

Lecture 1: Introduction

$$y = \alpha_0 x_0 + \alpha_1 x_1 + \dots + \alpha_n x_n + b$$

Data for initiating is the instance.

feature → only one of data set except last one

installs + runs of date set

Input + extract last column off data set

target → specific value

feature vectors now can be expressed as a single vector except table.

~~training~~ ~~prediction~~ ~~test~~ ~~data~~ + The data set from which we learn.

testing data → The data set which we test to check whether or not our training data is correct.

error made in training data is training error.

* Error made in testing data is testing error.

Hyperparameters → free variables which are independent of learning algorithm

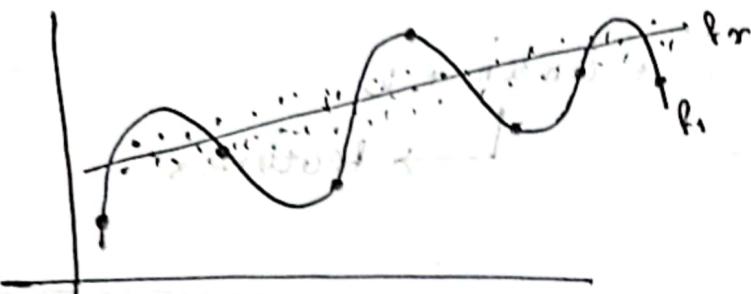
Validation Data → Testing the training data

Domains \rightarrow Each feature has a domain

Concept → It is what a learning algorithm tries to learn.

e.g. function.

$$(18) d \neq \frac{(18)^2}{19} = \frac{324}{19} \approx 17.05$$



Higher degree polynomial,
higher no. of number
of parameters,
more memory used.

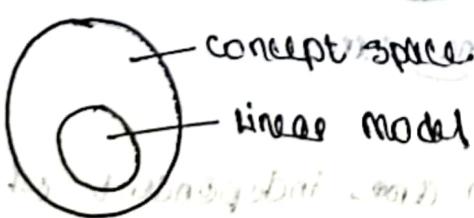
$f_t \rightarrow$ overfitting [data is less compared to capacity]
Overfitting: Giving an explanation that it is absolutely correct.

Underfitting: Unable to give an explanation.

[data is more, but capacity is less]

at least two different models fit the data with equal accuracy

Hypothesis: A guess / Assumption about the concept.



Inductive bias: Bias produced by hypothesis.

Empirical error: Error that occurs when our prediction does not match with the training data.

$$\text{Res} = \frac{1}{m} \sum_{i=1}^m |f(x_i) - h(x_i)|$$

Generalization Error: Error of our learning algorithm if we could test it against all the data instances in the domain.

Before $\frac{1}{m} \sum_{i=1}^m \text{loss}(y_i, \hat{y}_i)$

After $\frac{1}{m} \sum_{i=1}^m \text{loss}(y_i, \hat{y}_i) + \lambda \|\theta\|^2$

$$\frac{1}{m} \sum_{i=1}^m \text{loss}(y_i, \hat{y}_i) + \lambda \|\theta\|^2$$

\downarrow sum of weighted domain of data points

\downarrow weight

$P(x_1, x_2, \dots, x_n) = \prod P(x_i)$ $P(x_1, x_2, \dots, x_n | x_1, x_2)$

$$P(x_1) P(x_2) P(x_3) \dots P(x_n | x_1, x_2)$$

we cannot find real-time probability. Thus, cannot find generalization error of model

to solve using some method

more elegant

Classification: Assigning a label to each data instance.

classification with missing inputs: Some of the information will be absent. We can yet find out the data

using other available informations: ~~missing to missing~~

-~~one to one~~ ~~one to many~~

Regression: We don't assign label to data instance like we did in classification.

estimate a value which is real number

Transcription: Converts unstructured data into textual forms

-~~one to one~~ ~~one to many~~ speech to text, ~~join w sub~~

Ranking: Given a search prompt, a search engine generates a list of webpages which need to be sorted for user's convenience.

Clustering: Identifying groups of similar entities e.g. community detection in social media.

Manifold learning: We plot something which has a curved surface which need to be learned. Task of reducing "feature" space intelligently so that learning can occur more efficiently.

Anomaly detection: For any suspicious data points which is weird, anomaly detection is needed.

एसो इन वाले तुम्हें यह क्या करते होंगे?

Synthesis and Sampling: A new work is created from existing past works. (synthesis)

for example: AI is fed with Rabin德拉নath Tagore's songs and asked to write a song. It writes a song similar to that of Tagore's songs.

Imputation of missing values: Sometimes, there are missing values in a data set we can use our wisdom gained through learning to fill up these missing values.

Denoising: If a data x is observed to be corrupted \bar{x} due to noisy process, from experience, given \bar{x} , we can predict the original clean x .

Probability mass/density estimation: we are just learning a function $p(x)$ for given instance x . which is a measure of probability $p(x)$.

Supervised learning: Classification, ranking, regression are ~~with~~ ^{with} ~~different~~ ^{different} ~~task~~ ^{task}

• silent and supervised tasks - ~~with~~ ^{with} ~~input~~ ^{input} ~~output~~ ^{output} ~~task~~ ^{task}

Unsupervised learning: No column for labels present.

Clustering and dimensionality reduction

are some of the unsupervised learning

techniques. ~~with~~ ^{with} ~~input~~ ^{input} ~~output~~ ^{output} ~~task~~ ^{task}

~~with~~ ^{with} ~~input~~ ^{input} ~~output~~ ^{output} ~~task~~ ^{task} other features ~~task~~ ^{task}

Semi supervised learning: Is ideal when labeled data is expensive.
we can learn from given labels, utilise similarity measures as in unsupervised learning to make predictions.

Transductive inference while the semi-supervised learning

we get some labelled & some

unlabelled data.

- In case of induction, we learn from ~~with~~ ^{with} training set to make general predictions.
- makes no assumption about test data.
- Test data are unlabelled instances.
- Goal: minimize error in test samples only.

* Supervised learning, unsupervised learning and semi-supervised learning are inductive inference. We know the data set for learning, but the data set for test is unknown.

Online learning: ~~the algorithm continues to learn without being exposed to the environment~~

Offline learning: All the learning happens before testing.

It learns from the training training set

... learning before it is tested. ~~and it is exposed to a frequent~~

~~test set with the same inputs~~

~~which is different from the test set~~

Online learning: training and testing phases are checked.

~~Agent acts learned, prediction is made when a data~~

~~instance is encountered.~~

~~error if agent~~ learns from mistake.

~~agent starts with empty memory~~

~~acquisition of all necessary information~~

~~initialization of actions~~

Reinforcement learning: Environment doesn't give any kind of experience ~~as long as~~ supervision.

~~Agent gives reward/punishment for~~

~~an action.~~

~~Agent's goal: maximise rewards in future.~~

Active learning: Carefully choose data points for

~~you choose them, you can do better at~~

~~learning.~~

~~going to reach goal: get results comparable to the~~

~~supervised learning.~~

~~academic boundaries and initial goal~~

~~from original test in some remaining part~~

~~more than current performance, which requires the~~

~~most difficult and difficult situations the trained behaviour~~

~~to fit the latest with that, which not yet been used in the~~

~~environment~~

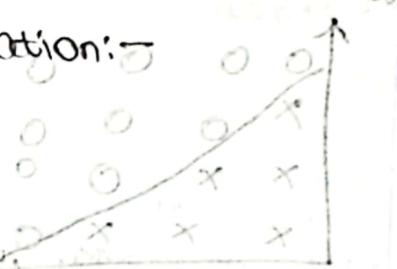
Lecture 7

I will!

Empirical risk minimization: Error needs to be minimized.

These can be risk or empirical risk minimization:-

- (i) Overfit
 - (ii) Underfit



* Hypothesis can introduce bias, named inductive bias.

* Inductive bias in RNN is ~~can~~ remedied to overfitting.

We have to reduce the search space. From the search space, we have to find the least square.

If we increase data, then errors will decrease.

Pro (RNGE) $\leq \frac{4c}{\epsilon}$ → no. of dose needed for accuracy
~~me~~ (or 0.01)
 ↓ hypothesis
 ↓ hypothesis
 ↓ error $[e_D, M] \approx 0$ $e_D + M = [e_D, M]$
 $[e_R, M] \approx 0$ constant with M
+ avengement - non

Learning Bound for Finite Inconsistent Hypothesis Classes

$$R_{\text{EN}} \leq R_{\text{SLN}} + \sqrt{\frac{\log(4H) + \log^2 n}{2m}}$$

target elimination error control

be found outside

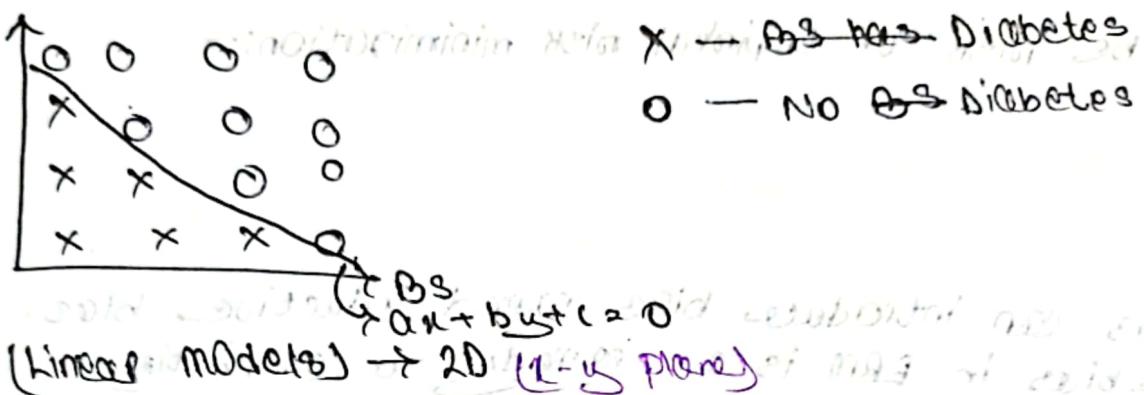
* Upper limit of generalization can be determined.

1) $H \downarrow$, $R_{HS} \downarrow$ \rightarrow Reactions \downarrow
 \rightarrow Size \downarrow \rightarrow $d \rightarrow$ r_{M-L}
 अली \downarrow \rightarrow r_{M-L} \downarrow
 \downarrow \rightarrow $\frac{1}{r_{M-L}}$ \downarrow \rightarrow $k \downarrow$

Lecture 6

Linear functions and Affine functions

w



$$\begin{matrix} \times & \times & | & 0 & 0 & 0 \\ & & & \text{BS} & & \end{matrix}$$

BS

1D (Line plane)

Hyperplane: Line, Plane. A plane whose eqn. is:-

$$w_1x_1 + w_2x_2 + \dots + w_nx_n + b = 0$$

$$f(x_1, x_2) = w_1x_1 + w_2x_2$$

(Affine function)

$w = [w_1, w_2]$

$x = [x_1, x_2]$

\hookrightarrow non-homogeneous

*Affine function \rightarrow linear function + something extra
(homogeneous)

$$f(x_1, x_2) = w_1x_1 + w_2x_2 + b$$

$w = [w_1, w_2, \dots, w_n]$

$$w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$

$$w_1, w_2, \dots, w_n + b$$

\downarrow

w weight

\downarrow

bias term

x_1, x_2

non-homogeneous condition

$$A = [k_{11} k_{12} \dots \dots k_{n1}]$$

(b) 

$$W_1 k_1 + W_2 k_2 + \dots + W_n k_n + W_{n+1} k_{n+1}$$

$$W_1 k_1 + W_2 k_2 + \dots + W_n k_n + W_{n+1} k_{n+1} + \text{excess weight}$$

$$\begin{cases} k_{n+1}^2 = 1 \\ W_{n+1} = b \end{cases}$$

2. LAGRANGE

derivative = 0



$$\left\{ d/dx (W_1 k_1 + \dots + W_n k_n) \right\} = 0$$

$$\left\{ d/dx (W_1 k_1 + \dots + W_n k_n + W_{n+1} k_{n+1}) \right\} = 0$$

method of undetermined coefficients

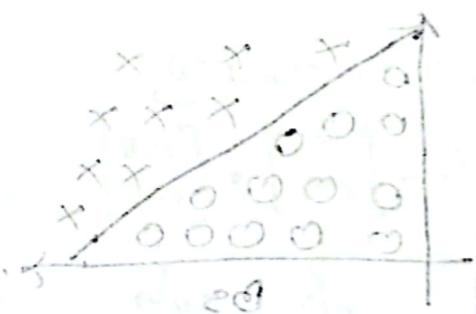
combined method of
undetermined coefficients

method of successive approximations

for unit A - step

out extrapolate with

step



[Extrapolate] $X \leftarrow O \leq 0.01$

[Extrapolate] $O \leftarrow O >$

initial value $\rightarrow 0.01$, 0.0015 half of difference
between X and O $\rightarrow 0.005$ \rightarrow $\frac{1}{2}(0.01 + 0.005)$

$O \leq 0.01$ \rightarrow step 0.001
 $O >$

extrapolate $\rightarrow 0.0005$ \rightarrow $\frac{1}{2}(0.001 + 0.0005)$

extrapolation and extrapolation (Turing) rule

[Extrapolate] $X \leftarrow 0.01$

[Extrapolate] $O \leftarrow 0.01$

$w_1x_1 + w_2x_2 + \dots + w_nx_n + b \rightarrow \text{linear model}$

$\Rightarrow \{w, x\}$

$$w = \{w_1, w_2, \dots, w_n, b\}$$

$$x = \{x_1, x_2, \dots, x_n, 1\}$$

non-homogeneous system

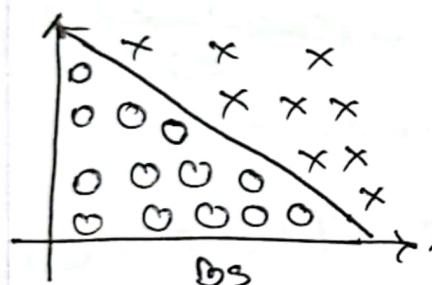
non-homogeneous coordinate system

univariate linear regression

$$y = ax + b$$

↓
one variable

w



↑ linearly separable
data — A line which separates two data



$(w \cdot x) > 0 \rightarrow x \text{ [Diabetes]}$

$< 0 \rightarrow o \text{ [No Diabetes]}$

$\{w_1, w_2, b\}$ To find these, we can
 $\{w_1, w_2, b\}$ use perceptron algorithm
or gradient descent

$w_1x_1 + w_2x_2 + b > 0$

< 0

So,

sign($w \cdot x$) → Hyperplanes and Hypotheses

+ve → x [Diabetes]

-ve → o [No Diabetes]

प्र

$$\sum_{i=1}^m$$

$$\text{sign}(w_1x_1 + w_2x_2 + b) \neq c_i$$

जहां तक है

$x_i = (x_{i1}, x_{i2}, \dots, x_{in})^T$,
 $w = (w_1, w_2, \dots, w_n)^T$

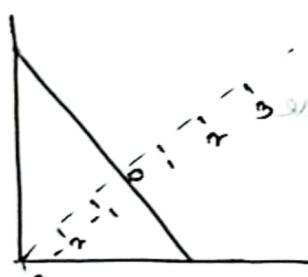
$$c_i = \begin{cases} 1 & \text{if } \dots \\ -1 & \text{if } \dots \end{cases}$$

if c_i matches, error = 1

if c_i doesn't match, error = 0

* We cannot use differentiation here.

gradient descent :-



$$\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} - \frac{\partial E}{\partial w_1} \quad \text{or} \quad \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} - \frac{\partial E}{\partial w_2}$$

* line अब ज्ञान 0 तक आए $(x, w) - P \leq \frac{1}{m}$

* line अब more -ve & line

अब more +ve

* error = 0 or 1 के लिए $\frac{1}{m+1} \times (error / 1 \cdot P)$ बेस्ट बेस्ट है

gradient = 1 \times (mean of error)

$$\text{sign}(w_1x_1 + w_2x_2 + b) = +1$$

$$c_i = \begin{cases} +1 \\ -1 \end{cases}$$

$$\sum_{i=1}^m c_i (w_1x_1 + w_2x_2 + b) \rightarrow$$

$\text{sign} = +1, c_i = +1, \text{sign} * c_i = +1$ $\text{sign} = -1, c_i = -1, \text{sign} * c_i = +1$ $\therefore \text{sign} * c_i (x, w) - P \geq 0$
--

