning • delta • signal
const turn strength

things

Sequential learning
Stochastic learning



EBPACN, n_m , n_h (number of nodes)



sent \Rightarrow lajhs keh lye hai

→ decide the network
(design choice : deciding how many nodes
we need)

hidden layer most important layer

transformation is layer ki dimensi

→ initial weights between small numbers

[-0.05 to 0.05] input layer is ko GI seh

why?

forward propagate inputs layer ja ye ja

how are we going to forward propagate?

weighted + threshold sum

for all neurons

in hidden layer

than in output layer

start off with first hidden layer

larger by larger propagation



use a non-linear continuous threshold
sigmoidal unit have effect
sigmoid

multinomial logistic regression

Softmax function

* back propagation algorithm

↳ direction of transferred

finding errors

by adjusting weights

also called forward propagation

* back propagate
or signals to forward propagate

x₁

f₁

x₂

f₂

input layer

sums of inputs to propagation level are not

(forward propagation)

o scalar o vector

G

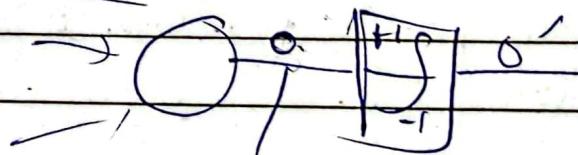
also called: error back propagation algo

Date: 1/10

this will converge till ~~T-O does isn't~~
minimum

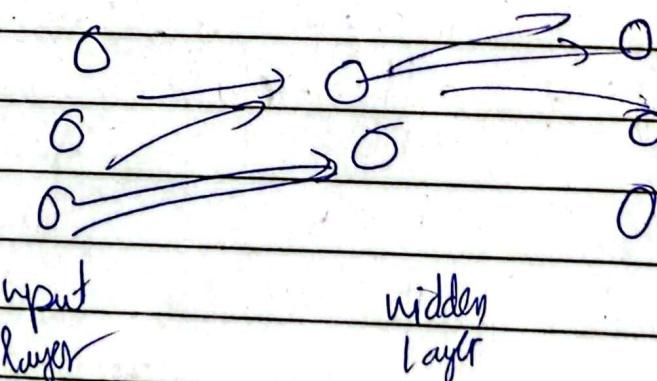
this will converge, but doesn't
mean it will converge & give correct
classification, but will give
you minimum error classification,
but normal perception doesn't converge.
~~converge~~

What is the difference b/w this
& a perceptron?



e.g.: squared features for XOR in logistic
how can we make powerful structures
from it

perception always outputs a line.



that's not a perceptron in MLP
after this

take a higher degree polynomial
square etc

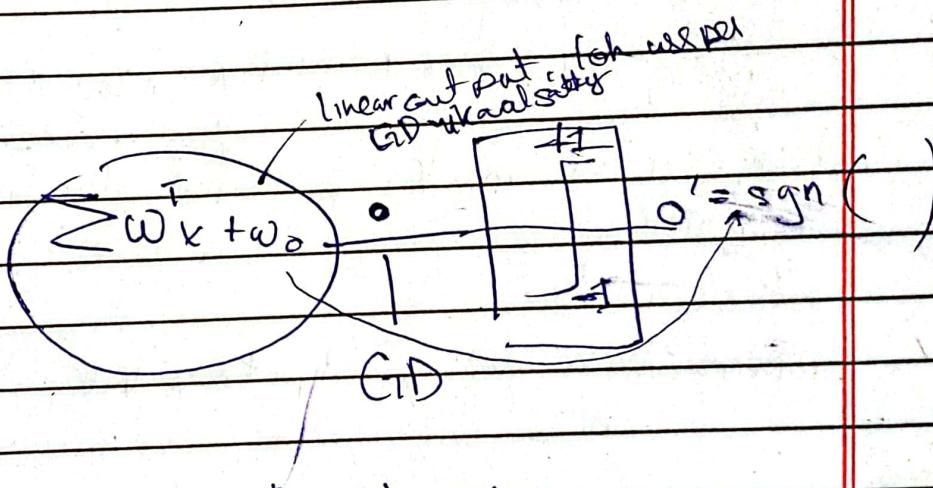
$$\sum w_i x^i + w_0 \quad \text{degree of polynomial}$$

but it's not differentiable
at Sigmoid vein perception ko manipulate
Kai den tak ho sakte?

how bhi line draw karay ga

but in a higher dimension

Can't apply GD to perception, derivative
bhi lena pata, can't take derivative
of step function



$$E(w) = \frac{1}{2} \sum_{i=1}^n (f^{(i)} - o^{(i)})^2$$

simpley

$$\frac{\partial E}{\partial w_i} = ?$$

$$w^T x - o_0$$

$o \rightarrow$ is a real value

keep it near 1 or 0

$$\text{Philosophy: } \Delta w_j \propto \frac{\partial E}{\partial w_i} = \frac{1}{2} \sum_{i=1}^n (f^{(i)} - o^{(i)}) \cdot x_j$$

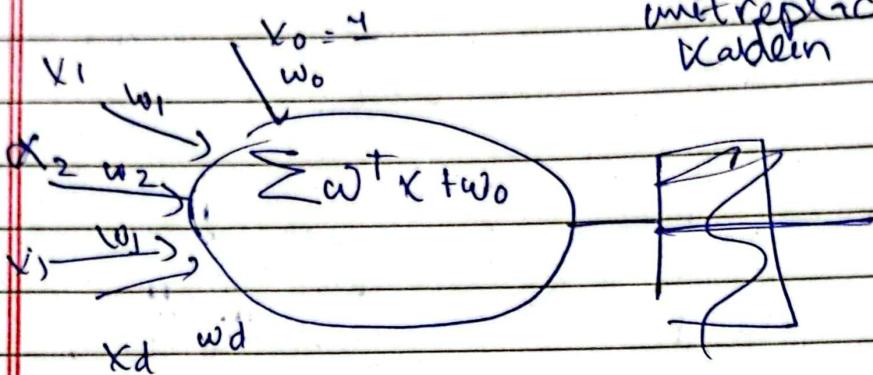
now, what's the difference b/w perception
in logistic regression: difference No softmax
var rates

Date: ___ / ___ / 20 ___

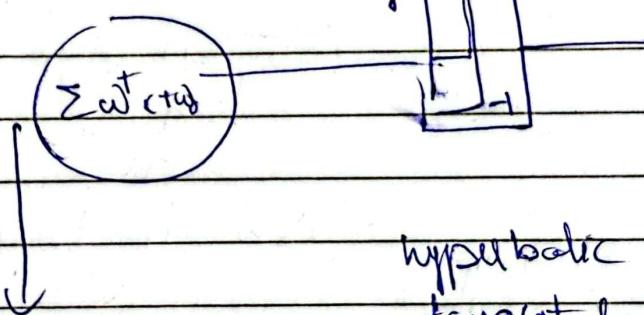
$$\vec{w} \leftarrow \vec{w} + \Delta \vec{w}$$

$x_j \rightarrow$ input vector

agar regression kamai tak perceptron
karna saath koi sajtaay, agar thresholding



for classification:
yeh logistic regression ban jaye !