Sigmoid function / Logistic function

~~W₁, W₂, b~~

$$L(w_1, w_2) = w_1 x_1 + w_2 x_2 + b$$

$$\frac{1}{1 + e^{-w_1 x_1 - w_2 x_2 - b}}$$

↳ (Diabetes A2
probability test 0.210 M%)

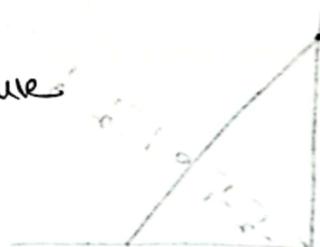
$$\log \left(\frac{1}{1 + e^{-w_1 x_1 - w_2 x_2 - b}} \right)$$

$$2(\ln(y_i)) = \log \left(\frac{1}{1 + e^{-w_1 x_1 - w_2 x_2 - b}} \right) \rightarrow y_i, h_w(x_i) \in \{-1, 1\}$$

$$L(w_1, w_2, b) = -\log -y_i \log h_w(x_i) - (1 - y_i) \log (1 - h_w(x_i))$$

$$\sum_{i=1}^m \frac{-y_i (w_1 x_1 + w_2 x_2 + b)}{1 + e^{-w_1 x_1 - w_2 x_2 - b}} \rightarrow \text{Perceptron Learning Rule}$$

$$\frac{1}{m} \sum_{i=1}^m \frac{-y_i (w_1 x_1 + w_2 x_2 + b)}{1 + e^{-w_1 x_1 - w_2 x_2 - b}}$$



$$\frac{1}{m} \sum (y - (w \cdot x))^2$$

$$P(Y=1 | w, x) = \frac{1}{1 + e^{-w_1 x_1 - w_2 x_2 - b}}$$

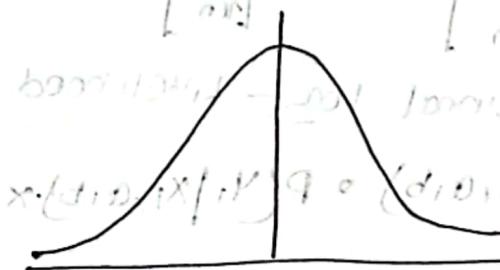
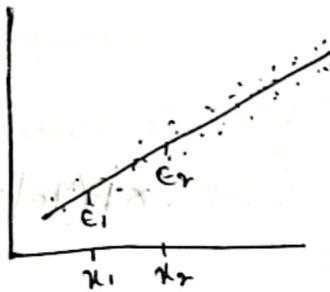
$$P(Y=0 | w, x) = 1 - h_w(x)$$

$$\frac{1}{1 + e^{-w_1 x_1 - w_2 x_2 - b}} = \frac{1}{1 + e^{-(d + \text{offset} + \alpha w \cdot x)}} = \frac{1}{1 + e^{-(d + \text{offset} + \alpha w_1 x_1 + \alpha w_2 x_2)}} = \frac{1}{1 + e^{-(d + \text{offset} + \alpha w_1 x_1 + \alpha w_2 x_2 + \alpha b)}}$$

$$\frac{1}{1 + e^{-(d + \text{offset} + \alpha w \cdot x)}} = \frac{1}{1 + e^{-(d + \text{offset} + \alpha w_1 x_1 + \alpha w_2 x_2 + \alpha b)}}$$

गणितीय विधि का अध्ययन करते हैं।
 $(\bar{x}) \in (\bar{A}) \subset (\bar{M}(\bar{A}))$

$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ Gaussian Distribution



Average error \bar{e}_2 का मूल्य ना बताया।
 मान ना बताया।

उसके बारे में जानकारी नहीं दी गई।

Error

* \bar{e}_2 और \bar{e}_1 का अपेक्षित मूल्य नहीं दी गई।

* \bar{e}_2 का अपेक्षित मूल्य नहीं दी गई।

$$\bar{e}_2 = \frac{\sum_{i=1}^m (x_i - \bar{x})}{m}$$

$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$

$e_i = x_i - \bar{x}$

$$\text{variance } \sigma^2 = \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n}$$

$$\text{standard deviation } \sigma = \sqrt{\frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n}}$$

$\bar{e}_2 = \frac{(e_1 + e_2 + \dots + e_m)}{m}$ average और अपेक्षित मूल्य

$$\left(\frac{(e_1 + e_2 + \dots + e_m)}{m} \right) = \frac{1}{m} \sum_{i=1}^m e_i = \frac{1}{m} \sum_{i=1}^m (x_i - \bar{x})$$

$$\frac{1}{m} \sum_{i=1}^m (x_i - \bar{x})$$

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}$$

* SIB Doctor independent \bar{x}_i

$$* P(A \cap B) = P(A)P(B)$$

Conditional Log - Likelihood

$$P(Y|X, \alpha, b) = P(Y_1|X_1, \alpha, b) \times P(Y_2|X_2, \alpha, b) \times \dots \times P(Y_n|X_n, \alpha, b)$$

$$\propto P(Y_1, Y_2, \dots, X_i = x_i | \alpha, b)$$

$$= P(E_1, E_2, \dots, E_n | X_1 = x_1, \dots, X_n = x_n, \alpha, b)$$

$$= 2P(E_1 | X_1 = x_1, \alpha, b)P(E_2 | X_2 = x_2, \alpha, b) \dots P(E_n | X_n = x_n, \alpha, b)$$

$$= 2C e^{-\frac{1}{2}\frac{\epsilon_1^2}{\sigma^2}} e^{-\frac{1}{2}\frac{\epsilon_2^2}{\sigma^2}} \dots e^{-\frac{1}{2}\frac{\epsilon_n^2}{\sigma^2}}$$

$$= 2C e^{-\frac{1}{2}(\frac{\epsilon_1^2}{\sigma^2} + \frac{\epsilon_2^2}{\sigma^2} + \frac{\epsilon_3^2}{\sigma^2} + \dots + \frac{\epsilon_n^2}{\sigma^2})}$$

maximise \leftarrow

\Rightarrow $\partial C / \partial \sigma^2$

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

\downarrow

ϵ_1

$$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\frac{(x_i-\mu)^2}{\sigma^2}}$$

$$\mu = \frac{\sum_{i=1}^m x_i}{m}$$

$$\sigma^2 = \sqrt{\frac{\sum_{i=1}^m (x_i - \mu)^2}{m}}$$

$$+ \frac{(x_i - \mu)}{m} + (x_i - \mu) = \frac{1}{m} \sum_{i=1}^m (x_i - \mu)$$

$$\ln J_2 = \ln C - \frac{1}{2\sigma^2} \left(\frac{\epsilon_1^2}{\sigma^2} + \frac{\epsilon_2^2}{\sigma^2} + \frac{\epsilon_3^2}{\sigma^2} + \dots + \frac{\epsilon_n^2}{\sigma^2} \right)$$

$$\ln J_2 = \ln C - \frac{1}{m\sigma^2} (\epsilon_1^2 + \epsilon_2^2 + \dots + \epsilon_n^2)$$

$$-\ln J_2 = -\ln C + \frac{1}{m\sigma^2} (\epsilon_1^2 + \epsilon_2^2 + \dots + \epsilon_n^2)$$

$$f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \rightarrow \text{Normal distribution}$$

14-41

* Error \Rightarrow यहाँ जाने वाली Gaussian distribution के निम्नलिखित गुणों का उपयोग करें।

BINARY CLASSIFICATION

$$\alpha = \{0, 1\}$$

CLASSIFICATION DISTRIBUTION

$$P(Y=0|w, b) = \text{hw}(w) \quad P(Y=1|w, b) = \frac{1}{1 + \text{hw}(w)}$$

if $w=0$,

$$P(Y=0|w, b) = 1 - \text{hw}(w)$$

$$P(Y=1|w, b) = \text{hw}(w)$$

P_0	P_1	P_2	P_3	P_4
0	0	0	1	0
0	0	0	0	1

$$\frac{4/10}{5/10} = \frac{8}{10}$$

$$P(Y|X, w) = P(Y_1|X_1, w) \cdot P(Y_2|X_2, w) \cdot \dots \cdot P(Y_m|X_m, w)$$

↑ patient's weight

$$= P(Y_1=0|X_1, w) \cdot P(Y_1=1|X_1, w) = \frac{1}{2}$$

$$J = \text{hw}(w)(1 - \text{hw}(w))^{1-w} \text{hw}(w)(1 - \text{hw}(w))^{1-w}$$

$$J = \text{hw}(w)(1 - \text{hw}(w))^{1-w} (1 - \text{hw}(w))^{1-w}$$

$$\log J = \sum_{i=1}^m \log [\text{hw}(w_i)(1 - \text{hw}(w_i))^{1-w_i}]$$

$$\text{Loss} = \sum_{i=1}^m \log \text{hw}(w_i)^{a_i} + \log (1 - \text{hw}(w_i))^{1-a_i} \rightarrow -\log \text{hw}(w_i)$$

$$\text{Loss} = -\sum_{i=1}^m a_i \log \text{hw}(w_i) \rightarrow (1-a_i) \log (1 - \text{hw}(w_i)) \rightarrow -H(a_i \log (1 - \text{hw}(w_i)))$$

error \Rightarrow minimize loss

regression | classification → total entropy error
 ↓
 sum of squared binary $\log_{2} m + 1$

$$\sum_{i=1}^m -\alpha_i \log(h_w(x_i)) - (1-\alpha_i) \log(1-h_w(x_i)) \rightarrow \text{binary classification}$$

କ୍ଷେତ୍ର
value
at A₂
value at
station 2
लାଗୁଣ୍ୟ

- binary classification

114 - 1

Value ai A3 அனாலி திருமதி கீழ்த்தான் பி சென்ட்
value A2 போன்ற விவரங்களை விடக்கூடிய
வாட்டு மதி
நாப்பி முடிவுகளை விடக்கூடிய

MULTINOMIAL CLASSIFICATION

$\hat{Y}_i = \beta_0 + \beta_1 X_i$ \rightarrow Linear Regression

x_1	R	check for color
x_2	G	

B_1	B_2	B_3	C_1	C_2
0	1	0	0	0
1	0	0	0	0

$w_{11} \rightarrow$ matrix
 \downarrow color 1 (C1)
 \downarrow feature 1 of C1

$$\rightarrow \text{score for } y$$

$$S_{1,2} w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

$$S_{K+} = W_{1K}X_1 + W_{2K}X_2 + \dots$$

$$S_{K-} = W_{1K}X_1 + W_{2K}X_2 + \dots$$

$$S_K = W_{1K}X_1 + W_{2K}X_2 + \dots$$

$$S_K = \frac{w_{1K}x_1 + w_{2K}x_2 + \dots}{w_{KK}x_K + b_K}$$

$$x_{ij} \rightarrow \text{color}(423|x_{ij})$$

1988.8.18.19

elbow ↴

$$= 8810\% \quad (43.31\%)$$

$$-\log(4\pi\beta/x_i)$$

9

P-2 Using Softmax function:-

$S_0^2 S_1, S_2, \dots$ Fin-SK $\xrightarrow{\text{in}}$ **Software Functions**
 $-1 -3$ $\text{Finest} \rightarrow \text{Coarsest}$ \Rightarrow **coarse** **soft** **area**

$$\frac{e^2}{m} + \frac{e^2}{m_0} + \frac{e^2}{m_0} = 1$$

$B_2 - \bar{e}^1 + \bar{e}^3$ first of \bar{e}^1

$$\sum -\alpha_i \log(P(x_i)) + (1-\alpha_i) \log(1-P(x_i)) = P(y_2 \text{ out})$$

TUESDAY

DATE: 20/09/23

Gradient Descent

As value directly $w = (x^T x)^{-1} x^T y$ for 23/09/23,
we are using offline learning.

$$J = (y_1 - \bar{y})^2 + (y_2 - \bar{y}_2)^2 + \dots + (y_m - \bar{y}_m)^2$$

$$J = \sum_{i=1}^m (y_i - \bar{y}_i)^2$$

$$\frac{\partial J}{\partial w} = \sum_{i=1}^m 2(y_i - \bar{y}_i) \frac{\partial}{\partial w_i} (y_i - \bar{y}_i)$$

$$= \sum_{i=1}^m 2(y_i - \bar{y}_i)(-x_{i1})$$

$$\bar{y}_i = \alpha_1 w_1 + \alpha_2 w_2 + \dots + \alpha_n w_n + b$$

* Gradient of J is the direction at which the value of function increases the fastest.

$$\frac{\partial J}{\partial w_1} = \sum_{i=1}^m 2(y_i - \bar{y}_i) \frac{\partial}{\partial w_1} (y_i - \bar{y}_i)$$

$$= \sum_{i=1}^m 2(y_i - \bar{y}_i)(-x_{i1})$$

* Gradient $\frac{\partial J}{\partial w_1}$ of J function As value $\frac{\partial J}{\partial w_1}$

* Gradient $\frac{\partial J}{\partial w_2}$ opposite $\frac{\partial J}{\partial w_1}$ function As value $\frac{\partial J}{\partial w_2}$

$$\frac{\partial J}{\partial w_2} = \sum_{i=1}^m 2(y_i - \bar{y}_i) \frac{\partial}{\partial w_2} (y_i - \bar{y}_i)$$

$$= \sum_{i=1}^m 2(y_i - \bar{y}_i)(-1)$$

$$\begin{bmatrix} \alpha'_1 \\ \alpha'_2 \\ \alpha'_3 \\ \vdots \\ b' \end{bmatrix} = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \\ \vdots \\ b \end{bmatrix} - \boxed{\lambda} \nabla J \quad \begin{array}{l} \text{learning vector} \\ \text{rate} \end{array}$$

$$\left. \begin{array}{l} \alpha'_1 = \alpha_1 + 2 \sum_{i=1}^m (y_i - \bar{y}_i)x_{i1} \\ \alpha'_2 = \alpha_2 + 2 \sum_{i=1}^m (y_i - \bar{y}_i)x_{i2} \\ \vdots \\ \alpha'_n = \alpha_n + 2 \sum_{i=1}^m (y_i - \bar{y}_i)x_{in} \\ b' = b + 2 \sum_{i=1}^m (y_i - \bar{y}_i) \end{array} \right\} \quad \begin{array}{l} \text{learning with} \\ \text{gradient descent for linear regression} \end{array}$$

Learning with gradient descent [for linear classification]

$$\bar{y}_i = \alpha_0 x_{i0} + \alpha_1 x_{i1} + \alpha_2 x_{i2} + \dots + \alpha_n x_{in} + b$$

\downarrow actual value

$x_i \rightarrow$ given value

$$\begin{cases} \bar{y}_i > 0 \rightarrow \text{correct} \\ \bar{y}_i \leq 0 \rightarrow \text{incorrect} \end{cases}$$

\downarrow minimise

* $\bar{y}_i & y_i$ (same sign) error = 0

* $\bar{y}_i & y_i$ (different sign) error $\neq 0$

$$-1 * -1 = 1$$

+ given value

10 $\bar{y}_i & y_i$ have same sign

value $\neq 0$

$$\min(-1, 0) = 0$$

$$\max(-\bar{y}_i, 0)$$

$$J = \sum_{i=1}^m \min(-\bar{y}_i, 0)$$

$$\frac{\partial J}{\partial \alpha_j} = \sum_{i=1}^m \left\{ \begin{array}{l} 0 \\ -\bar{y}_i x_{ij} \end{array} \right\}$$

$\bar{y}_i x_{ij} \neq 0$

$$\frac{\partial J}{\partial \alpha_j} = \sum_{i=1}^m \left\{ \begin{array}{l} 0 \\ -\bar{y}_i x_{ij} \end{array} \right\}$$

$$\alpha'_1 = \alpha_1 - \lambda \frac{\partial J}{\partial \alpha_1}$$

$$\alpha'_2 = \alpha_2 - \lambda \frac{\partial J}{\partial \alpha_2}$$

$$\text{grad } J = \frac{\partial J}{\partial \theta}$$

$$\text{grad } J = \frac{\partial J}{\partial \theta}$$

$\bar{y}_i \rightarrow$ given data

if \bar{y}_i & y_i have same sign
then the error has to be 0
But $-(\bar{y}_i)$ gives a negative
value. Thus, we use
 $\max(-\bar{y}_i, 0)$ to get the
correct err.

- minimize J \rightarrow sub. gradient descent!

(0, ..., 0, 0) \rightarrow initial

step proportion \rightarrow rate of 0.1

now ($0.1 \times \text{grad } J$) \rightarrow $\text{new } \alpha$
following step \rightarrow $\text{new } \alpha$
 $\alpha = \alpha + \text{new } \alpha$ \rightarrow $J = \text{new } J$

* $(\bar{y}_i) \neq 0$, then $\frac{\partial J}{\partial \alpha_j}$ $\neq 0$

$\bar{y}_i \neq 0, \bar{y}_i \neq 0$

* $\bar{y}_i \neq 0$, then $\frac{\partial J}{\partial \alpha_j}$ $\neq 0$
value $\bar{y}_i \neq 0, \bar{y}_i \neq 0$

$$J = \sum_{i=1}^m \min(0, -\bar{y}_i) \rightarrow \text{offline learning}$$

$$J = \min(0, -\bar{y}_i) \rightarrow \text{online learning}$$

YiXi_i < 0

$$\frac{\partial J}{\partial w_i} = -YiXi_i$$

$$\frac{\partial J}{\partial w_0} = -YiXi_0$$

∴ $w_0 = \text{constant}$

$$\frac{\partial J}{\partial w_i} = -YiXi_i$$

$$\frac{\partial J}{\partial b} = -Yi^+ \text{ sum } \text{some } \text{add } \text{some}$$

\rightarrow $b = \text{sum } \text{some } \text{add } \text{some}$

\rightarrow $b = \text{sum } \text{some } \text{add } \text{some}$

\rightarrow $b = \text{sum } \text{some } \text{add } \text{some}$

Perceptron algorithm:-

input: A training set (x_i, y_i) ...

initialize: $w^{(0)} = (0, \dots, 0)$

for $t=1, 2, \dots$

if $(\exists i, s.t. \underbrace{w^{(t-1)} \cdot x_i}_{\text{dot product of } w^{(t-1)} \text{ & } x_i} + b^{(t-1)} \leq 0)$ then

$w^{(t)} = w^{(t-1)} + y_i x_i$

else:

output $w^{(t)}$,

if $y_i \neq 0$ then $0.5 \cdot \text{sum } \text{some}$

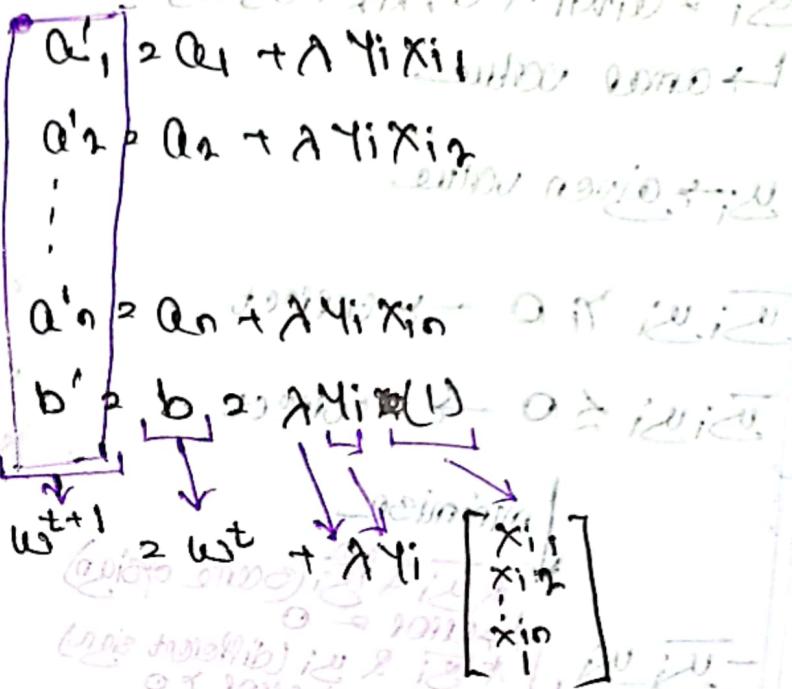
else $0.5 \cdot \text{sum } \text{some}$

\rightarrow $w^{(t+1)} = w^{(t)} + 0.5 \cdot \text{sum } \text{some}$

\rightarrow $w^{(t+1)} = w^{(t)} + 0.5 \cdot \text{sum } \text{some}$

\rightarrow $w^{(t+1)} = w^{(t)} + 0.5 \cdot \text{sum } \text{some}$

\rightarrow $w^{(t+1)} = w^{(t)} + 0.5 \cdot \text{sum } \text{some}$



*For perceptron algorithm,
 λ is always 1

After iteration

$0 \cdot (0, F_1 - 1, 0, \dots)$

$0, F_1 - 1, 0, \dots$

[incorrect F_1 value, update F_1]

$0.5 \cdot F_1$

$0.5 \cdot F_1 + k_{11} \rightarrow 1 \text{ number house}$

1 number feature

$w \rightarrow \text{weight}$

$n \rightarrow \text{feature}$

$\frac{F_1}{n} \cdot 1 \rightarrow 10 \cdot 10$

CLASSIFIERS

→ 1 Naive Bayes Classifier

0 Decision Trees

1 Boosting

$$P(A \cap B) = P(A)P(B|A)$$

NAIVE BAYESIAN CLASSIFIER

$$P(B)P(A|B) = P(A \cap B) = P(A)P(B|A)$$

$$P(A|B) = \frac{P(A)P(B|A)}{P(B)} \rightarrow \text{Baye's Theorem}$$

$$\begin{array}{c} x_1, x_2, x_3 \\ \text{---} \\ 0, 1, 2, 0 \end{array} \quad \left| \begin{array}{c} u \\ \text{---} \\ \text{Yes/No} \end{array} \right. \quad B = \{0, 1, 2, 0\}$$

$$P(u = \text{Yes} | x_1 = 0, x_2 = 1, x_3 = 0) = \frac{P(u = \text{Yes}) + P(x_1 = 0, x_2 = 1, x_3 = 0) | u = \text{Yes}}{P(x_1 = 0, x_2 = 1, x_3 = 0)}$$

* Naive Bayesian classifier is used for discrete, categorical data.
 * Discrete cannot be used directly, even though they are discrete. We need to group.

10 - 20 → child

Leg: d9, T=0, d=1, 0=T, E=0

Out: d9, T=0, d=1, 0=T, E=0

* Get distinct column or date or group values as a group
 2013 तक, जो values same be 21 unique अन्य ग्राहक,

$$P(Ph = Yes | O_2 S, T = C, H = h, W = T) = \frac{P(Ph = Yes)P(O_2 S, T = C, H = h, W = T | Ph = Yes)}{P(O_2 S, T = C, H = h, W = T)}$$

$$P(Yes) = \frac{9}{9+3}$$

$$P(O_2 S, T = C, H = h, W = T | Ph | Ph = Yes)$$

$$= P(T = S | Ph = Yes) + P(T = C | Ph = Yes)$$

$$+ P(H = h | Ph = Yes) + P(W = T | Ph = Yes)$$

$$\Rightarrow \left(\frac{3}{9}\right) * \left(\frac{3}{9}\right) * \left(\frac{3}{9}\right) * \left(\frac{3}{9}\right)$$

$P(A \cap B) = P(A)P(B) \rightarrow$ Absolutely independent

$P(A \cap B | C) = P(A | C) * P(B | C) \rightarrow$ Conditionally independent

$$P(O_2 S, T = C, H = h, W = T | Ph = Yes) = P(O_2 S, T = C, H = h, W = T, Ph = Yes)$$

$$+ P(O_2 S, T = C, H = h, W = T, Ph = No)$$

$$P(O_2 S, T = C, H = h, W = T | Ph = No) = P(T = S | Ph = No)$$

$$P(O_2 S, T = C, H = h, W = T, Ph = Yes) = P(Ph = Yes)P(O_2 S, T = C, H = h, W = T | Ph = Yes)$$

$$P(Ph = Yes | O_2 S, T = C, H = h, W = T) = \frac{P(O_2 S, T = C, H = h, W = T, Ph = Yes)}{P(O_2 S, T = C, H = h, W = T, Ph = Yes) + P(O_2 S, T = C, H = h, W = T, Ph = No)}$$

$$P(\text{Ph} = \text{Yes} | O_2 S, T_2 t, H_2 h, W_2 T) = \frac{P(O_2 S, T_2 t, H_2 h, W_2 T | \text{Ph} = \text{Yes}) / P(\text{Ph} = \text{Yes})}{P(\text{Ph} = \text{Yes}) P(O_2 S, T_2 t, H_2 h, W_2 T | \text{Ph} = \text{Yes}) + P(\text{Ph} = \text{No}) P(O_2 S, T_2 t, H_2 h, W_2 T | \text{Ph} = \text{No})}$$

$$P(\text{Ph} = \text{No} | O_2 S, T_2 t, H_2 h, W_2 T) = \frac{P(O_2 S, T_2 t, H_2 h, W_2 T | \text{Ph} = \text{No}) / P(\text{Ph} = \text{No})}{P(\text{Ph} = \text{Yes}) P(O_2 S, T_2 t, H_2 h, W_2 T | \text{Ph} = \text{Yes}) + P(\text{Ph} = \text{No}) P(O_2 S, T_2 t, H_2 h, W_2 T | \text{Ph} = \text{No})}$$

Frequent list approach: If probability of event A = 1 and probability of event B = 0, then frequent list approach occurs.

* Using more data, we can get probability of B.

LAPLACE SMOOTHING

$\text{Ph} = \text{Yes}$	R	O	S
9	$\frac{3}{9}$	$\frac{6}{9}$	$\frac{0}{9}$
12	$\frac{1}{12}$	$\frac{7}{12}$	$\frac{1}{12}$

↓

$\text{Ph} = \text{Phes}$	R	O	S
9	$\frac{3}{9}$	$\frac{6}{9}$	$\frac{0}{9}$
$9+3k$	$\frac{3+k}{9+3k}$	$\frac{6+k}{9+3k}$	$\frac{0+k}{9+3k}$

* Smoothing 2017 value
0 2018 etc.

* 2017 value missing 2018 etc.

* 2017 value missing 2018 etc.
value + k 2018 (that missing value + all values in table)

* For Laplace smoothing
value of k is always 1

$P(H=1|T=1) = P(T=1|H=1)P(H=1) / P(T=1)$
 $P(H=1|T=0) = P(T=0|H=1)P(H=1) / P(T=0)$
 $P(T=1|H=1) = P(H=1|T=1)P(T=1) / P(H=1)$
 $P(T=0|H=1) = P(H=1|T=0)P(T=0) / P(H=1)$

$$P(H) = \frac{9}{10}, P(T) = \frac{1}{10}$$

$$\frac{P(H=1|T=1)P(H=1) / P(T=1)}{P(H=1|T=0)P(H=1) / P(T=0)}$$

likelihood prior
P(H=1|T=1) / P(H=1|T=0)

$$= \frac{10}{1} \cdot \left(\frac{1}{10}\right)^9 \left(\frac{1}{10}\right)^1 \left(\frac{9}{10}\right)^9$$

$\int_0^1 10 \cdot \left(\frac{1}{10}\right)^9 \left(1 - \frac{1}{10}\right)^1 d\pi$ for update
 $\pi = 0$ is a prior to π
 $\pi = 1$ is a prior to π

* previous idea (prior) 20/20 new observation (likelihood) 33
 20/20, 33 new observation main date 0, update 20/23
 0.87. This is ~~not~~ online learning.

DATA AND PREDICTION

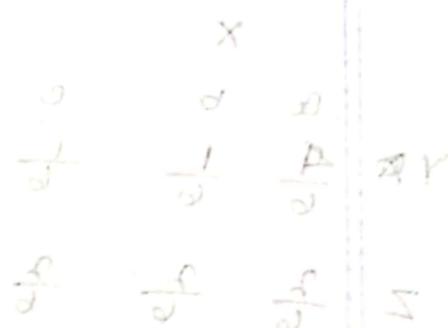
E	O	A	P
0	0	0	P
1	0	1	P
0	1	0	P
1	1	1	P

E	O	A	P
0	0	0	P
1	0	1	P
0	1	0	P
1	1	1	P

SUNDAY

DATE: 09/07/23

ENTROPY

 $H = T$ $S_0 \rightarrow \text{uncertainty about weather}$ $90 \rightarrow 10$ $10 \rightarrow 90$ 

Dice :-

1	2	3	4	5	6
$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

\rightarrow outcome
 \rightarrow P is equally divided among all outcomes
 $\rightarrow H = -\sum p_i \log p_i = [P]H$

Probability Distribution

$$H = -\sum p_i \log p_i = -f_1 \log f_1 - f_2 \log f_2 - f_3 \log f_3 - f_4 \log f_4 - f_5 \log f_5 - f_6 \log f_6 = [P]H$$

\rightarrow $\frac{1}{p}$ event ২০৫৩ chance এফি, তার p আৰু
 value এফি, $\frac{1}{p}$ আৰু value বিটা
 \rightarrow $\frac{1}{p}$ event ২০৫৩ chance ২০২৫, তার p
 আৰু value ২০২৫,
 $\frac{1}{p}$ আৰু value অন্তক আৰু

$\log(\frac{1}{p}) \rightarrow$ To decrease a very large value of $\frac{1}{p}$ if p is very small

$\rightarrow -\log(p)$

* $I(X) = -\log P(X)$

\hookrightarrow self information

$$E[\text{outlook}] = \frac{3}{14}(-\log \frac{3}{14}) + \frac{1}{14}(-\log \frac{1}{14}) + \frac{1}{14}(-\log \frac{1}{14})$$

\rightarrow Probability of event raining
 \rightarrow self information of raining

* Entropy is the expected value and is denoted by $H[\text{outlook}]$

$$P(\text{Rainy}) = \frac{3}{14}$$

$$P(\text{Overcast}) = \frac{1}{14}$$

$$P(\text{Sunny}) = \frac{10}{14}$$

Entropy of 2 min

RAMADEV

RAMADEV

X

	a	b	c
y	$\frac{4}{6}$	$\frac{1}{6}$	$\frac{1}{6}$
z	$\frac{2}{6}$	$\frac{2}{6}$	$\frac{2}{6}$

* Uncertainty of

T	H
0.1	0.2
0.1	0.2
0.1	0.1

$$H[Y] = \frac{4}{6}(-\log \frac{1}{6}) + \frac{1}{6}(-\log \frac{1}{6}) + \frac{1}{6}(-\log \frac{1}{6}) = 1.28$$

$$H[Z] = \frac{2}{6}(-\log \frac{2}{6}) + \frac{2}{6}(-\log \frac{2}{6}) + \frac{2}{6}(-\log \frac{2}{6}) = 1.58$$

* Uncertainty of Z > Y as $H[Z] > H[Y]$

* we need small uncertainty so entropy need to be small

Decision

DECISION TREE

Overcast

ON 44

* No entropy
as we are
sure we
can play

if Y... then

Rainy

3N 2Y

Sunny

3Y 2N

* Some entropy [Entropy + 0]

$$\frac{3}{5}(-\log \frac{3}{5}) + \frac{2}{5}(-\log \frac{2}{5})$$

$$\approx 0.97$$

Leads to

* We are
finding out
entropy of
playoff with

respect to
outlook

outcomes

* Out entropy

of playoff w.r.t
outlook is

the smallest

compared to

Temp, Humidity,

windy.