→ you can only draw a linear boundary

n

while logistic regression can draw a non-linear boundary:

how can we draw a non-linear boundary in logistic regression

perception needs multiple perceptrons
to make non-linear boundary

Date: 1/120

MON TUE WED THU FRI SAT SUN
0000000

2 layer kehsaath you can recognise all patterns. with perception

$\|x\|^2 + 1 \rightarrow$ adding this to weight
jab galti hoti model se
 \downarrow
input vector

$$e = (t - o)$$

\downarrow \downarrow \rightarrow actual output

target

Ch 4 → textbook

$$\Delta w_j \propto e \quad \text{error ziyada}$$
$$\Delta w_j \propto x_j \quad \text{input strength}$$
$$\Delta w_j = \eta (t - o) \cdot x_j \quad \text{ziyada}$$
$$\Delta w_j = \eta (t - o) \cdot x_j$$

↑
eta

way ①: update after every iteration
(instantaneous adjustment)

(online) offline immediately apply it

$$\eta = (0, 1]$$

↑
learning constant

$$\Delta \vec{w} = \eta (t - o) \cdot \vec{x}$$

calculated got

$$+1 - (-1)$$

$$= 2$$

$$12 - (+1)$$

$$e = -2$$

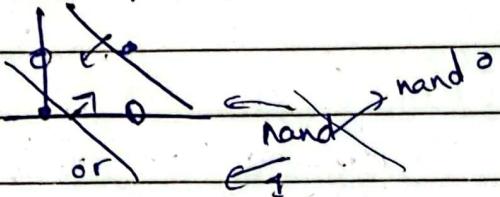
way ②: accumulate hony dan, E-had aik

batch adjustment

vector sumayein &
aggregate suma dan

$$\Delta w^{(1)} + \Delta w^{(2)} + \dots + \Delta w^{(n)}$$

need 3 perceptron to calc make XOR.

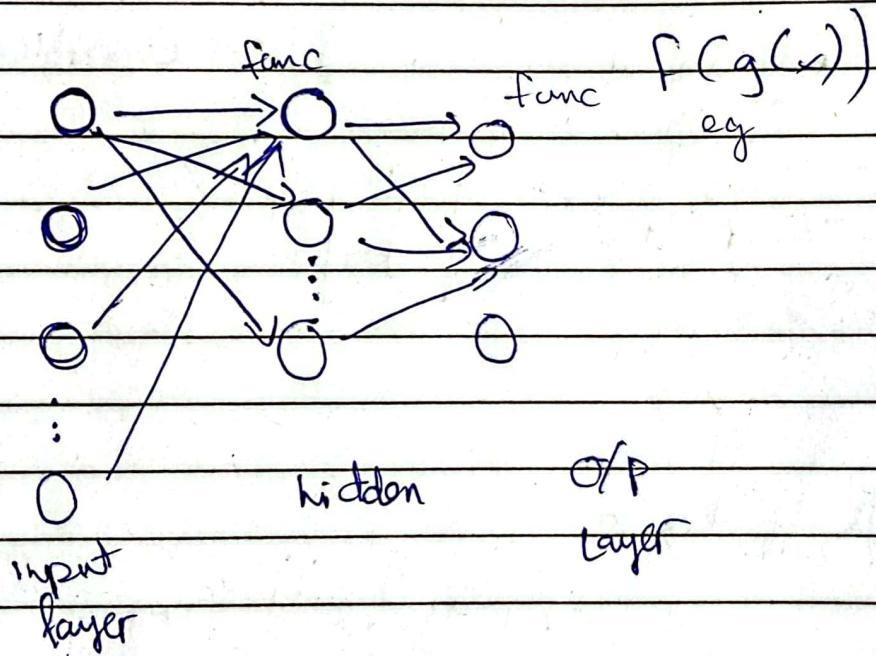


in dono koi firing ko kon decide
karay ga? another perceptron (and)

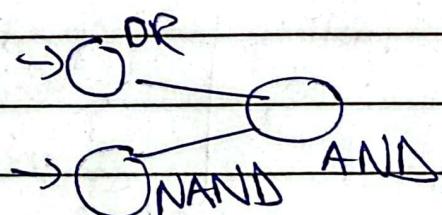
Multi Layered Perception

→ we'll study one only

→ Ifeet forward



Recurrent network: multiple hidden layers



We need a multilayered perception
with learning algorithm

Perception (D_n, T)

$$w = [0 \ 0 \ \dots \ 0]^T; w_0 = 0$$

for $t=1$ to T # of epochs (iterations in GA)

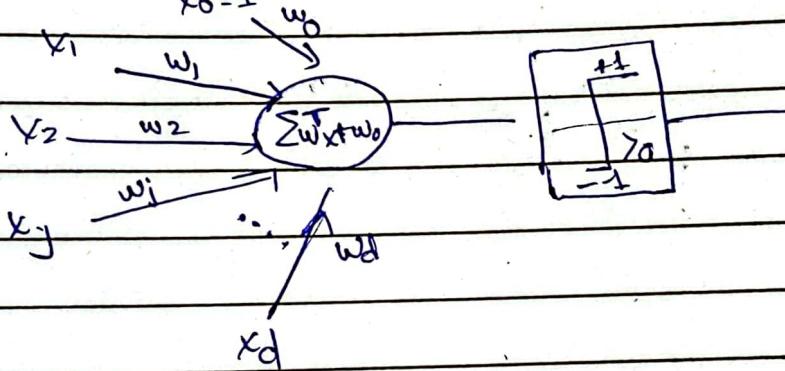
for $i=1$ to n

$$\text{if } y^{(i)} (w^T x + w_0) \leq 0$$

$$w \leftarrow w + y^{(i)} \cdot x^{(i)}$$

$$w_0 \leftarrow w_0 + y^{(i)} \quad (\text{scalar update})$$

return w, w_0



starts off at zero

$$w \leftarrow w + y^{(i)} \cdot x^{(i)} \rightarrow \text{vector}$$

$$w_j \leftarrow w_j + x_j^{(i)} \rightarrow \text{scalar}$$

w_1, w_2, \dots

x_1, x_2

w_d, w_0

x_d, x_0

If a pattern cannot fit the data, then job take some iteration

explore no hope says toh dega, ϵ in the end will give misclassification

$$\begin{bmatrix} = \\ - \\ \vdots \\ = \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} + 1 \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