Windy
 T F
 04 2N 34 ON
 24 IN 14 IN
 mild cool
 entropy 20 O 10 D

330 atm
 (A) 01 - J

qnet
 L f
 1000 800 500 Pa
 11 Pa 48 Pa 40 Pa
 (13.0) (10.0) (1)

1000 Pa
 40 Pa
 (10.0) (1)

March 20 2020 11:10
 (O₂-F) + (F₂)
 H₂O

(O₂-F) + (F₂)
 H₂O

Information gain:

$$\begin{matrix} \text{Q. Yes} & \text{No} \\ 9 & 5 \end{matrix}$$

	OUTLOOK	WEATHER
K	✓	✓
100% 60%	WINDY	WINDY

	OUTLOOK
K	T

* Dataset \Rightarrow split
lowest entropy \Rightarrow Outlook

Overall entropy of dataset

$$= \frac{9}{14}(-\log \frac{9}{14}) + \frac{5}{14}(-\log \frac{5}{14})$$

$$= 0.94$$

OUTLOOK		
Sunny	↓	Rainy
34 2N		24 3N
<u>Ent(0.97)</u>		<u>(0.97)</u>

OUTLOOK		Outlook
Sunny	Rainy	Outlook
34 2N	24 3N	44 0N
<u>Ent(0.97)</u>	<u>(0.97)</u>	<u>(0)</u>

Temp		
Hot	mild	Cool
24 2N	44 2N	34 1N
<u>(1)</u>	<u>(0.97)</u>	<u>(0.81)</u>

Overall entropy of outlook dataset

$$= \left(\frac{3}{14} * 0.97\right) + \left(\frac{6}{14} * 0.97\right) + \left(\frac{5}{14} * 0\right)$$

$$= 0.69 \rightarrow \text{outlook from dataset}$$

\Rightarrow split \Rightarrow Outlook
Entropy \Rightarrow 0.69

Overall entropy of dataset

$$= \left(\frac{1}{14} * 1\right) + \left(\frac{6}{14} * 0.92\right) + \left(\frac{7}{14} * 0\right)$$

$$= 0.911$$

* \Rightarrow data \Rightarrow different values in a column, so entropy smallest \Rightarrow But, \Rightarrow we use \Rightarrow \Rightarrow reason. For that reason, we use information entropy.

Information gain = amount of reduction in entropy

$$= \text{Previous entropy} - \text{split entropy}$$

$$\text{e.g. } \text{Information Gain (Outlook)} = 0.94 - 0.69$$

$$\text{Information Gain (Temp)} = 0.94 - 0.911 \quad \text{O/T}$$

Split Entropy

* If there are many values in a column, then the frequency of that value will be 1. So, split entropy will be highest.

$$\begin{aligned} & -\frac{1}{n} \log \frac{1}{n} - \frac{1}{n} \log \frac{1}{n} - \dots \quad \text{(approx)} = -\frac{1}{n} \log \frac{1}{n} \quad \text{Dilog} \quad \text{Dilog} \\ & = n * \frac{1}{n} \log n \\ & = \log n \end{aligned}$$

* Gain ratio ମଧ୍ୟ ସଫାର
ତାଙ୍କେ splitting A
consider, 2023.

$$\text{Gain Ratio}(x) = \frac{\text{Information Gain}(x)}{\text{split Entropy}(x)}$$

$$P_1 + P_2 + \dots + P_n = 1$$

$$-P_1 \log P_1 - P_2 \log P_2 - \dots - P_n \log P_n$$

$$0 \leq -P_1 \log P_1 - P_2 \log P_2 - \dots - P_n \log P_n \leq \log \left(\frac{1}{P_1} + \frac{1}{P_2} + \dots + \frac{1}{P_n} \right) = H$$

Gini Impurity

↳ Index

$$P_1^2 + P_2^2 + P_3^2 + \dots + P_n^2$$

$$1 - (P_1^2 + P_2^2 + P_3^2 + \dots + P_n^2)$$

$$= 1 - \left(\frac{1}{n^2} + \frac{1}{n^2} + \dots + \frac{1}{n^2} \right)$$

$$= 1 - \frac{1}{n}$$

→ usually n is value 2

$$= \frac{n-1}{n}$$

$$\lim_{n \rightarrow \infty} \frac{n-1}{n} = 1$$

[if n is very large]

$$P_i \leq P_i \mid -P_i^2 \geq -P_i$$

$$1 - (P_1^2 + P_2^2 + P_3^2 + \dots + P_n^2)$$

$$\Rightarrow 1 - P_1^2 - P_2^2 - P_3^2 - \dots - P_n^2$$

$$\Rightarrow 1 - (P_1 + P_2 + P_3 + \dots + P_n)$$

$$1 - (P_1 + P_2 + \dots + P_n)$$

$$1 - 0.5 = 0.5$$

$$\text{impurity} \leq 0.5$$

Maximum likelihood estimation

Linear Classification \rightarrow Perceptron Algo

Linear Regression \rightarrow Perceptron Algo

\rightarrow Gradient descent
 \rightarrow Least square Sol?

Naive Baye's Classifier

Decision Trees

Titanic data -

Supervised learning - classification, ranking, regression

Classification - assigning a label to each data instances

Regression - linear → used for continuous values

- logistic → almost linear

Categorical Data - used for prediction

- whenever we use categorical data
it is classification

- Random Categories
- e.g. 0/1, Yes/No

Ordinal Value - when one value is better than other

Rent for a house is linearly dependent on all the factors influencing the rent.

$$Y = \alpha_1 K_1 + \alpha_2 K_2 + \alpha_3 K_3 + \alpha_4 K_4 + \dots + \alpha_b K_b + \dots + \beta + \phi$$

↳ linearity dependency [power of all x are 1]

* variable 2051, equation Rfirst - linear dependent 23/25
change Rfirst

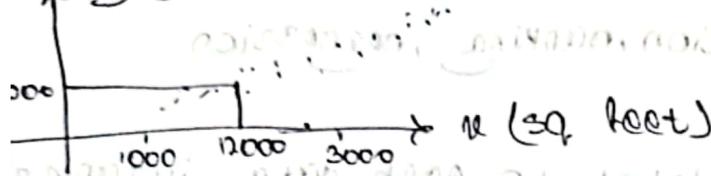
* variable Rfirst, equation 2051 - no sol?

* variable & equation are equal - linear dependent 21/25
0113 तात्पर्य, $X^T X = I$

$$X^T X (2 \times 2) = 4 \times I$$

$$y = \alpha_1 x + b$$

y (rent)



- non diff

coefficient, gradient descent - gradient descend

gradient descent - gradient descend

n = data points / features / $\alpha_1, \alpha_2, \alpha_3$ = weights / variables / parameters

feature dimensions for $\alpha_1, \alpha_2, \alpha_3$ = 1000, 2000, 3000 → 1000, 2000, 3000

$$y_1 = \underbrace{\alpha_1 x_{1,1}}_{\text{Data 1}} + \underbrace{\alpha_2 x_{1,2}}_{\text{for 2nd}} + \underbrace{\alpha_3 x_{1,3}}_{\text{for 3rd}} + \dots + b_1 * \underbrace{\text{noise}}_{\text{1st data}}$$

A	B	C
a_1	$a_2 b_1$	c_1

add 1000, 2000, 3000 → 1st data
multiplication of 1st data

$$y_2 = \alpha_1 x_{2,1} + \alpha_2 x_{2,2} + \alpha_3 x_{2,3} + \dots + b_2 * 1$$

$$y_m = \alpha_1 x_{m,1} + \alpha_2 x_{m,2} + \alpha_3 x_{m,3} + \dots + b_m * 1$$

x → feature matrix

w → feature vector

y_1, y_2, \dots, y_m → target values

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,n-1} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m,1} & x_{m,2} & \dots & x_{m,n-1} \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \\ b \end{bmatrix}$$

homogeneous to make it non-homogeneous
for b variable, we have to append a column with 1 which makes the matrix non-homogeneous from homogeneous.

non-homogeneous

target values

$$y = Xw$$

This is off-line learning

$$\Rightarrow Xw = y$$

→ $X^T X w = X^T y$

$$\Rightarrow X^T X w = X^T y$$

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

→ least squares rule

* This is off-line learning

$\begin{bmatrix} 3 & 6 \\ 2 & 4 \end{bmatrix}$ Determinant = $(3 \cdot 4) - (6 \cdot 2) = 0$, so can't invert

$$\begin{bmatrix} 3 + 0.003 & 6 \\ 2 & 4 + 0.003 \end{bmatrix}$$

$$\underbrace{\begin{bmatrix} 3 & 6 \\ 2 & 4 \end{bmatrix}}_X + \begin{bmatrix} 0.003 & 0 \\ 0 & 0.003 \end{bmatrix} + 0.003 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

\downarrow
identity matrix

$$x + 0.003]$$

$$x w = y$$

$$\Rightarrow x^T x w = x^T y$$

$$\Rightarrow (x^T x + \lambda I) w = (x^T y)$$

Error:-

$y \rightarrow$ actual

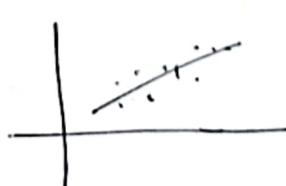
$\bar{y} \rightarrow$ approximate

$$(y_1 - \bar{y}_1)^2 \text{ or } |y_1 - \bar{y}_1|$$

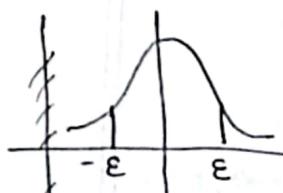
\downarrow
better

$$(y_1 - \bar{y}_1)^2 + (y_2 - \bar{y}_2)^2 + \dots + (y_K - \bar{y}_K)^2$$

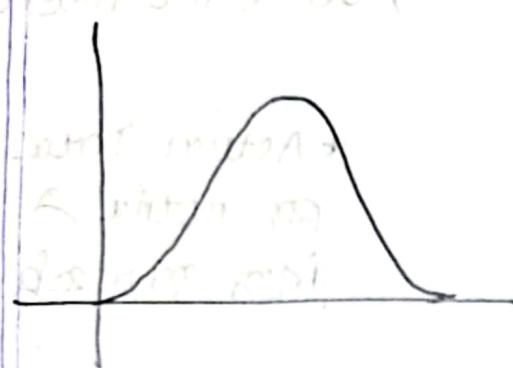
\downarrow
mean square error



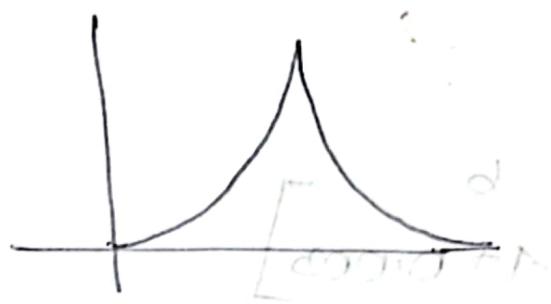
$$\epsilon_1 = y_1 - \bar{y}_1$$



+ve & -ve
error are of
same
probability



Gaussian Distribution



Laplace Distribution

$$\int_{-\infty}^0 e^{-|x|/\lambda} dx + \int_0^\infty e^{-|x|/\lambda} dx$$

$$\left[e^{-|x|/\lambda} \right]_{-\infty}^0 + \left[e^{-|x|/\lambda} \right]_0^\infty$$

cancel with

$$[e^{-|0|/\lambda} + e^{-|0|/\lambda}]$$

$$\begin{aligned} P &= \\ P^T x &= c x \\ (P^T x) &= c_0(I b + x T) \end{aligned}$$



$y = 1$

$$1/2 - 1/2 = 0$$

$$1/2 - 1/2 = 0$$

LAB

2

SATURDAY

DATE: 10/06/23

$$\text{Q2} \sum_{i=1}^m (\hat{y}_i - y_i)^2 \rightarrow \text{mean square error (MSE)}$$

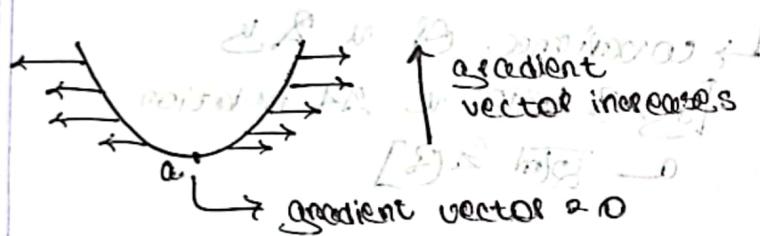
$\hat{y}_i = ax + b$

$$J = \sum_{i=1}^m (\hat{y}_i - y_i)^2 \rightarrow \text{sum of squared error}$$

$\hat{y}_i = ax + b$

$$J = \sum_{i=1}^m (\hat{y}_i - (ax_i + b))^2$$

$$\text{gradient vector: } \nabla J = \begin{bmatrix} \frac{\partial J}{\partial a} \\ \frac{\partial J}{\partial b} \end{bmatrix}$$



$$\nabla J_2 = \left[\frac{\partial J}{\partial a}, \frac{\partial J}{\partial b} \right]$$

$$\text{Total change} = \Delta J = \|\nabla J\| \cdot \Delta \vec{v}$$

$$= |\Delta \vec{v}| |\nabla J| \cos \theta$$

$$\text{Total change} = \Delta J = \nabla J \cdot \Delta \vec{v}$$

$$= \frac{\partial J}{\partial a} \Delta a + \frac{\partial J}{\partial b} \Delta b$$

$$\text{At maximum/minimun point, } \frac{\partial J}{\partial a} = 0, \frac{\partial J}{\partial b} = 0$$

$$\frac{\partial J}{\partial a} = \sum_{i=1}^m 2(\hat{y}_i - ax_i - b) (-x_i) = 0$$

$$\frac{\partial J}{\partial b} = \sum_{i=1}^m 2(\hat{y}_i - ax_i - b) (-1) = 0$$

* $A \cdot B = |A||B|\cos\theta$
 * Gradient ∇J
 direction $\nabla J / |\nabla J|$
 maximum value of $\cos\theta$
 মানুষ মানুষ
 দেখতে চায়

* Gradient ∇J opposite
 direction $-\nabla J / |\nabla J|$
 minimum value of $\cos\theta$
 মানুষ মানুষ

$$\frac{\partial J}{\partial a} = \sum_{i=1}^m (x_i y_i - ax_i^2 - bx_i) = 0 \rightarrow \frac{\partial J}{\partial a} = \sum_{i=1}^m (x_i y_i - a x_i^2 - b x_i) = 0 \rightarrow a \sum_{i=1}^m x_i y_i - a \sum_{i=1}^m x_i^2 - b \sum_{i=1}^m x_i = 0$$

$$\frac{\partial J}{\partial b} = \sum_{i=1}^m (x_i y_i - ax_i - b) = 0 \rightarrow \frac{\partial J}{\partial b} = \sum_{i=1}^m (x_i y_i - a x_i - b) = 0 \rightarrow a \sum_{i=1}^m x_i + b m = \sum_{i=1}^m y_i$$

Using Cramer's Rule,

$$(1) \begin{vmatrix} \sum x_i y_i & \sum x_i \\ \sum x_i & \sum m \end{vmatrix} \rightarrow \frac{+ve - \cancel{m} \rightarrow \text{L.R. } \uparrow, \downarrow \quad \cancel{m}}{-ve \rightarrow \text{L. } \downarrow, \downarrow} \quad \frac{a_2 m \sum x_i y_i - (\sum x_i)(\sum y_i)}{m \sum x_i^2 - (\sum x_i)^2}$$

coefficient of $a \& b$
~~[$\text{L.R. BMATL R. L.R. relation}$]~~
 $a \rightarrow \frac{\partial J}{\partial a} \in \{3\}$

$$b = \frac{\begin{vmatrix} \sum x_i y_i & \sum x_i y_i \\ \sum x_i^2 & \sum y_i \end{vmatrix}}{\begin{vmatrix} \sum x_i^2 & \sum x_i \\ \sum x_i & m \end{vmatrix}} \rightarrow b = \frac{\sum x_i y_i - (\sum x_i)^2 (\sum y_i)}{(\sum x_i^2) m - (\sum x_i)^2}$$

Question 2
Solve the given Q.

Method of gradient descent
~~gradient descent~~
~~gradient descent~~
 Method of gradient descent
 Gradient descent
 Gradient descent
 Gradient descent

* gradient descent
~~gradient descent~~
~~gradient descent~~

$$\alpha \cdot \left(\frac{\partial J}{\partial a} \right) \left(\text{forward pass} - \text{backward pass} \right) + \frac{\partial J}{\partial a} = \frac{\partial J}{\partial a}$$

$$\alpha \cdot \left(\frac{\partial J}{\partial a} \right) \left(\text{forward pass} - \text{backward pass} \right) + \frac{\partial J}{\partial a} = \frac{\partial J}{\partial a}$$

$$\begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_n \\ b \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_n \\ b \end{bmatrix} + \lambda \begin{bmatrix} \frac{\partial J}{\partial a_1} \\ \frac{\partial J}{\partial a_2} \\ \frac{\partial J}{\partial a_3} \\ \vdots \\ \frac{\partial J}{\partial a_n} \\ \frac{\partial J}{\partial b} \end{bmatrix}$$

\downarrow
New a_i

Old a_i

b ~~Don't do gradient
Update and start again~~

$$\lambda \frac{\partial J}{\partial a}$$

must be very small so for maximum

Find max. min.

Max and min
position of gradient
step function

Max step gradient
Effect same as

$$\begin{array}{|c|c|} \hline & 1 & \\ \hline 0 & & 1 \\ \hline 1 & & 1 \\ \hline \end{array}$$

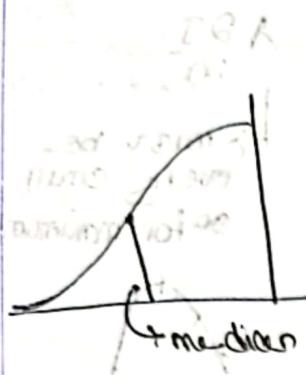
Initial position
Step function
Max step function
minimum step function

$$P(X^T Q X) \leq 0$$

P step function has a
minimum because P

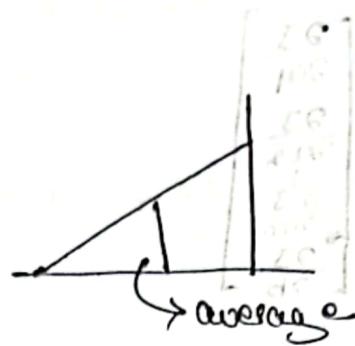
$$Q \geq 0$$

Step function minimum must be a global



ମେଡିଆ ଗ୍ରାଫ ହୁଅ,
median କେ ବନ୍ଦୀ
better for missing
inputs.

Average କେ କଣ୍ଠରେ
କ୍ଷେତ୍ର କେ ଆହୁ



ଅଭିନ୍ନ ଗ୍ରାଫ ହୁଅ,
average କେ ସମ୍ଭା
better for missing
inputs.

$$\hat{w}_0 + (\mathbf{x}^T \hat{\mathbf{w}})^T \mathbf{x}^T \mathbf{y}$$

\mathbf{x} → all columns except \mathbf{y}

\mathbf{y} → survived column

pandas → data frame manipulate ୧୦୯

\hat{w}_0	Not vector
$\hat{\mathbf{w}}$	Not vector

m	f
1	0
0	1

ଅଭିନ୍ନ କାଳେ
ହାତରିଲେ ୨୨୪
 $m = 1$ & $f = 0$
୨୭୩ ପାଇଁ ନାହିଁ

FINAL



LEARNER

POINTWISE PREDICT.

- * Learned: an algorithm which learns from a dataset

$$y = w + \epsilon$$

ϵ error
from normal distribution

- * A learner can find out the hyperplane giving optimal accuracy.

- * Realizability assumption: If a hyperplane can divide a dataset into separating (0 and x), it is realizable.

- * % of women in the world population \rightarrow True estimate

- * % of women in a state \rightarrow Sample estimate

$R_t(n) \rightarrow$ Error calculated over all data in the world
 \downarrow
 true

$R_s(n) \rightarrow$ Error calculated over all data in a given area.
 \downarrow
 sample

Weak learner: accuracy $\frac{1}{n} (2/20)$ (After for binary classification) error $\frac{1}{n} (2/20) (20/2)$

- * It is not possible to find out $R_t(n)$ directly.

DATA SHEET

DATA SHEET

Interval Function

DATA SHEET

STATEMENT OF THE PROBLEM OF INTERVAL FUNCTION



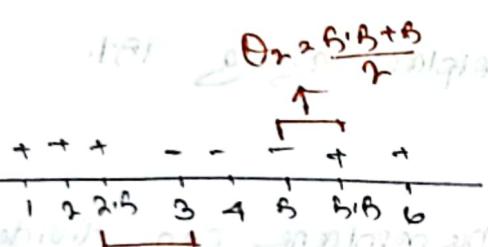
$\theta_1 < x < \theta_2$

ZONE +
BETWEEN MARK
CONTINUOUS

ZONE

LINE AND ZONE

++ - - - + + + -



STATEMENT OF THE PROBLEM OF INTERVAL FUNCTION

* NOT realizable, as there is - after θ_2 interval. So, cannot be solved using interval function

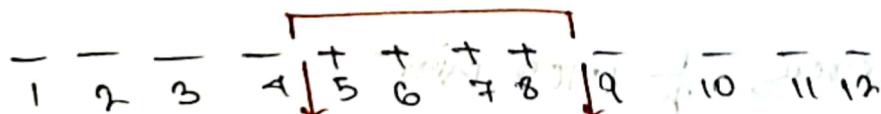
PROBLEMS WITH FOLLOWING PLANE

FUNCTION ASSUMED AS $f(x) = x$ IN RANGE $\theta_1 < x < \theta_2$

$$h_{\theta_1, \theta_2}(x) = \begin{cases} -b & \text{if } \theta_1 \leq x \leq \theta_2 \\ +b & \text{otherwise} \end{cases}$$

$$\theta_1 < x \leq \theta_2 = +1$$

$$-b = +1$$

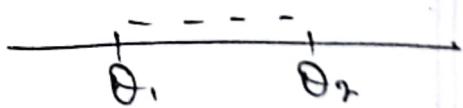


$$h_{4, 9, 8, 12, -1}(x)$$

$$\text{for } -b = +1$$

STATEMENT OF THE PROBLEM OF INTERVAL FUNCTION

$$\text{for } -b = -1$$



BOOSTING

INTERVAL FUNCTION

* Boosting : UNILAD learned 2012, UNILAD opinion 2012, 321200
relevant functions from decision tree

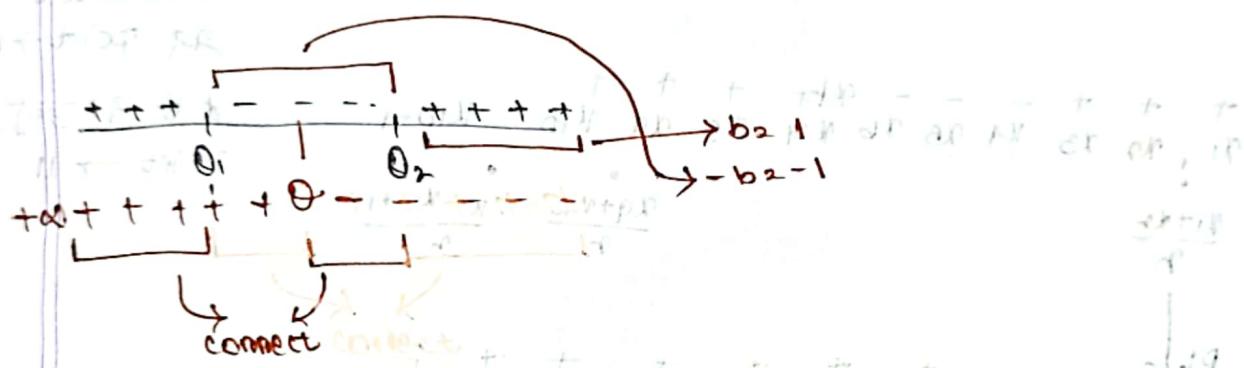
interval function

threshold function

$$f_{0,b}(x) = \text{sign}(\theta - x)$$

$$b \in \{-1, +1\}$$

accuracy



* Accuracy is not 100%, so we have to choose θ accurately.

$$\frac{\theta + \theta_1}{\theta_1} \quad \frac{\theta - \theta_2}{\theta_2} \quad [\text{interval function}]$$

$$+\infty \quad \theta = \infty$$

$$-\infty \quad \theta = -\infty$$

$$-\infty \quad \theta = 0$$

$$+\infty \quad \theta = \infty$$

$$\frac{\theta - \theta_1}{\theta} \quad \frac{\theta + \theta_2}{\theta} \quad \rightarrow \text{Here, accuracy is } 67.1\% \text{ which is poor.}$$

This is weak learner.
(1) same as standard, if

* If we use many weak learners; then we can find a better accuracy \rightarrow Ensemble learning.

(1) same as standard, if

* NL (to) increasingly
soot ဆွဲတဲ့ အမှုပါ။
[soot in descending
order]

+ + + - - - + + + +
x1 x2 x3 x4 x5 x6 x7 x8 x9 x10

Older people, 70-79 years, 30-49% of 1990

D1 D2

2018 (F) 0128

ଏହା କିମ୍ବା ଏହା କିମ୍ବା ୧୯୨୩ ମୁଦ୍ରଣ କରିଲା

ଶ୍ରୀମଦ୍ଭଗବତ

important

animal husbandry

original document

* Dataset

23 point → 10

~~7~~ 88-83-1

ପାତ୍ରବିଦ୍ୟା

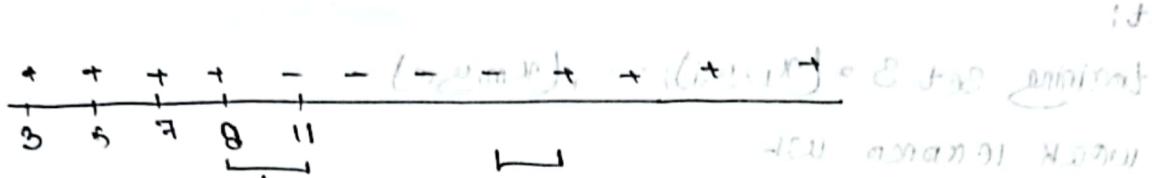
point → II
right right + O + f + t + out

$$\begin{array}{c|c|c} a_{j+1} & a_i & a_{i+2} \\ \hline & + & \\ & . & \\ & \theta_i & \theta_{i+1} \\ \hline \vdash & - & - \end{array}$$

[2015-16] + + + - - - + + +

* Actual output is +ve and after checking using θ_i , we get -ve. By shifting θ_i to θ_{i+1} , we get +ve for θ_i . So, there's no error (Δ_i).

* Actual output of n is +ve and after using D_1 , n gets -ve. By shifting D_1 to D_{1+V} , we get +ve for n , so, there's an error (D_1)



$\Theta = 9.5$ & review $f(x) = \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$

For threshold function,

$$h_{\theta, b}(x) = b \cdot \text{sign}(\theta - x)$$

(θ, b) for each neuron

$$\left[\sum_{i=1}^m w_i x_i + b \right] / \sum_{i=1}^m w_i = f_\theta \text{ output}$$

if $x > \theta$,
 $\text{sign}(\theta - x)$
is -ve.

$$\text{ReLU}(x) = \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

$$\frac{\partial L}{\partial \theta} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial \theta} = \frac{\partial L}{\partial z} \cdot \text{ReLU}'(z) = \frac{\partial L}{\partial z} \cdot \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

$\theta_i = \infty, b = 1 \rightarrow$ for threshold
 $\theta_i = \frac{x_i + x_{i+1}}{2}$

* To increase the accuracy, we have to find θ_i upto maximum errors are resolved.

To solve $L(x_i(z)) + \lambda - \lambda \cdot \text{ReLU}'(z) = 0$

$$\lambda \leq L'(x_i(z)) + \lambda - \lambda \cdot \text{ReLU}'(z)$$

What is $L'(x_i(z))$?

ReLU'(x) = 1

ReLU'(x) = 0

0 part

What is $L'(x_i(z))$?

ReLU'(x) = 1

ReLU'(x) = 0

0 part

22/07/2021

10:00 AM

ADABOOST

Input:

training set $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$

weak learners W_t

number of rounds $T \rightarrow$ no. of weak learners
→ ~~no. of iteration~~

initialize $D^{(0)} = (\frac{1}{m}, \frac{1}{m}, \dots, \frac{1}{m})$

for $t = 1, \dots, T$:

invoke weak learners $h_t = W_t(D^{(t)}, S)$

compute $\epsilon_t = \sum_{i=1}^m D_i^{(t)} I_{[y_i \neq h_t(x)]}$

let $w_t = \frac{1}{2} \log(\frac{1}{\epsilon_t} - 1)$

update $D_i^{(t+1)} \leftarrow D_i^{(t)} e^{-w_t y_i h_t(x_i)}$
→ next iteration

data

$$\sum_{j=1}^m D_j^{(t)} e^{-w_t y_j h_t(x_j)}$$

for all $i = 1, \dots, m$

Output the hypothesis $h_S(x) = \text{sign}(\sum_{t=1}^T w_t h_t(x))$

$$D_i^{(t+1)} = D_i^{(t)} \left[e^{-y_i h_t(x_i) w_t} \right] \quad \begin{matrix} \text{Actual} \\ \text{weight of} \\ \text{data} \end{matrix}$$

$$\text{Error } \sum_{i=1}^m (-y_i h_t(x_i) w_t) D_i^{(t)} \quad \text{if } 0$$

$$\text{if } \frac{1}{\epsilon_t} - 1 < 0 \\ \downarrow \\ \text{must be non-ve}$$

w_t must be 0 or greater than 0.

* classifier (A) data weight ≥ 0
BCE, OK data to OR classifier

$$D_i^{(t+1)} \leftarrow D_i \stackrel{t}{\leftarrow} e^{-h_i h_t(x_i) w_t}$$

$$\frac{1}{m} \sum_{i=1}^m \frac{(x_i + \theta)^{-}}{e^{(x_i + \theta)^{-}} + e^{(x_i + \theta)^{+}}} = \frac{1}{m} \sum_{i=1}^m \frac{1}{e^{(x_i + \theta)^{-}} + 1}$$

$$D_1^{(t+1)} \leftarrow D_1 \stackrel{t}{\leftarrow} e^{-h_1 h_t(x_1) w_t}$$

Total of all $D_i = \frac{1}{m} \sum_{i=1}^m \frac{1}{e^{(x_i + \theta)^{-}} + 1} = \frac{1}{m} = \frac{1}{m} = \dots$

$$h_1(11) = + \quad h_2(11) = - \quad h_3(11) = +$$

$$(+)10 + (-5)(-1) + (+7)(+1) = +4$$

\downarrow
step weight
 $n(t)$ w_t

$$\frac{(x_i + \theta)^{-}}{(x_i + \theta)^{-} + e^{(x_i + \theta)^{+}}} = \frac{1}{e^{(x_i + \theta)^{-}} + 1}$$

$$R_S(h_{AB}) = \frac{1}{m} \sum_{i=1}^m I_{h_{AB}(x_i) \neq y_i} \stackrel{-2\lambda T}{\leftarrow} e^{-2\lambda T} \rightarrow \text{error } 2\lambda T \quad R_S(h_{AB}) \approx 1$$

$$\frac{\partial}{\partial \theta} \frac{(x_i + \theta)^{-}}{(x_i + \theta)^{-} + e^{(x_i + \theta)^{+}}} = \frac{1}{e^{(x_i + \theta)^{-}} + 1} \quad [\text{increasing } T, \text{ decreases } R_S(h_{AB})]$$

$$\rightarrow R_S(h_{AB}) = 1 \quad (\text{for error})$$

$$f_T(x) = w_0 h_0(x) + w_1 h_1(x) + \dots + w_T h_T(x) \rightarrow \text{total error}$$

$$h_0(x) = 0$$

$$h_1(x) = w_1 h_1(x)$$

$$h_2(x) = w_2 h_2(x) + w_3 h_3(x)$$

$$\frac{(x_i + \theta)^{-}}{(x_i + \theta)^{-} + e^{(x_i + \theta)^{+}}} = \frac{1}{e^{(x_i + \theta)^{-}} + 1} \quad \frac{(x_i + \theta)^{+}}{(x_i + \theta)^{-} + e^{(x_i + \theta)^{+}}} = \frac{e^{(x_i + \theta)^{+}}}{e^{(x_i + \theta)^{-}} + 1}$$

$$\hat{y}_t = \frac{1}{m} \sum_{i=1}^m e^{-\gamma_i f_t(x_i)} \gamma_i \frac{1}{m}$$

$$f_t(x) = \frac{\sum_{i=1}^m e^{-\gamma_i f_t(x_i)}}{\sum_{i=1}^m e^{\gamma_i f_t(x_i)}}$$

$$f_t(x) = \frac{\sum_{i=1}^m e^{-\gamma_i f_t(x_i)}}{\sum_{i=1}^m e^{\gamma_i f_t(x_i)}}$$

$$\frac{1}{m} \sum_{i=1}^m e^{-\gamma_i f_t(x_i)} + \frac{1}{m} = \frac{\frac{1}{m} \sum_{i=1}^m e^{-\gamma_i f_t(x_i)}}{m} + \frac{1}{m} = \text{id} \text{ in } \mathbb{R} \text{ MGF}$$

$$F = (W) \otimes \dots \otimes (W) \otimes \dots \quad r = (W) \otimes \dots$$

$$F + e^{(W)(F)} + (1 - e^{(W)}) + (W)(F)$$

\downarrow
Simplifying
 \downarrow
 \downarrow

$$D_i^{(t+1)} = \frac{e^{-\gamma_i f_t(x_i)}}{\sum_{j=1}^m e^{-\gamma_j f_t(x_j)}}$$

$$D_i^t = \frac{e^{-\gamma_i f_{t-1}(x_i)}}{\sum_{j=1}^m e^{-\gamma_j f_{t-1}(x_j)}}$$

$$1 \leq \frac{1}{m} \sum_{i=1}^m e^{-\gamma_i f_{t-1}(x_i)} \leq 1 \quad \text{and} \quad \sum_{i=1}^m \frac{1}{m} \cdot (\text{const}) \otimes$$

$$e^{-w_t \gamma_i f_t(x_i)} D_i^t = \frac{e^{-w_t \gamma_i f_t(x_i)}}{\sum_{j=1}^m e^{-\gamma_j f_{t-1}(x_j)}}$$

$$f_t(x_i) = f_{t-1}(x_i) + w_t \gamma_i f_t(x_i)$$

target \leftarrow const +

$$f_t(x) = \frac{w_{\min}(x) + \dots + w_{\max}(x)}{m} + f_{t-1}(x)$$

$$e^{-w_t \gamma_i f_t(x_i)} D_i^t = \frac{e^{-\gamma_i f_t(x_i)}}{\sum_{j=1}^m e^{-\gamma_j f_t(x_j)}}$$

$$\frac{f_{t-1}(x_i)}{f_t(x_i)} = \frac{w_{\min}(x_i) + \dots + w_{\max}(x_i)}{w_t \gamma_i f_t(x_i)}$$

$$(x_i) \text{ and } (x_i) \cdot$$

$$(x_i) \text{ and } (x_i) \cdot$$

SUPPORT VECTOR MACHINE

- binary classification

 $A \rightarrow A \text{ or not } A \rightarrow \text{classifying}$

[Dataset 2 S1(GT STAT 2013)]