epochs

Date: 1/20

always include it, is!
 ↑
 MON TUE WED THU FRI SAT
 ○○○○○○○

x_2	.	x_0	x_1	x_2	y
.	.	1	1	3	1
.	.	1	3	1	1
.	.	1	-1	1	-1
.	.	1	1	-3	1

$$1) \text{sgn } (\underline{w^T x})$$

$$w = [1 \ -1 \ 1]$$

$$\begin{bmatrix} 1 & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \\ 1 \end{bmatrix} = 1 - 3 + 1 = -1$$

$$\text{sgn } (\underline{w^T x}) = -1$$

↓
signum function

write this way

$$\begin{array}{|c|c|c|} \hline & x_1 & x_2 & x_0 \\ \hline \end{array}$$

$$(2) 1 \cdot -1 \leq 0$$

(3)

$$\begin{bmatrix} w \cdot \\ 1 \\ -1 \\ 1 \end{bmatrix} + \begin{bmatrix} 1 \\ 3 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix}$$

new weight vector is $\begin{bmatrix} 2 & 2 & 2 \end{bmatrix}$ $\xrightarrow{\text{iter 2}}$

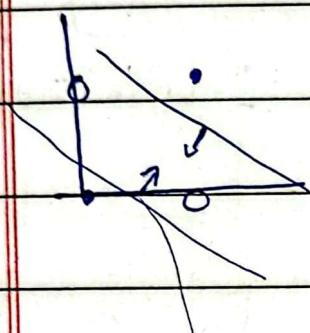
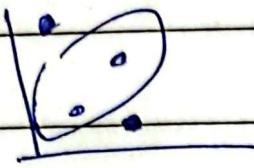
$$\begin{bmatrix} 2 & 2 & 2 \end{bmatrix} \begin{bmatrix} 3 \\ 1 \\ 1 \end{bmatrix} \quad \begin{bmatrix} 2 & 2 & 2 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix} = f_2$$

completely run it
 check where it stops
 if it does

$$-1(2) < 0$$

$0 \rightarrow w$ now
 in reference book

You can use logistic regression to solve this XOR problem, as we can create non-linear boundaries.



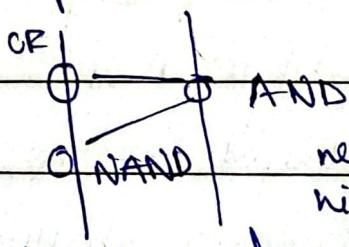
Can solve it easily with 2 perceptions

firing decision j ab done 1, 1-
fire and this class \rightarrow 0

\rightarrow multiple perceptions seh hum koi

whi boundary bana saktey, like circular
C can make any pattern, specially 2 layer
wala perception:

for XOR



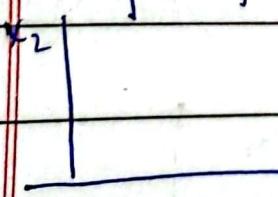
practical difficulty.

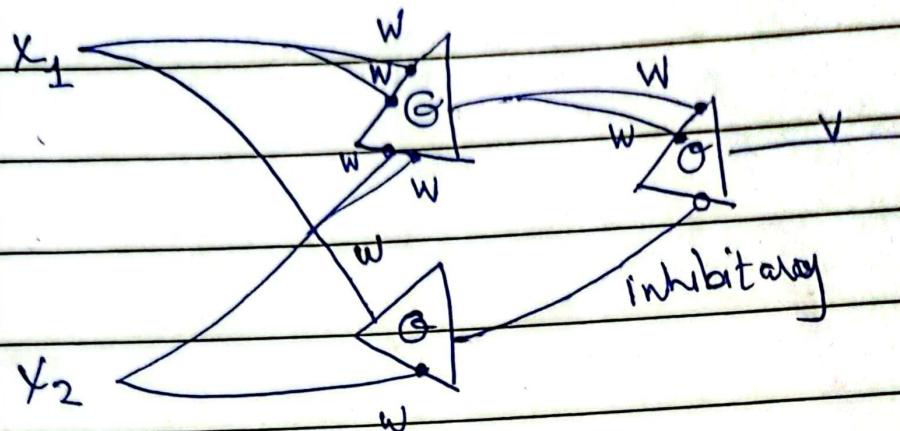
need 2 layers, aka
hidden

\hookrightarrow phir emage wa deep neural networks

X ₁	X ₂	Y
1	3	1
3	1	1
-1	1	-1
1	-3	1

$$W = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_0 \end{bmatrix}$$

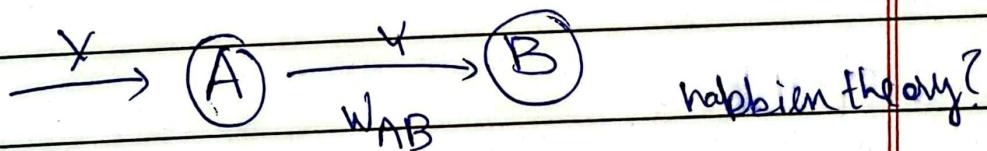




all weights (' w ') = 1 & $\theta = 2$

x_1	x_2	v
0	0	0
0	1	1
1	0	1
1	1	0

* special condition, inhibitory
lock ka dega.



$$W_{AB} \leftarrow W_{AB} + \alpha X Y$$

problem of this model:

* learning hai, but learning is unbounded

* Neuron Networks ki problem: learning
nhi hoti

Perception - Neural Model
- animal eye demonstration

McCulloch Pitt Newton

Solution for XOR

M-P characteristics:

(i) They are binary devices, i.e.

$$O/P \in \{1, 0\}$$

(ii) Each neuron has a fixed threshold ' Θ '

(iii) The neuron receives inputs from excitatory synapses, all with 'identical' wts.

There may be more than one weight from a single source.