A

B

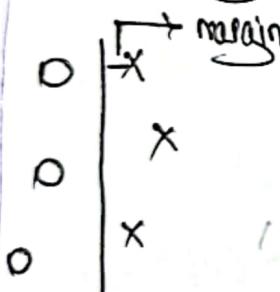
C

D



Margin = $\min_{i=1}^n \frac{|w_i^T x_i + b|}{\|w\|}$

SVM - margin to maximise

 $\forall i \in [m], y_i(w_i^T x_i + b) \geq 1$

OR

(With $w_i^T w_j = 0 \Rightarrow w_i^T w_j + b = 0$)So $y_i(w_i^T x_i + b) \geq 1$ for all i in the training set.

1 is the margin

$$y_1(w_1x_{11} + w_2x_{12} + \dots + b) = \gamma_1$$

$$y_2(w_1x_{21} + w_2x_{22} + \dots + b) = \gamma_2$$

$$y_j(w_1x_{j1} + w_2x_{j2} + \dots + b) = \gamma_j$$

(for j=1, 2, ..., m)

$$y_m(w_1x_{m1} + w_2x_{m2} + \dots + b) = \gamma_m$$



$$y_1\left(\frac{w_1}{\gamma_j}x_{11} + \frac{w_2}{\gamma_j}x_{12} + \dots + \frac{b}{\gamma_j}\right) = \frac{\gamma_1}{\gamma_j} \gamma_1$$

$$y_2\left(\frac{w_1}{\gamma_j}x_{21} + \frac{w_2}{\gamma_j}x_{22} + \dots + \frac{b}{\gamma_j}\right) = \frac{\gamma_2}{\gamma_j} \gamma_1$$

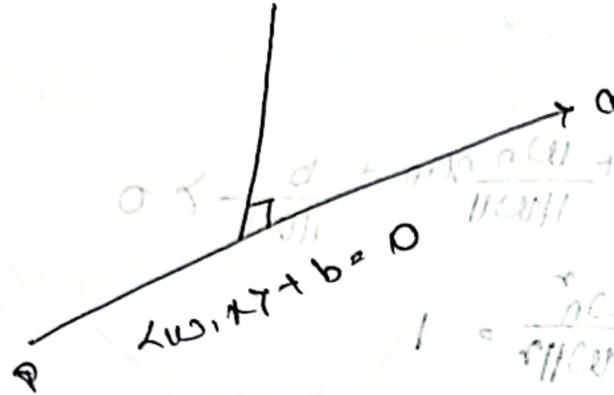
$$y_j\left(\frac{w_1}{\gamma_j}x_{j1} + \frac{w_2}{\gamma_j}x_{j2} + \dots + \frac{b}{\gamma_j}\right) = 1$$

$$y_m\left(\frac{w_1}{\gamma_j}x_{m1} + \frac{w_2}{\gamma_j}x_{m2} + \dots + \frac{b}{\gamma_j}\right) = \frac{\gamma_m}{\gamma_j} \gamma_1$$

$$y_i(w_{ii}x_{ii} + b) \leq 1$$

$\Rightarrow (d + \epsilon \sin \alpha) \leq 1$

* Divide by γ_j for smallest value. If $w_1 > 0$, values greater than smallest value would be greater than 1.



$$\langle w, p \rangle + b = 0$$

$$\langle w, q \rangle + b = 0$$

从图中得 $\|w\| \cdot \|q\| \cos \theta + b = 0$ 且 $\|w\| \cdot \|p\| \cos \theta + b = 0$

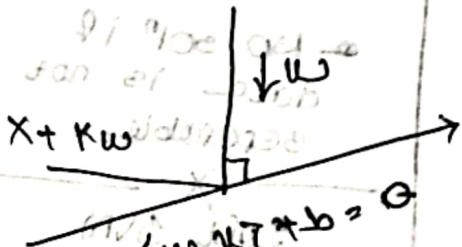
$$\langle w, q \rangle = \langle w, p \rangle$$

$$l = l(d + \langle x, w \rangle)$$

$$w_1 p_1 + w_2 p_2 + \dots + w_n p_n = w_1 q_1 + w_2 q_2 + \dots + \frac{b}{\|w\|} = \frac{l(d + \langle x, w \rangle)}{\|w\|}$$

$$\langle w, \overrightarrow{pq} \rangle = 0$$

$$w_1 w_2 \dots$$



$$d^2 = \langle Kw, Kw \rangle$$

$$d^2 = K^2 \langle w, w \rangle$$

$$d^2 = \frac{K^2 \|w\|^2 \|w\|^2 |b + \langle w, x \rangle|}{\|w\|^2} \Rightarrow d^2 = \frac{|b + \langle w, x \rangle|}{\|w\|^2}$$

$$\langle w, x + Kw \rangle + b = 0$$

$$\langle w, x \rangle + \langle w, Kw \rangle + b = 0$$

$$\langle w, x \rangle + \langle w, w \cdot K \rangle + b = 0$$

$$\langle w, x \rangle + K \|w\|^2 + b = 0$$

$$K \|w\|^2 = -b - \langle w, x \rangle$$

$$\|w\| = \sqrt{w_1^2 + w_2^2 + \dots + w_n^2} = 1$$

$$\frac{w_1}{\|w\|}x_{11} + \frac{w_2}{\|w\|}x_{12} + \dots + \frac{w_n}{\|w\|}x_{1n} + \frac{b}{\|w\|} > 0$$

$$\frac{w_1^2}{\|w\|^2} + \frac{w_2^2}{\|w\|^2} + \dots + \frac{w_n^2}{\|w\|^2} = 1$$

* w and x are scaled so that margin is always 1.

$$|w \cdot x + b| = 1$$

for

$$d = \frac{|w \cdot x + b|}{\|w\|} = \frac{1}{\|w\|}$$

margin

$$(p.w) = 1 - d$$

Hard SVM

linearly separable

points are

separable from it

- points are separable

- NO sol if data is not separable

Soft SVM

- gives sol even if data is inseparable

- points are not greater than 1 so we use $1 - \xi_i$

$$\frac{w \cdot x + b + \xi_i}{\|w\|}$$

$$\frac{w \cdot x + b - \xi_i}{\|w\|}$$

Soft SVM

$$y_i(\mathbf{x}_i^T \mathbf{w} + b) \geq 1 - \xi_i \quad \text{Hard SVM} \quad \text{if } \xi_i = 0 \\ \rightarrow \text{Soft SVM} \quad \text{if } \xi_i > 0$$

$$\lambda \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i$$

normal + violations
to minimize

$$\frac{\lambda \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i}{\sum_{i=1}^m \xi_i}$$

to minimize

• Constrained optimization

$$\text{def method} \leftarrow \nabla(\bar{F} - \mu \bar{V}) \sum_{i=1}^m \frac{1}{m} \cdot \mathcal{L}$$

• Minimizing \bar{F} with respect to \mathbf{w} → def gradient

• Minimizing \bar{V} with respect to b → def gradient

• Minimizing \bar{V} with respect to ξ → def gradient

• def gradient

gradient descent [the first method] • def gradient

$$\nabla(\bar{F} - \mu \bar{V}) \sum_{i=1}^m \frac{1}{m} \cdot \mathcal{L}$$

• def gradient

$$\nabla(\bar{V} - \mu \bar{F}) \sum_{i=1}^m \frac{1}{m} \cdot \mathcal{L}$$

• def gradient

• def gradient

Error \Leftrightarrow ઓછું માન $(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b))$ એ કેપ્રેસિંગ
ટૉર.

$$\text{man}(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b)) =$$

$\overbrace{\quad \quad \quad}$
 \downarrow
error

$$\text{man} \sum_{i=1}^m \frac{1}{m} + \text{Manifolds}$$

\downarrow
positive + negative
bias deviation

$$J = \frac{1}{m} \sum_{i=1}^m (y_i - \bar{y}_i)^2 \rightarrow \text{Batch GD}$$

Batch GD — સાથે સાથે એવી પ્રક્રિયા વિના

Minibatch GD — રાખી રાખી ડાટા

Stochastic GD — એક એક પ્રક્રિયા વિના

← Stochastic GD → [અનુભૂતિ માટે ના]
estimated GD

$$J = \frac{1}{m} \cdot \boxed{m} (y_{1n} - \bar{y}_{1n})^2$$

\downarrow
m ટાઈ માટે ગુણી

Minibatch GD

$$J = 10 \sum (y_{1n} - \bar{y}_{1n})^2$$

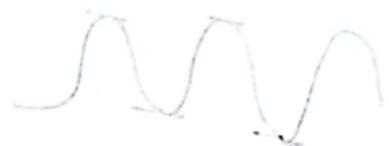
(10 ડાટા
બિલ્ડ)

* Adv: એક એક લોપ એક લોપ
ના
પુરા લોપ રાખી જોગશીલ

SVM Objective function

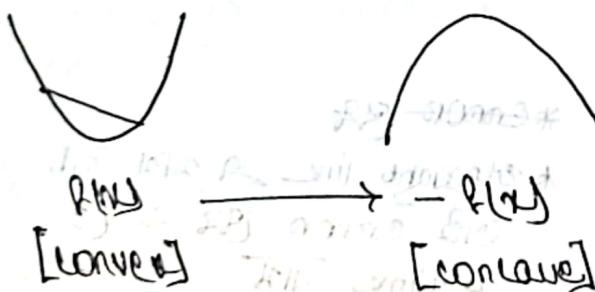
$$\text{Hinge loss} = \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b))$$

$y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \rightarrow \text{correct}$
 \downarrow
 given predicted



$$0 \leq 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)$$

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b))$$

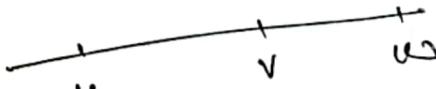
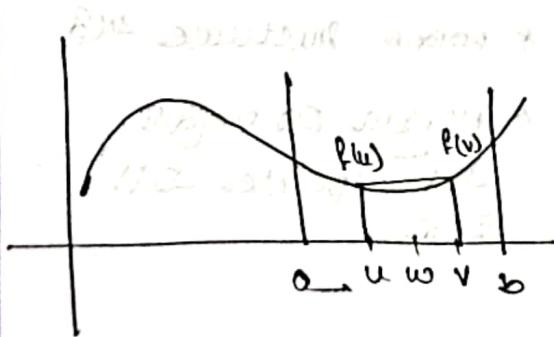


convex \rightarrow minimise
 [Zigzag, concave maximise]

concave \rightarrow

neither convex nor
 concave

* Convex is represented by a domain



$$\lambda u + (1-\lambda)v \rightarrow \text{convex combination}$$

$$\frac{1}{2}u + \frac{1}{2}v = \frac{u+v}{2}$$

$$\lambda f(u) \rightarrow \lambda f(u) + (1-\lambda)f(v)$$

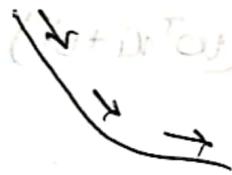
$$f(\underbrace{\lambda u + (1-\lambda)v}_w) \leq \lambda f(u) + (1-\lambda)f(v) \rightarrow \text{convex}$$

* Convex হলো gradient র প্রবলেম ২/৩।



\rightarrow gradient descent changes

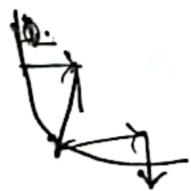
* Batch GD \rightarrow



* ৫১৮ data points
use - সকল মাত্র, সু

error কর এবং time
মাত্র

* step size \rightarrow



* error রেখা

* straight line এ মাত্র না,
তাই corner কর এবং প্রস্তুত
to time মাত্র

* প্রেস্ট পয়েন্ট কর GD

(কর এবং তারি স্ট্রেইট
লাইন এ মাত্র না কর)

* error রেখা fluctuate কর

* weight এর অপ্রকা

বাব update করা,
মাত্র

বেবে

$$w_{t+1} = w_t - \eta(\lambda - 1) + \eta X^T y$$

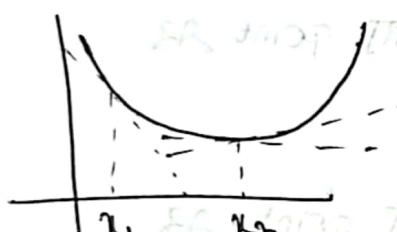
$$w_{t+1} = w_t + \eta X^T y - \eta \lambda X^T X$$

* An unique tangent lies at the point of a function.

தொடர்ச்சி என்றால் கீழே கொண்டு முடியும்

* tangent curve touch $y = f(x)$

* sub-tangent curve touch $y = f(x) - \frac{f'(x_0)}{2}(x - x_0)$



* மூலக் காலத்தில்

நீண்ட காலம்,

எனவே போன்று,

நீண்ட காலம்,

எனவே போன்று,

எனவே போன்று,

எனவே போன்று,

$$\frac{\partial J}{\partial w_j} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) x_{ij}$$

Sub-gradient descent

$$J(w) = \frac{1}{2} \|w\|^2 + C \cdot \max(0, 1 - y_i w_i)$$

$$\frac{\partial J}{\partial w_j} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) x_{ij} + C \cdot \max \begin{cases} 0 & \text{if } y_i w_i \geq 1 \\ -y_i x_{ij} & \text{else} \end{cases}$$

$$\begin{aligned} & \frac{\partial J}{\partial w_j} = w_j^T w \\ & = w_{j1}^2 + w_{j2}^2 + \dots \\ & + w_{jn}^2 \end{aligned}$$

$$\begin{aligned} & \frac{\partial J}{\partial w_j} = w_j^T w \\ & = w_j^2 + w_{j+1}^2 + \dots \\ & + w_{j+n}^2 \end{aligned}$$

* learning rate decreases over time, due to slow progress towards minima

* sum of different learning rates use α_t for each iteration

* $\alpha_t = \frac{1}{t}$

* learning rate plot: two point (0,0) and point A
jump to $\frac{1}{t}$

* learning rate plot: two point (0,0) and point A
jump to $\frac{1}{t^2}$

* time varying learning rate solution

$$w_j^{(t+1)} = w_j^t - \gamma_t \frac{\partial J}{\partial w_j}$$

$$\left| \frac{\partial J}{\partial w_j} = w_j + CN - y_i x_i \right|$$

$$w_j^{(t+1)} = w_j^t - \gamma_t w_j = w_j (1 - \gamma_t) \rightarrow \frac{\partial J}{\partial w_j} = 0$$

$$w_j^{(t+1)} = w_j^t - \gamma_t (w_j + CN - y_i x_i) \rightarrow \frac{\partial J}{\partial w_j} = -y_i x_i$$

$$w_j^{(t+1)} = w_j^t - \gamma_t (w_j) + CN \gamma_t y_i x_i$$

$$w_j^{(t+1)} = (1 - \gamma_t) w_j + \gamma_t \left[CN y_i x_i \right]$$

$$\gamma_t = \frac{40}{1+t}$$

start at 0.01

end at 0.4

ratio 4:1

SUM

\rightarrow Null or Least

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_i \max(0, 1 - y_i(w^T x_i + b))$$

regularization term

- maximize the margin

hyper-parameter
• controls trade off bet' large margin & small hinge loss

FORMAL NOTATION

$$w_0 \leftarrow w$$

$$w \leftarrow [w_0, b]$$

$$q_i \leftarrow [x_i, 1]$$

Empirical Loss

- Minimizes empirical loss
- Penalizes weight vectors that make mistakes

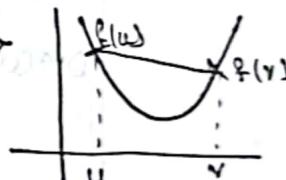
$$(w^T x_i + b) = \hat{y}_i$$

Empirical

* This function is convex in w, b .

* A function is convex when the domain is:-

$$f(\lambda u + (1-\lambda)v) \leq \lambda f(u) + (1-\lambda)f(v) \quad \forall \lambda \in [0,1]$$



Necessary condition:-

- ① x needs to be minimum
- ② $\nabla f(x) = 0$

* Tangent lies below a convex function with minimum

Convex functions

$\nabla f(0) = 0$ definition

* HOW to plot convex function:

① Definition of convexity

② 2nd derivative is +ve \rightarrow Upward shift of the function

③ 2nd derivative is +ve \rightarrow Convex function

semi definite \rightarrow vector functions

$f(x) = -x$	$f(x) = x^2$	$f(x) = \max(0, x)$
Convex	Convex	Convex

$\nabla^2 f(0) = 2I$ \rightarrow positive definite

$\nabla^2 f(0) = 0$

GRADIENT DESCENT

[To solve sum optimization problem]

* Start with an initial guess, w_0

* Compute gradient of $J(w)$ at w^t .

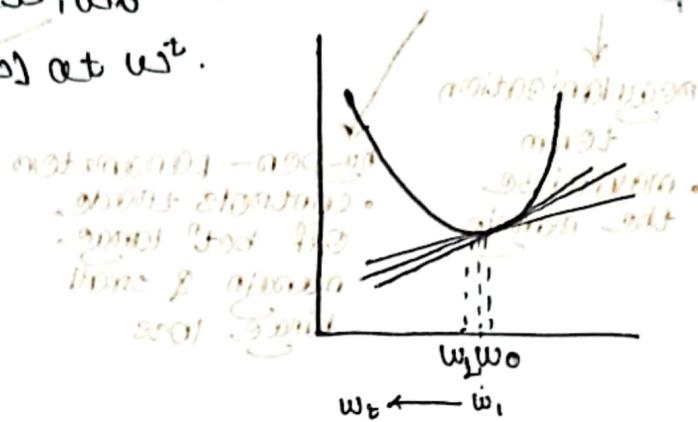
$$\rightarrow \nabla J(w^t)$$

* Update w^t to w^{t+1}

by taking a step
in the opposite
direction.

$$w^{t+1} = w^t - \eta \nabla J(w^t)$$

+ learning rate



Disadv: Gradient of sum requires summing of entire training set \rightarrow SLOW

STOCHASTIC GRADIENT DESCENT [preferable for huge data]

* A training data is given as $\{(x_i, y_i)\}$ & learning rate

* Initialize $w^0 = 0 \in \mathbb{R}^n$

* for epoch = 1, ..., T:

① Pick a random example (x_i, y_i) from S

② Repeat (x_i, y_i) to make a full dataset &

take derivative of sum objective at the whole set

③ Set w^{t+1} to be $\nabla J^t(w^{t+1})$

$$J^t(w) = \frac{1}{2} w_0^T w + C \cdot N \max(0, 1 - y_i w^T x_i)$$

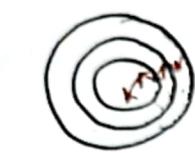
\hookrightarrow Hinge Loss

④ Update: $w^t \leftarrow w^{t+1} - \gamma_t \nabla J^t(w^{t+1})$

\hookrightarrow learning rate

⑤ Return w

Adv: Can converge to minimum of J if γ_t is too small.



stochastic GD



GD

Stochastic GD has many more updates than GD which are less expensive.

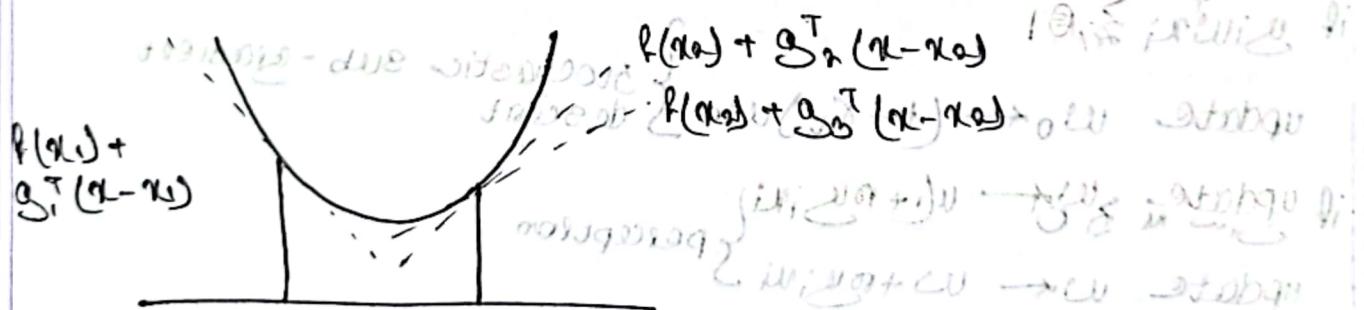
SUB GRADIENTS



A sub-tangent at a point is any line lies below the function at the point.

A sub-gradient is the slope of that line.

$$f(y) \geq f(x_0) + g_i^T(y - x_0) \text{ for all } y$$



* f is differentiable at x_0 , so there's a tangent at x_0 . $\therefore g_i$ is the gradient of f at x_0 (and a subgradient).

* g_2 and $w_2 g_3$ are subgradients

$$J^t(w) = \frac{1}{2} \|w\|_2^2 + C \cdot \max(0, 1 - y_i w^T x_i)$$

* solve max $\nabla J^t = \begin{cases} [w_0; 0] & \text{if } \max(0, 1 - y_i w^T x_i) = 0 \\ [w_0; 0] - C \cdot N y_i x_i & \text{otherwise} \end{cases}$

* compute gradient

STOCHASTIC SUB-GRADIENT DESCENT FOR SVM

Given a training set $S = \{(x_i, y_i)\}_{i=1}^n, x \in \mathbb{R}^d, y \in \{-1, 1\}$

① Initialize $w^0 = 0 \in \mathbb{R}^d$

② For epoch = 1 to T : need to shuffle

③ For each training example $(x_i, y_i) \in S$ for each epoch: minimizing cost function

if $y_i w^t \leq 1 \rightarrow$ learning rate

$$w \leftarrow (1 - \gamma_t)[w^t; 0] + \gamma_t C y_i x_i$$

else

$$w \leftarrow (1 - \gamma_t)w^t$$

④ Return w

PERCEPTRON V/S STOCHASTIC SUB-GRADIENT DESCENT

if $y_i w^t \leq 1$

update $w \leftarrow (1 + \gamma_t)w^t$ stochastic sub-gradient descent

if $y_i w^t \geq 0$

update $w \leftarrow w + \gamma_t y_i x_i$ perceptron

$$\text{Lperceptron}(y, w^t x) = \max(0, -y w^t x) \rightarrow \text{No regularization}$$

$$\text{L hinge}(y, w^t x) = \max(0, 1 - y w^t x) \rightarrow \text{Regularization}$$

(SVM optimizes hinge loss)

* Perceptron and stochastic has similar framework with different objectives.

* Perceptron cannot maximize margin width. Can have a margin.

CONVERGENCE AND LEARNING RATES

$$\delta_t = \frac{\delta_0}{1 + \frac{\delta_0 t}{C}} \Rightarrow \delta_t = \frac{\delta_0}{1+t}$$

↓
learning rate
 $t \downarrow, \delta_t \uparrow$

the GTR has many ways

AD = Gradient Descent Method
GD with step size δ_0 / C is called
SGD or Stochastic Gradient Descent



Quiz 2 Solution

(1, +), (2, +), (3,

$(1, +5), (2, +5), (3, +5), (6, -5), (7, -5), (9, -5), (11, +5)$

$$a=0.6, b=1$$

$$\begin{array}{r}
 & + & + & + & - & - & - & + \\
 \hline
 & 1 & 1 & 3 & 6 & 7 & 9 & 11 \\
 & 1 & 2 & 3 & 6 & 7 & 9 & 11 \\
 \hline
 + & - & - & - & - & - & - & - \\
 \hline
 & 1 & - & - & - & - & - & -
 \end{array}$$

CONCORDE = A fast supersonic jet
airplane.

$$a = \frac{b+b}{2} = b, b, b = 1$$

$$\begin{array}{r}
 \text{P} \quad \text{y} \quad \text{a} \quad \text{b} \quad \text{c} \\
 \hline
 + & + & + & - & - & + \\
 \hline
 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\
 \hline
 + & + & + & + & - & - & - \\
 \hline
 \end{array}$$

$\text{C}_{\text{min}} = 1 \text{ mg/mL}$

possible values of $a = 8$

minimum error = 1

Value of $a^2 B B, b^2 l$ for minimum error

$$a = \frac{9+11}{2} = 10, b = 1$$

$$+ + + - - - +$$

$$+ + + + + + -$$

$$0.0000 = 9$$

$$a = 11, b_1 = 1$$

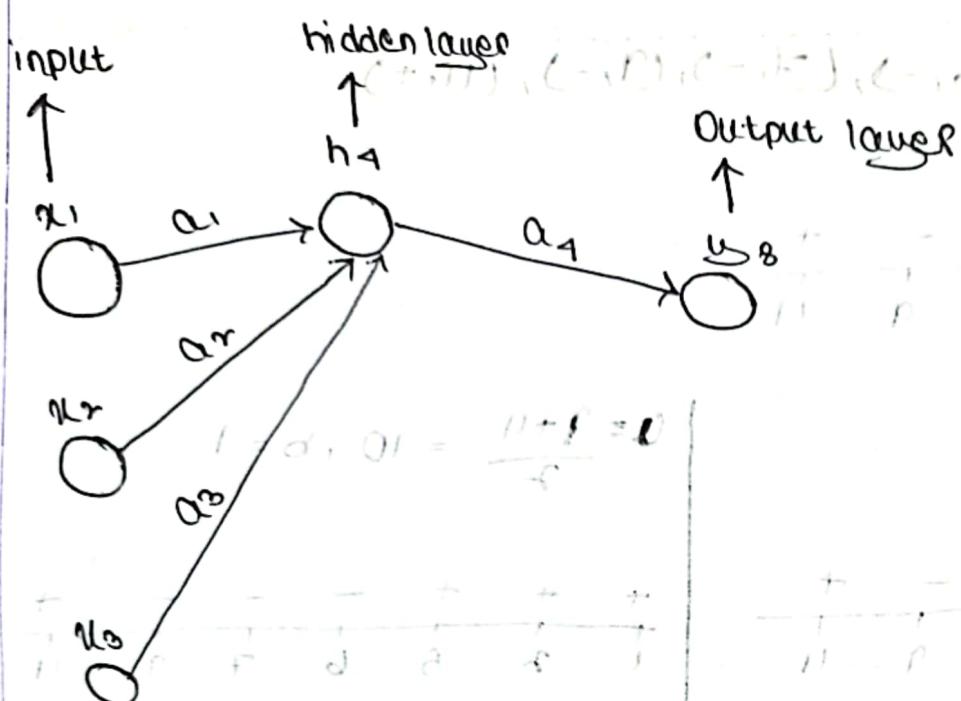
$$\overline{\mathcal{E}^0 \mathcal{G} + f} = (4) \mathcal{E}^0 \mathcal{G}^0$$

தமிழ் மூலத்தில் காலாக விடு
கோடி = ८

1) $\text{f}(x) = x$ $\text{f}(g(x)) = g(x)$

minimum error

NEURAL NETWORK



4th node of hidden activation a_4

$\text{in}_4 = w_{1,4}x_1 + w_{2,4}x_2 + w_{3,4}x_3 + b_4 \rightarrow$ for hidden layer

$a_4 = \frac{1 + e^{-\text{in}_4}}{1 + e^{\text{in}_4}}$

$a_4 = g(\text{in}_4) = \frac{1}{1 + e^{-\text{in}_4}}$

$a_4 = g(\text{in}_4) = \frac{1}{1 + e^{-\text{in}_4}}$

(logistic function)

from complex activation write $g(\text{in}_4)$

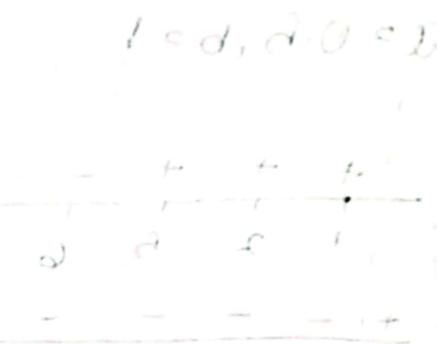
$\text{in}_4 = w_{4,8}x_4 \rightarrow$ for output layer

$$\frac{1}{1 + e^{-\text{in}_4}}$$

for minimum loss $f = a_4 - y_8$

* Hidden layer
যোগাত ও যোগ
না ও যোগ

* a_1, a_2, a_3, a_4
→ Activation

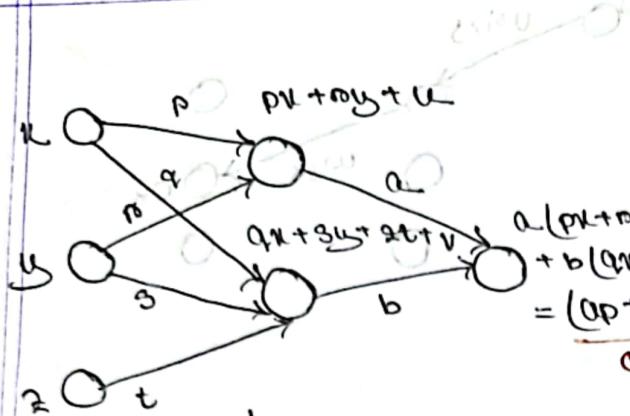


* Input layer (প্রবেশ)
Hidden layer (যোগাত)
গোটা মুক্তি a_1
 a_2 value in A2
value a_1, a_2 \neq 0
 $a_3 = x_3$ কিন্তু
identity function
[f(x) = x] একই মান

* Hidden layer (যোগ)
Output layer
গোটা মুক্তি a_4
 $\frac{1}{1 + e^{-\text{in}_4}}$ [logistic regression]
মান

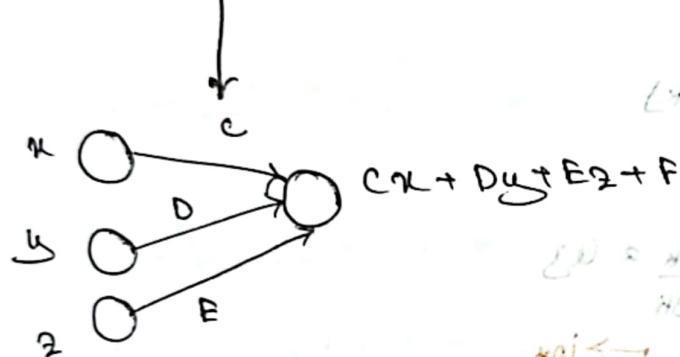
TUESDAY

DATE: 22/08/23



$$\begin{aligned} \text{For } a_1 &= p(x) + q(y) + r(z) + u \\ \text{For } a_2 &= s(x) + t(y) + b(z) + v \end{aligned}$$

$$\begin{aligned} a_1(p(x) + q(y) + r(z) + u) \\ + b(s(x) + t(y) + b(z) + v) + v \\ = \frac{(ap+aq)x}{c} + \frac{(ar+as)y}{d} + \frac{(ar+ab)z}{e} + \frac{(au+bv+v)}{f} \end{aligned}$$



$$ED = \frac{4916}{4800}$$

$$\frac{1}{4800} = \frac{4916}{4800}$$

$$\frac{4916}{4800}, \frac{4916}{4800}$$

For activation, we use :-

- ① logistic function for binary classification
- ② softmax for multinomial
- ③ identity function for linear
- ④ threshold function for perceptron

* Neural network is more useful for multiple layers inputs.