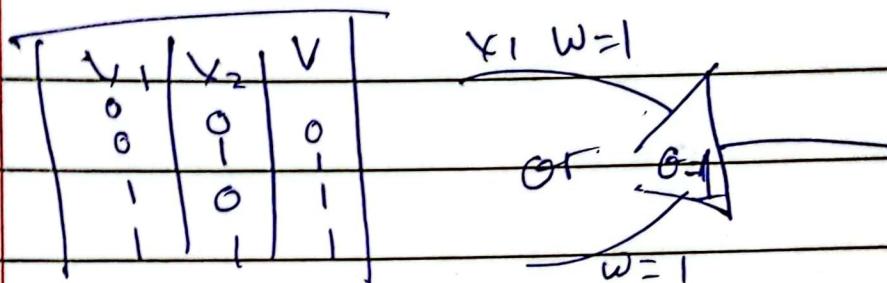
(iv) Inhibitory inputs have an absolute veto power over any excitatory input

(v) At each time step, the neurons are synchronously updated by summing the weighted excitatory inputs & outputting as follows:

$$V = \begin{cases} 1 & \sum Wx_j \geq \Theta \text{ and NO INHIBITION} \\ 0 & \text{otherwise} \end{cases}$$

M-P $V = \begin{cases} 1 & \sum_j w x_j \geq 0 \text{ and there} \\ & \text{is no inhibition} \end{cases}$

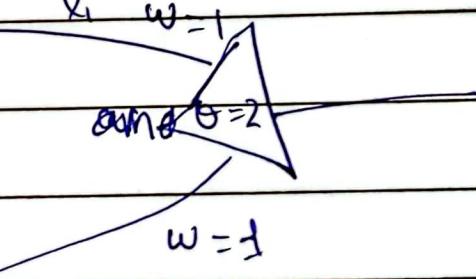
It is equivalent to boolean gate (NOT)



* multinomial logistic regression

convert to and

* XOR ka bana
Koi dikhayein
next time



selection

Quiz: linear regression
at Thursday +
logistic regression

Loss function ka derivative dekhna hai.

Mon Tue Wed Thu Fri Sat

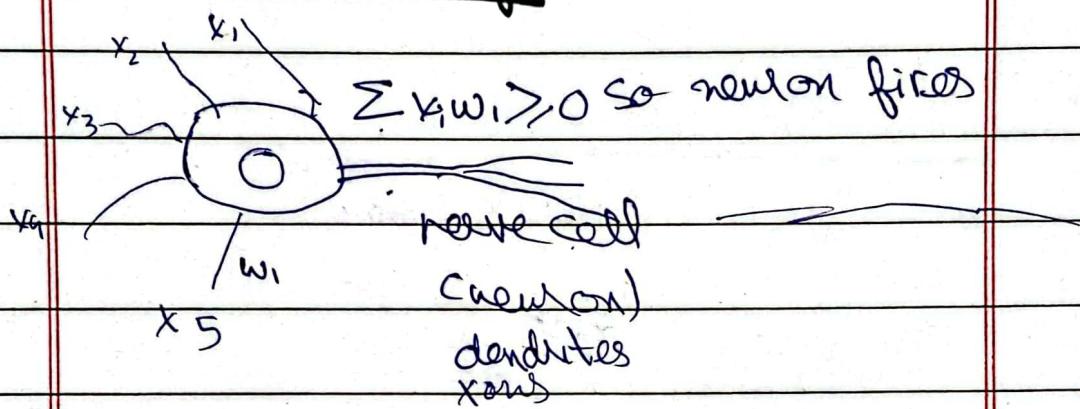
Date: ___/___/20___

$$\frac{d \mathcal{L}(\theta)}{d\theta_i} = \vec{\theta}_i - \vec{\theta}_{\text{target}}$$

$$\vec{\theta} = \begin{bmatrix} \vec{\theta}_1 \\ \vec{\theta}_2 \\ \vdots \\ \vec{\theta}_n \end{bmatrix}$$

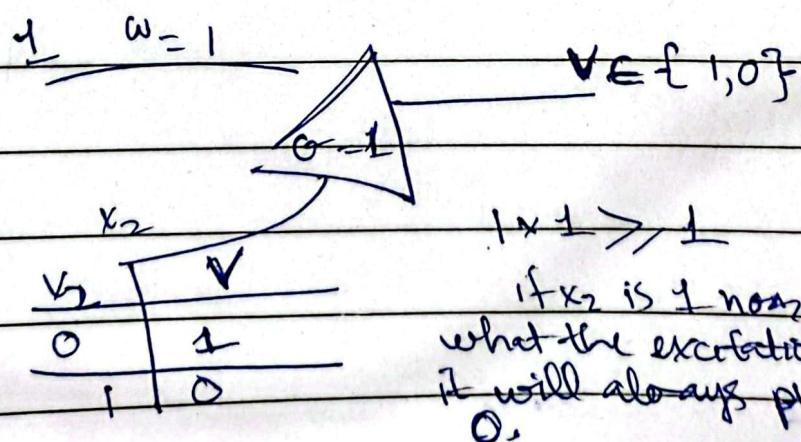
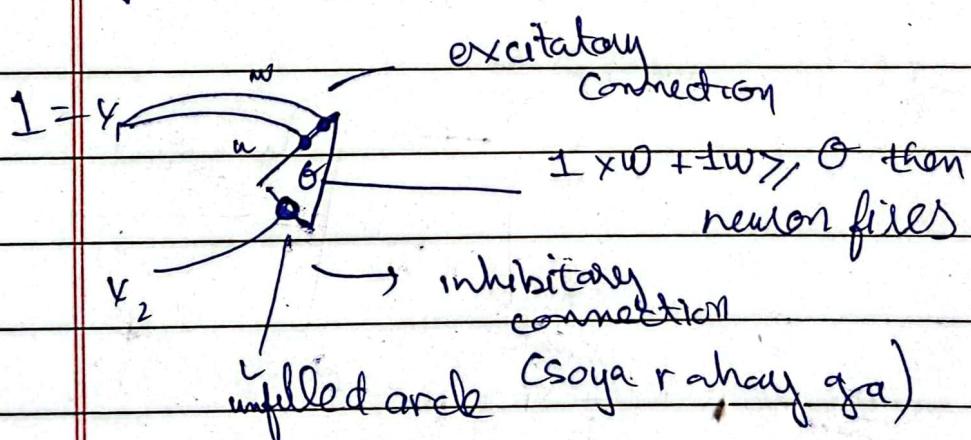
1 - Hot coding