* Identity function; input π , output π



* Identity function with other functions non-linearity shows π .

$$\begin{aligned} \text{For } \pi(\pi(\pi(\pi(\pi(\pi(4916 - 4800) \cdot 1) + 4916 - 4800) \cdot 1) + 4916 - 4800) \cdot 1) + 4916 - 4800 \\ = 4916 - 4800 = 116 \end{aligned}$$

CS 7.007

BACK PROPAGATION

\rightarrow गुणी त्रैतीय

$$E = \sum (t_k - o_k)^2 \rightarrow \text{माना predict } o_k$$

$$in_k = \sum w_{jk} w_{jk} (o_j) o_k \circ g(\text{link})$$

$$w_{jk} = w_{jk} - \lambda \frac{\partial E}{\partial w_{jk}}$$

$$\frac{\partial E}{\partial w_{jk}} = 2(t_k - o_k) \rightarrow \frac{\partial}{\partial w_{jk}} (-o_k)$$

$$\frac{\partial o_k}{\partial w_{jk}} \circ \frac{\partial g(\text{link})}{\partial w_{jk}}$$

$$= \frac{\partial in_k}{\partial w_{jk}} \cdot \frac{\partial g(\text{link})}{\partial in_k}$$

$$= o_j (1 - o_j) g(\text{link})$$

$$\frac{\partial in_k}{\partial w_{jk}} = o_j$$

$$\frac{\partial g(\text{link})}{\partial in_k} = \frac{\partial}{\partial x} \frac{1}{1 + e^{-x}}$$

$$= \frac{\partial(1 + e^{-x})}{\partial x} \frac{\partial}{\partial in_k} = \frac{e^{-x}}{(1 + e^{-x})^2}$$

$$= \frac{e^{-in_k}}{(1 + e^{-in_k})^2}$$

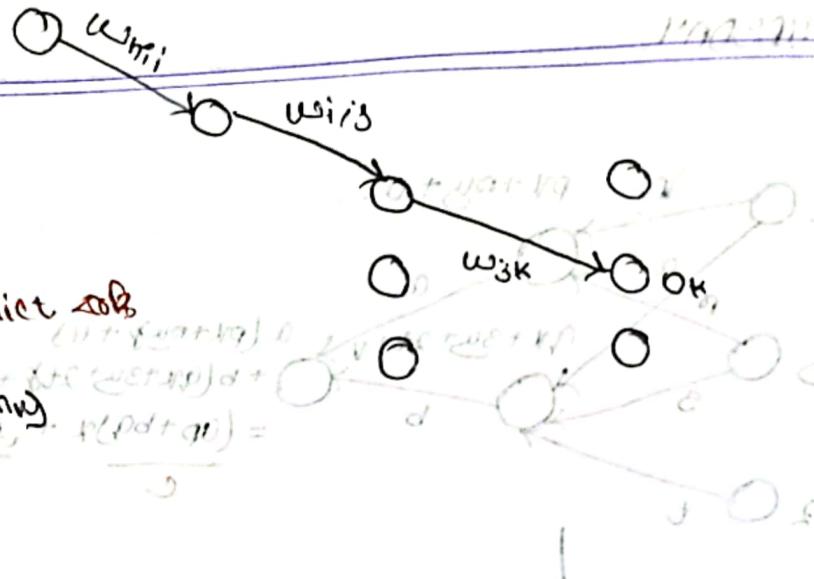
$$\frac{\partial E}{\partial w_{jk}} = o_j (1 - o_j) g(\text{link})$$

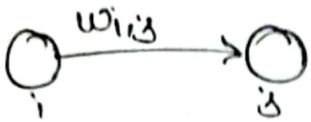
$$w_{jk} = w_{jk} - \lambda \frac{\partial E}{\partial w_{jk}}$$

$$w_{jk} = w_{jk} - \lambda (2(t_k - o_k)) o_j (1 - o_j) g(\text{link})$$

$$g(\text{link}) = (1 - o_j) g(\text{link}), \Delta k = \frac{\partial E}{\partial w_{jk}} (t_k - o_k) o_j g(\text{link})$$

$$w_{jk} = w_{jk} - \Delta k o_j$$





$$E = \sum (t_k - o_k)^2$$

$$\frac{\partial E}{\partial w_{ij}} = -2 \sum (t_k - o_k) \frac{\partial o_k}{\partial w_{ij}} = -2 \sum (t_k - o_k) \left(\frac{\partial o_k}{\partial w_{ij}} \delta_{ij} \right)$$

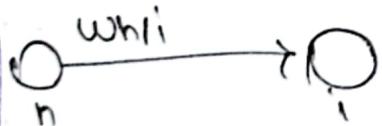
$$\frac{\partial E}{\partial w_{ij}} = -2 \sum \Delta K w_{ij} \delta_{ij} \frac{\partial \delta_{ij}}{\partial w_{ij}} = -2 \sum \Delta K w_{ij} \delta_{ij} \frac{\partial w_{ij}}{\partial w_{ij}} \delta_{ij}$$

$$= -2 \sum \Delta K w_{ij} \delta_{ij} \delta_{ij}$$

$$\frac{\partial E}{\partial w_{ij}} = -2 \alpha_i \delta_{ij} \sum \Delta K w_{ij} \delta_{ij}$$

$$\Delta_{ij} = \delta_{ij} \sum \Delta K w_{ij} \delta_{ij}$$

$$w_{ij} = w_{ij} + \lambda \alpha_i \Delta_{ij}$$



$$\frac{\partial E}{\partial w_{nii}} = -2 \alpha_n \delta_{ni} \sum \Delta K w_{nii} \delta_{ni}$$

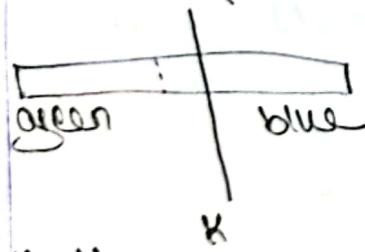
$$\Delta_i = \delta_{ni} \sum \Delta K w_{nii} \delta_{ni}$$

$$w_{nii} = w_{nii} + \lambda \alpha_n \Delta_i$$

K-NEAREST NEIGHBOURS (KNN) on 21/10/2012

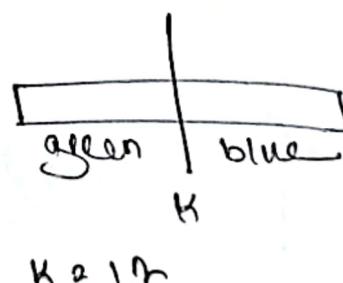
- * K as value. 06 + 38 = 00 + 36
54 + 1 = 52 + 1
- * training data set & testing data 21/10/2012
- * loop 21/10/2012 from 1 to no. of rows in 21/10/2012
training data set to find out the distance 1 + 3, 4 + 5
between training data and test data.
- * Then distance परिमाण से असर नहीं soft असर 2/2
- * K value यहाँ 'distance' जल्दी पर्त हो (e.g. K=5)
- * If 3 blue & 2 green then vote for 3
under 1 test data यहाँ, 3, for
- * for example, 3 for +ve/yes/1 and 2 for -ve/no/0
उत्तर, so test data is +ve.

for example:



$K=1$

- * 7 blue points
- * 6 green points
- So, test data lies under \rightarrow blue points.



$K=1$

- * 6 blue points
- * 6 green points
- so, either decision
by 51 or
randomly चुनौती
pick 50% for

distance $\Rightarrow d_i = \sqrt{\sum_{j=1}^n (x_{ij} - a_{ij})^2}$

→ training data

→ test data

$\begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix}$

\rightarrow 2D array \Rightarrow 1D array শালাতী

$\rightarrow K$ - ক্লাসিফাই লেবেল মান

\rightarrow প্রত্যেক প্লেট এর ক্ষেত্রে score শালাতী হবে

$s_1 = w_{11}x_1 + w_{12}x_2 + \dots + w_{1m}x_m + b_1$

$s_2 = w_{21}x_1 + w_{22}x_2 + \dots + w_{2m}x_m + b_2$

⋮

s_K

মোটা weight array use 2048 (512 weight same)
ইন্সি মান,

so weight র মতো use 2048

$\rightarrow X$ রয়েছে image ক্লাসিফাই করে flatten করে প্রক্রিয়া, তারপর probability এর মতো লাভাতী।

$$\frac{s_1}{s} \quad \frac{s_2}{s} \quad \dots \quad \frac{s_K}{s} \rightarrow \text{probability}$$

$$\frac{1}{4} \quad 2 \quad p_1 \quad p_2 \quad \dots \quad \dots \quad p_K$$

Prediction

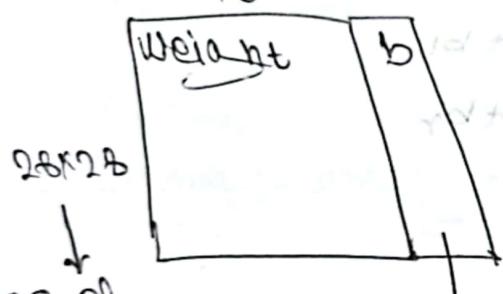
$$Y = \underset{\text{total data}}{\underset{n}{\sum}} \underset{\text{total label}}{\underset{k}{\sum}} y_{ij} \log p_{ij}$$

* weights \Rightarrow update 2020 2B

$$\frac{\partial J}{\partial w_{p,j}} = - \sum_{i=1}^m \sum_{k=1}^K y_{i,k} \text{ (DATA)} \text{ (label of } i \text{)}$$

\downarrow
क्रम
एक्स

\rightarrow no. of class
10



28x28
 \downarrow
no. of
weights

one D array
(10 के लिए)

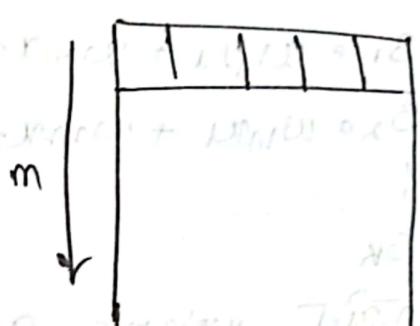
$\frac{\partial J}{\partial w_{1,1}} w_{1,2}, w_{1,3}, \dots$

\downarrow
प्रति वर्ष
में

2000 में 70000 - 4 +

2000 में 70000 - 4 +

$K = 28 \times 28$



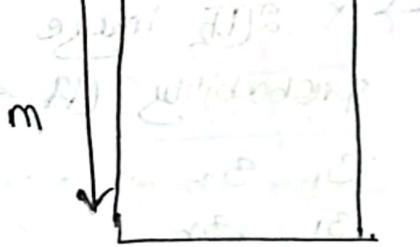
m

One-hot vector

$$W = W - \frac{1}{m} * np.dot$$

आ जा जाना है prediction

$$b = b - \frac{1}{m} * np.sum [pred - \text{onehot}]$$



m

गणना करने की prediction

* X \Rightarrow reshape 2D

* label \Rightarrow one-hot
vector 41 से 2B

पहली रेंज का गुणनक
दूसरी रेंज का गुणनक

पहली रेंज का गुणनक
दूसरी रेंज का गुणनक

Training data: $\{(x_i, y_i)\}$

Test data: x

$$\text{soft}(x, w_i) = \frac{e^{-w_i^T x}}{S}$$

$$\text{soft}(x, w_1) + \frac{e^{-w_2^T x}}{S} + \dots + \frac{e^{-w_m^T x}}{S}$$

Doteno.

x_1	1-12-012
x_2	1-12-012
\vdots	\vdots
x_m	\vdots

$$1 - mB = mC$$

$$\sum_{k=1}^m \text{soft}(x, w_k) y_k \geq 0 \rightarrow +1 \text{ if } \text{soft}(x, w_k) > 0 \text{ else } -1 \text{ if } \text{soft}(x, w_k) < 0$$

* Data \Leftrightarrow centralise \Leftrightarrow ~~min max~~

$$x' = \frac{x - \text{mean}}{\text{std}}$$

* Classifying \Leftrightarrow $\text{soft}(x, w_k) \geq 0$ for all k

$$\text{mean} = \frac{1}{n} \sum x_i$$

$$\text{std} = \sqrt{\frac{1}{n-1} \sum (x_i - \text{mean})^2}$$

$$w = \frac{1}{n} \sum x_i y_i$$

$$w = \frac{1}{n} \sum x_i y_i$$

$$\text{mean} = \frac{1}{n} \sum x_i$$

$$\text{std} = \sqrt{\frac{1}{n-1} \sum (x_i - \text{mean})^2}$$

$$w = \frac{1}{n} \sum x_i y_i$$

$$\alpha_1 = \frac{\alpha_1 - \mu}{\sigma}$$

$$\alpha_2 = \frac{\alpha_2 - \mu}{\sigma}$$

⋮
⋮
⋮
⋮

$$\alpha_m = \frac{\alpha_m - \mu}{\sigma}$$

No. of test

Fit the data $\alpha_1, \alpha_2, \dots, \alpha_m$ with weight $\alpha_1, \alpha_2, \dots, \alpha_m$

$\alpha \rightarrow \text{test}$

$$e^{-\frac{1}{2}(\|\alpha_1 - \alpha\|^2 + \|\alpha_2 - \alpha\|^2 + \dots + \|\alpha_m - \alpha\|^2)}$$

Euclidean distance α is zero weight α is minimum JPT -

$$-\|\alpha_1 - \alpha\|^2 - \|\alpha_2 - \alpha\|^2 - \dots$$

$$e^{-\frac{1}{2}(\|\alpha_1 - \alpha\|^2 + \|\alpha_2 - \alpha\|^2 + \dots + \|\alpha_m - \alpha\|^2)} = \frac{e^{-\frac{1}{2}\|\alpha_1 - \alpha\|^2}}{S}$$

$$\begin{aligned} & -\|\alpha_m - \alpha\|^2 \rightarrow \text{score} \\ & \frac{-\|\alpha_m - \alpha\|^2}{S} \rightarrow \text{softmax} \end{aligned}$$

α_1	
α_2	
\vdots	
α_m	

Training

α_1	
α_2	
\vdots	

Test

$$\frac{-\frac{1}{2} \|x_1 - \alpha\|^2}{S}$$

$$\|x_1 - \alpha\|^2 = \|x_1\|^2 + \|\alpha\|^2 - 2x_1^T \alpha$$

↓ ↓
 centralise 0 നാണ്യാഭാസ
 ഏകി same
 പ്രധാന നേര
 constant
 വാച ലിത്തു പാടി
 $\Rightarrow -2x_1^T \alpha$

$$\frac{-\frac{1}{2} \|x_1 - \alpha\|^2}{S} \Rightarrow \frac{x_1^T \alpha}{S}$$

$$\frac{x_1^T \alpha}{S} \quad \frac{x_2^T \alpha}{S} \dots \frac{x_n^T \alpha}{S}$$

score(1) score(2) ... score(n)

$(\text{score}(1) \times 1) + (\text{score}(2) \times -1) + \dots + (\text{score}(n) \times -1) \rightarrow$ എഴി + ve
 അഭാസ
 survived,
 -ve അഭാസ
 not survived.

vector-normalization