Neural Computing



Unit
Combination
 $= \text{signal} \times \text{wei ght}$

Mechanic Pitts



difference b/w stochastic & incremental

↓ ↓

randomly
at ith value,
you do the next
val

first ith based
2nd stochastic
phr 3rd

also extend to K classes:

$$y \in \{1, 2, \dots, k\}$$

$$\vec{\phi}_1, \vec{\phi}_2, \dots, \vec{\phi}_k \quad \rightarrow \text{ } \} \text{ vectors}$$

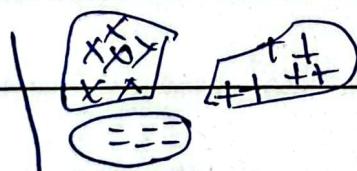
$$\phi_i \in \mathbb{R}^{d+1} \quad \text{k scalars}$$

$$\rightarrow \text{SOFTMAX}(\phi_1^T x, \phi_2^T x, \dots, \phi_k^T x)$$

$$f: \mathbb{R}^k \rightarrow \mathbb{R}^k$$

$$\begin{aligned} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_k \end{aligned} \leftarrow \text{K entries} \quad \frac{e^{\phi_1^T x}}{\sum_j e^{\phi_j^T x}}$$

$$\sum \phi_i = 1 \quad \text{exponential function}$$



3 classes
jth highest probability,
yeh wala class.

$$= \frac{1}{1+e^{-x}} - \left(\frac{1}{1+e^{-x}} \right)^2$$

$$= \frac{1}{1+e^{-x}} \left(1 - \frac{1}{1+e^{-x}} \right)$$

now back to binary cross-entropy:

$$= - \left(\frac{y}{h_\theta(x)} + \frac{(1-y)}{1-h_\theta(x)} \right) h_\theta(x)(1-h_\theta(x))$$

$$= - \left(\frac{y}{h_\theta(x)} - \frac{(1-y)}{1-h_\theta(x)} \right) h_\theta(x)(1-h_\theta(x)) \quad x \text{ is a vector}$$

$$\begin{aligned} J'(\theta) &= - (y(1-h_\theta(x)) - (1-y)h_\theta(x)) x_j \quad \frac{\partial \theta^T x}{\partial \theta_j} \\ &= - (y - yh_\theta(x) - h_\theta(x) + yh_\theta(x)) x_j \quad \frac{\partial (\theta_0 x_0 + \theta_1 x_1 + \dots + \theta_d x_d)}{\partial \theta_j} \\ J(\theta) &= -(y - h_\theta(x)) x_j \quad \theta_{d-1} x_{d-1} + \dots + \theta_1 x_1 + \theta_0 \end{aligned}$$

$$J'(\theta) = (h_\theta(x) - y) x_j \quad = x_j \quad \text{scalar } \theta \text{ s no j entry except } J \text{ wala term}$$

derivative term is exactly
Same as linear regl, but
this has different cost function

$$GD: \quad \theta_j = \theta_j - \alpha (J'(\theta))$$

$$\theta_j = \theta_j - \alpha (h_\theta(x) - y) x_j$$

$$\theta_j = \theta_j + \alpha (y - h_\theta(x)) x_j$$

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

batch made mean summation:

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

$$\therefore \frac{d}{dx} \log x = \frac{1}{x} \frac{d(x)}{dx}$$

Date: 1/120

Mon Tue Wed Thu Fri Sat

$$\begin{aligned}
 &= - \left(\frac{\partial (y \log(h_\theta(x)))}{\partial \theta_j} + \frac{\partial (1-y) (\log(1-h_\theta(x)))}{\partial \theta_j} \right) \\
 &= - \frac{\partial \log(h_\theta(x))}{\partial \theta_j} \\
 &= \frac{1}{h_\theta(x)^k} \cdot \frac{\partial (h_\theta(x))}{\partial \theta_j} \\
 &= - \left(\frac{2(y \log(h_\theta(x)))}{\partial \theta_j} + \frac{2(1-y) (\log(1-h_\theta(x)))}{\partial \theta_j} \right) \cdot \frac{(1-y) \cdot \frac{1}{1-h_\theta(x)}}{1-h_\theta(x)} \cdot \frac{\partial}{\partial \theta_j} (1-h_\theta(x)) \\
 &= - \left(\frac{y}{h_\theta(x)} - \frac{(1-y)}{1-h_\theta(x)} \right) \frac{\partial h_\theta(x)}{\partial \theta_j}
 \end{aligned}$$

$$f(x) - \sigma(x) = \frac{1}{1+e^{-x}} - \frac{1}{1+e^x} =$$

$$\frac{d}{dx} \sigma(x) = \frac{d}{dx} \left(\frac{1}{1+e^{-x}} \right)$$

$$\begin{aligned}
 &\frac{d}{dx} (1+e^{-x}) = \\
 &\frac{d}{dx}(1) \cdot \frac{d}{dx}(e^{-x}) \\
 &\approx e^{-x}
 \end{aligned}$$

$$\begin{aligned}
 &= 1(e^{-x}) \cdot \frac{d(1)}{dx} - (1) \frac{d(1+e^{-x})}{dx} \\
 &\quad (1+e^{-x})^2
 \end{aligned}$$

$$\begin{aligned}
 &= - \frac{(-e^{-x})}{(1+e^{-x})^2} = \frac{e^{-x}}{(1+e^{-x})^2}
 \end{aligned}$$

$$\begin{aligned}
 &= \frac{e^{-x} + 1 - 1}{(1+e^{-x})^2} = \frac{1+e^{-x} - 1}{(1+e^{-x})^2} = \frac{e^{-x}}{(1+e^{-x})^2}
 \end{aligned}$$

Cost function: किसी misclassification
(not good)

- correct category seh

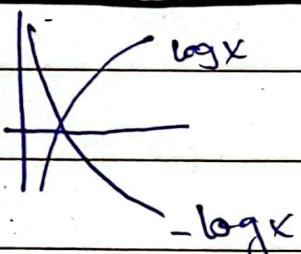
कोठा दूँदा,

humay continuous label

Banana hai.

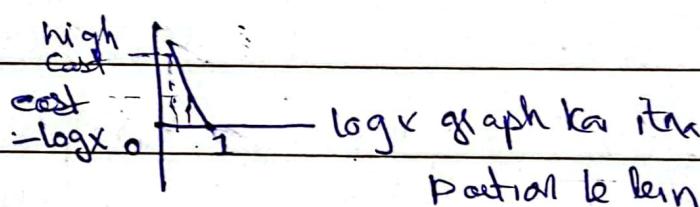
Binary class entropy

$$J(G) = -y \log(h_G(x)) - (1-y) \log(1-h_G(x))$$

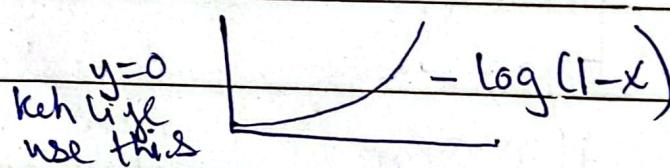
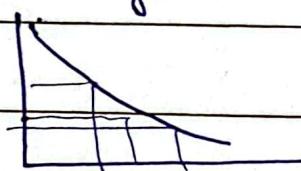


Is composite function,
done ka combine
Kai raha

ab gradient
descent
lagasaltay



near 0, cost ziyada, aavitoh
& chahye thi



$$\Theta_j = \Theta_j - \alpha \frac{\partial J(G)}{\partial \Theta_j}$$

$$J(G) = -y \log(h_G(x)) - (1-y) \log(1-h_G(x))$$

$$J'(G) = \frac{\partial (J(G))}{\partial \Theta_j}$$

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

y, \vec{x}, \vec{G}

$$J(\theta) = \frac{1}{2} \sum_i (h_{\theta}(\vec{x}_i; \theta) - y_i)^2$$

out sample
tech type

$$\theta = \theta + \alpha (y - h(\vec{x}; \theta)) \vec{x}$$

→ we can use squared error for classification,
but arr no.

$$h_{\theta}(x) = \theta^T x \quad (\text{regression})$$

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} \rightarrow \text{for logistic regression}$$

(classification)

$$J(\theta) = \frac{1}{2} \sum_i (h_{\theta}(x_i) - y_i)^2$$

↓
non-linear

b/w 0 & 1 & it is continuous

swigly little function

, can't lower it
with gradient descent

(non-convex)
problem

Can't use it

→ it is not guaranteed
to converge, or
converge into the
best possible way

decision boundary

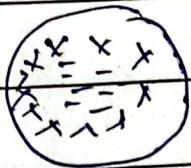
if $g(\theta^T x) \geq 0.5 \rightarrow y = 1$ else $g(\theta^T x) < 0.5 \rightarrow y = 0$

$$\theta^T x < 0$$

 x_2  x_1

$$\begin{aligned} \theta^T x = & \theta_0 + \theta_1 x_1 + \theta_2 x_2 \\ & + \theta_3 x_2^2 + \theta_4 x_1^2 \end{aligned}$$

- We can get confused in logistic regression with linear Regression
- can draw all sorts of decision boundaries
sum of sq. errors won't work here.



classification ki tareef
sani chahiye.

class (y)

1
0.5
0

regression line $h(\mathbf{x}) > 0.5 \rightarrow y = 1$

$h(\mathbf{x}) < 0.5 \rightarrow y = 0$

is it an outlier?

no

(\times) feature

luxury house or not

classifies these as economy
although it is luxury

hypothesis galat hogaya

1D

| xxx xx x |
1 0^T $\mathbf{x} > 0$

we can draw a line
but not a trend line

but a separator

2D

class 2
class 1
 $0^T \mathbf{x} < 0$

won't be a trend
line but a separator

line par agar koi value
aaye toh uski value kya
hogi?

logistic Regression

$$\mathbf{0}^T \mathbf{x} = 0$$

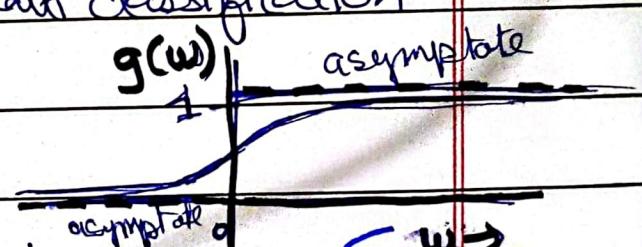
→ linear regression but Katti classification

$$h_{\mathbf{G}}(\mathbf{x}) = g(\mathbf{0}^T \mathbf{x})$$

↓
function

$$g(w)$$

$$g(w)$$



$$\text{for } \mathbf{0}^T \mathbf{x}$$

$$\text{domain } w \in (-\infty, +\infty)$$

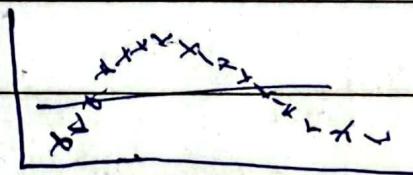
$$\text{range } 0 \leq g(w) \leq 1$$

→ can take it as
probability

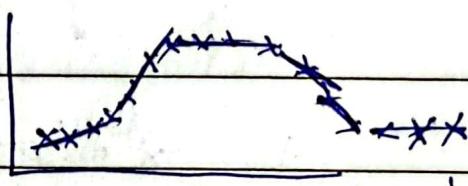
Agar trend kuch aur ho jaye toh
learning rate acha nahi



→ trend dekh kar, we can set the best learning rate, otherwise change it.



linear regression ka algo Karlega, but we know trend isn't fit.



logically administered linear regression?
google the term

large learning

linear Regression vs classification?

feature space

what is essence of classification?

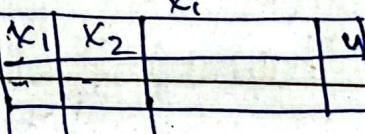
piecewise approximation

x_2

\vdots

classification mein discrete values
regression mein continuous values

x_1



x_2



values is range mein toh class

feature space

decision tree is

narrowing your search

- You can use decision trees for regression; Regression trees.

decision

trees classification

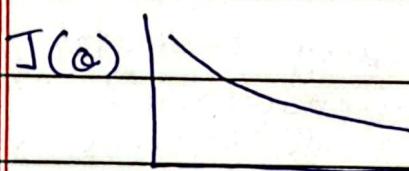
cation.

$$\alpha \in (0, 1]$$

- bohot chota alpha
bohot baat convergence
(not)
- bohot bala alpha, bali
convergence karta

hyperparameter $\rightarrow \alpha$

why is it a hyperparameter?
why not parameter? \hookrightarrow you set it.
 \hookrightarrow learned from
your algorithm.



\checkmark x aescunhi
jaana

$J(\theta)$ = sum of squared residuals

$h_{\theta}(x^T \theta) \rightarrow$ hypothesis (predicted)
 $y \rightarrow$ known value

$$(y - \theta^T x)^2$$

kitne feature vectors? n

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (y - \theta^T x_i)^2$$

x_1	x_2	...	x_d	x_0
-------	-------	-----	-------	-------

-	-	-	-	4
-	-	-	-	3

Artificial Intelligence

types? or

at linear regression solutions:

- normal eq

- gradient descent

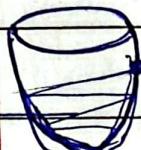
types of gd:

- sequential

- stochastic

- mini batches

$\frac{nk}{5} \cdot \frac{n}{K}$ tells this times we have to update parameter number



α = learning rate

converged in GD?

what does it mean?

it is a problem of GD

Solution to it:

$$\left(\begin{array}{l} \text{range} \\ \rightarrow 0 \end{array} \right) x_{\text{new}} = x_{\text{old}}$$

max(x)

or

$$\left(\begin{array}{l} \text{range} \\ -1 \rightarrow 1 \end{array} \right) x_{\text{new}} = \frac{x_{\text{old}} - x_{\text{mean}}}{\text{max}(x)}$$

Scaling only matters

in GD based solution,
not normal equations

$$\begin{aligned} & \min \\ & \max \\ & \text{normalization: } x_{\text{new}} = \frac{x_{\text{old}} - x_{\text{mean}}}{\text{max}(x) - \text{min}(x)} \end{aligned}$$

Learning rate kyun hona

chahiye:

→ used in GD, not
in normal eq

$$\nabla = \frac{x_{\text{old}} - x_{\text{new}}}{z_s - z_{\text{old}}}$$

$$\theta_j \leftarrow \theta_j - \alpha \nabla \theta_j$$

- resign to go downwards, opposite to gradient

↓
rate of change
step size

why do we scale our features?
→ to make convergence faster

$$\nabla_{\theta} J(\theta) = 0$$

$$(f_\theta - y)(x_\theta - y)$$

gradient symbol (nabla)

$$\theta_j = \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j} \quad \text{vector}$$

$$\theta = \theta + \alpha \sum_{i=1}^n \underbrace{\theta^T x^{(i)} - y^{(i)} \cdot x^{(i)}}_{\text{scalar}} \Delta \theta$$

$$\theta = \theta + \alpha \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)}) x^{(i)} \quad \begin{matrix} \text{batch} \\ \text{gradient} \end{matrix}$$

per una
theta

cumulative
adjustment

descent

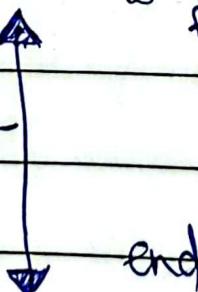
Sequential GD (online, instantaneous, ...)

while
for i=1 to n

Flavors of GD —

$$\theta = \theta + \alpha (y^{(i)} - \theta^T x^{(i)}) \cdot x^{(i)}$$

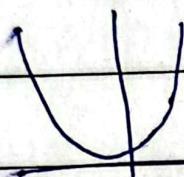
end



Stochastic GD

agar maximizing ki problem thi toh
max lagrangian

$$\theta^* = \arg \min_{\theta} J(\theta)$$



$R[a, b]$

↳ find min max
over this interval

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

dimensions:
 $(d+1 \times n) \times (n \times d+1)$
 $d+1 \times d+1 \quad (d+1) \times n$

we got this

after inv
still:
 $d+1 \times d+1$

$(d+1) \times 1$

formula
(least squared)
from 1A ch 6

$d+1 \times 1$

X matrix \rightarrow design matrix

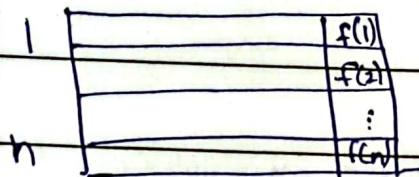
→ Problem: agar $d+1 \times d+1$ big matrix
so inverting is hard.

polynomial regression \rightarrow general term

linear regression \rightarrow specific term

zaroori whi keh sinf line guzaarein

Price



Which hypothesis is best/good for linear regression?

\rightarrow jismen cost function ki

value kam aaye

σ value... jo minimum hua (optimal)

\rightarrow sum of sq residuals criterie

batayein gaay

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

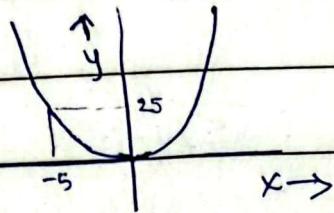
$J: \mathbb{R}^{D+1} \rightarrow \mathbb{R}$

vector \downarrow vector \rightarrow vector

scalar \downarrow scalar \rightarrow scalar

$$y^{(i)} = \begin{bmatrix} -x_1 \\ -x_2 \\ -x_3 \\ \vdots \end{bmatrix}$$

dot product is scalar



$$f(x) = x^2$$

$$x_0 = -5$$

$$\frac{d(x^2)}{dx} = 2x$$

$x^{(1)}$	x	$f(x)$
$x^{(1)} \rightarrow 1$	$x^{(1)} = -5$	25
$x^{(1)} \rightarrow 2$	$x^{(2)} = -4$	16

GD:

$$\begin{aligned}
 x_{\text{new}} &= x_{\text{old}} - \alpha (f'(x_{\text{old}})) \\
 &= -5 - 0.1 (2(-5)) \\
 &= -5 - \frac{1}{10} (-10) \\
 &= -5 + 1 \\
 &= -4
 \end{aligned}$$

gradient descend:

opposite to the
direction of gradient

gradient ascend:

$$x_{\text{new}} = x_{\text{old}} + \alpha (f'(x_{\text{old}}))$$

d+1 dimensional
input

$$\theta^T x + \theta_0 = \sum_{j=0}^d \theta_j x_j = \theta^T x$$

inner product

$$[\theta_1 \ \theta_2 \ \theta_3 \ \dots \ \theta_d \ \theta_0] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

$$x_3 = -1.25$$

$$x_4 = -1.25 - 0.25(-1.25) = 0.625$$

$$x_4 = 0.625$$

$$x_5 = 0.625 - 0.25(3.75) = -0.3125$$

x	f(x)
-5	75
2.5	18.75
-1.25	4.6875
0.625	1.171875
-0.3125	0.29296875
0.15625	0.073242

$$x_5 = 0.625 - 0.3125$$

$$x_6 = -0.3125 - 0.25 \cancel{(-0.3125)}_{6x} = 0.15625$$

6x

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

for complete vector

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

batch gradient descent

Q. $f(x) = 3x^2$ minimize using
gradient descent
iteratively

x	f(x)	GD ; $\alpha = 0.25$
$x^{(1)} = -5$	75	$v = x - \alpha f'(x)$

$$f'(x) = \frac{d}{dx} (3x^2) = 6x$$

$$\begin{aligned}
 x_1 = -5: \quad x_{n+1} &= x_n - \alpha f'(x_n) \\
 &= -5 - 0.25 \cdot 6(-5) \\
 &= -5 - 0.25(-30) \\
 &= -5 + 7.5 \\
 &= 2.5
 \end{aligned}$$

$x = 2.5:$

$$\begin{aligned}
 x_3 &= 2.5 - 0.25(6(2.5)) \\
 &= 2.5 - 3.75 = -1.25
 \end{aligned}$$

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

↓
alpha
(learning rate/const)
(step size)
[0, 1]

$$\alpha > 0 \\ [0.1, 0.2]$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = ?$$

20% consider
karna chah rakhay

$$= \left[\frac{1}{n} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2 \right] \quad 0.8 \rightarrow 80\% \text{ consider}
karna chah rakhay$$

consider we have single example

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

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

$$\frac{\partial (h_\theta(x))}{\partial \theta_j} = \frac{\partial \theta^T x}{\partial \theta_j} = \cancel{\theta^T} \cdot x_j$$

$$= \cancel{\theta^T} d^T d^{-1} d^{-1} \circ$$

$$+ (\theta_0 + \theta_1 x_1 + \dots + \theta_d x_d)$$

$$\frac{\partial (h_\theta(x))}{\partial \theta_j} = x_j$$

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

$$j^{\text{th}} \text{ component } G_j = \theta_j + \alpha (y - h_\theta(x)) \cdot x_j$$

namay salay vector mein update find karna

$$\mathbf{v}^T \mathbf{y} = \mathbf{x}^T \mathbf{y}$$

$$\hat{\mathbf{y}} = \mathbf{x}\theta$$

$$\boxed{\mathbf{x}^T \mathbf{x}\theta = \mathbf{x}^T \mathbf{y}}$$

LA Ch 6

normal equation

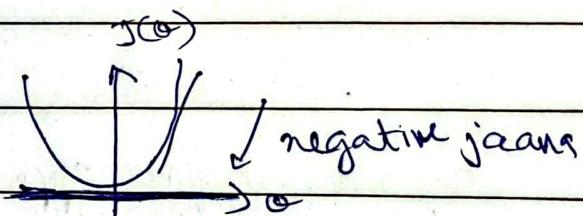
(sys of linear eq)

$$\frac{\partial J(\theta)}{\partial \theta}$$

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

Visit LA chap 6

Assume a vector of ' θ ' having specific values and then iteratively improve it to minimize the cost function for $J(\theta)$

 $J(\theta)$ is a quadratic func

$$\theta_j^{new} = \theta_j^{old} - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

gradient descent

move in opposite direction of gradient...

opposite direction to gradient

matrix x

$$X \underset{\text{rows}}{n \times (d+1)} \underset{\text{cols}}{O \underset{(d+1) \times 1}{= Y} \underset{n \times 1}{}}$$

matrix . vector =

$$\boxed{X O = Y}$$

$O = ?$ can't find O

$$A x = b \quad (\text{LA})$$

$$x = A^{-1} b$$

as $n \gg d$

X sq matrix hi
nhai

is it under determined
or over determined system?
under determined hai
since

Q over determined system
baray ga, uska
solution kya?

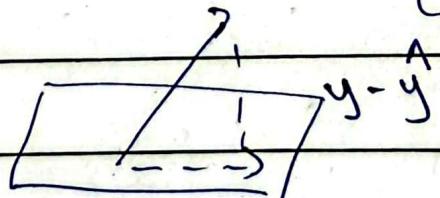
in general, no sol

$$x = A^{-1} b$$

$$\boxed{x = A^{-1} b} \rightarrow \text{col}(x)$$

so we will find approximate
solution.

(least squared approximation)



$$x^T (y - \hat{y}) = 0$$

$$x^T y - x^T \hat{y} = 0$$

$$\Theta = \begin{bmatrix} \Theta_2 \\ \Theta_3 \\ \vdots \\ \Theta_d \\ \Theta_0 \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \\ x_0 \end{bmatrix} \rightarrow \text{constant 1}$$

$$y = \sum_{j=0}^d \Theta_j x_j = \Theta^T x$$

$$h_\Theta(x) = \Theta^T x \leftarrow f(x) = x^2$$

Parameterized by Θ like

compact hypothesis representation

Cost function

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

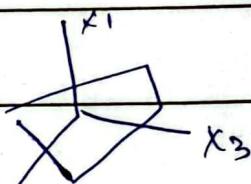
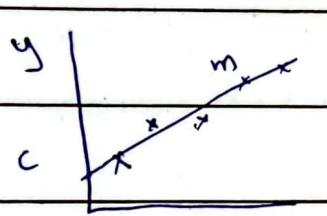
↓ ↓ ↓
predicted actual label subscript

mean sq error / sum of sq. residuals

$$\Theta^* = \underset{\Theta}{\operatorname{argmin}} J(\Theta)$$

Lec 1 $D: \mathbb{R} \rightarrow \mathbb{R}$ $D: \mathbb{R}^{d+1} \rightarrow \mathbb{R}$

$x^{(1)}$	x	y	y_1	x_2	x_3	\dots	x_d	y
$x^{(2)}$	-	-	$y^{(1)}$					$y^{(1)}$
-	-	-	$y^{(2)}$					
:	:	:						
$x^{(n)}$	-	-	$y^{(n)}$					$y^{(n)}$



$$y = mx + c$$

Residuals & sum of square errors

minimizing Kallen

Sum of sq. residuals

$$y = mx + c$$

$$y = \theta_1 x + \theta_0 \quad \text{new notational convenience}$$

agar ziyada values keh liye likhni:

$$y = \theta_d x_d + \theta_{d-1} x_{d-1} + \dots + \theta_1 x_1 + \theta_0$$

$$= \sum_{j=0}^d \theta_j x_j$$