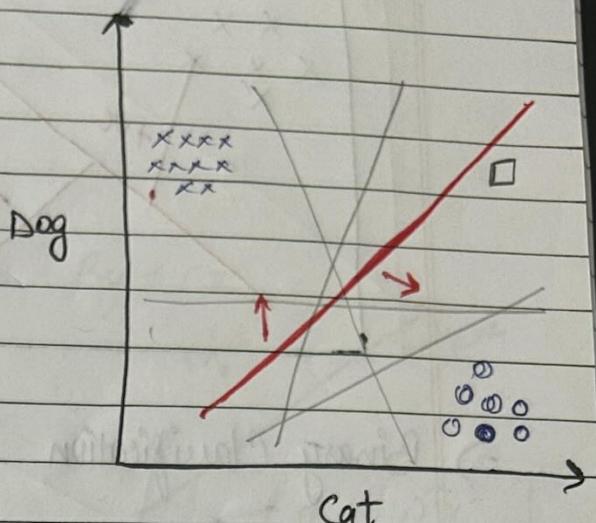
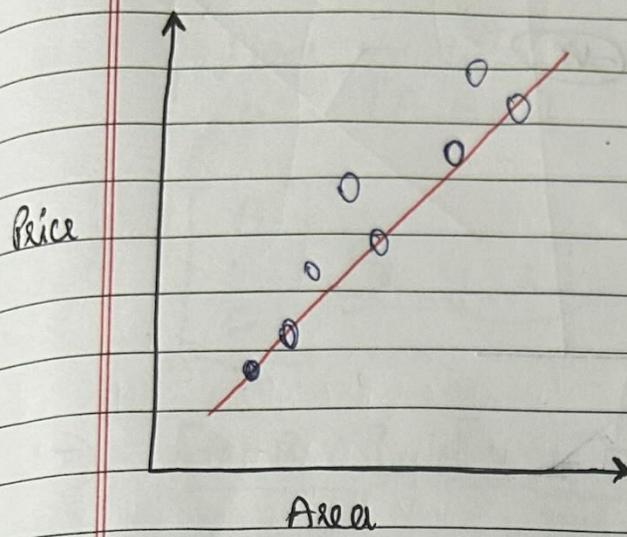


18/1/23

Date / /
Page No.

LOGISTIC REGRESSION

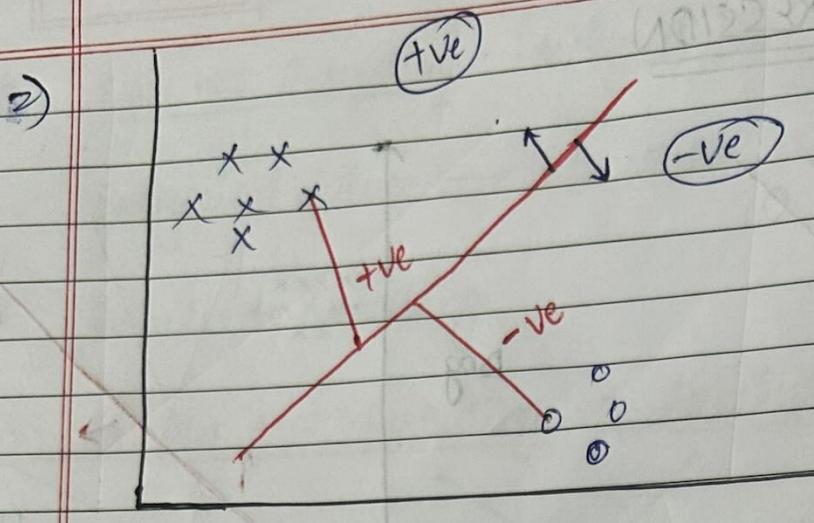


→ Supervised Data → Regression
 → Classification

Topics -

- Logistic Regression
- Loss function for logistic regression
- Why logistic regression called regression
- Precision | Recall | F1-score | Confusion Matrix
- Sigmoid Function | Log function | Step function | Exponential function
- Ridge | Lasso
- Bias-variance trade off
- Why we can't use linear ~~regression~~ instead of logistic regression?
- Overfitting | Underfitting

* In logistic we have to find which line perfectly separates the classification data.



→ Binary Classification

→ Logistic regression is a binary classification algorithm which can separate the data for two diff. classes with a single line.

→ It depends on data, which separation needs to be done.

→ Ques How to identify the +ve class with distance?

→ 2D → line to divide

3D → (x, y, z) → 3 dimension | Plane

↑
to divide data

more than 3D → Hyperplane
(to divide data)

→ $3x - y = 0 \rightarrow$ line

$Ax + By + C = 0 \rightarrow$ line eqⁿ

$Ax + By + Cz + D = 0 \rightarrow$ plane

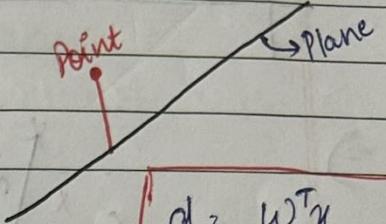
* Intuition of plane & distance using vectors.

$$\text{Plane} \rightarrow w^T u$$

$$w = [A, B, C]$$

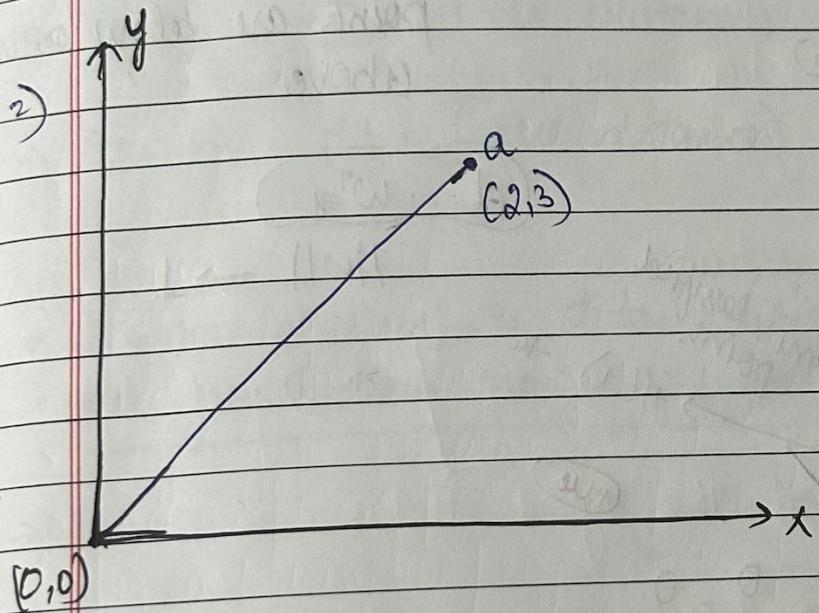
$$\begin{bmatrix} A \\ B \\ C \end{bmatrix} \cdot \underbrace{\begin{bmatrix} u & y & z \end{bmatrix}}_u = Ax + By + Cz$$

$$\Rightarrow \underbrace{[A \ B \ C]}_w \underbrace{\begin{bmatrix} u & y & z \end{bmatrix}}_u$$



$$d = \frac{w^T u}{\|w\|}$$

distance of a point
from a plane



$$\|a\| = \sqrt{2^2 + 3^2} = \sqrt{13}$$

~~19/01/23~~

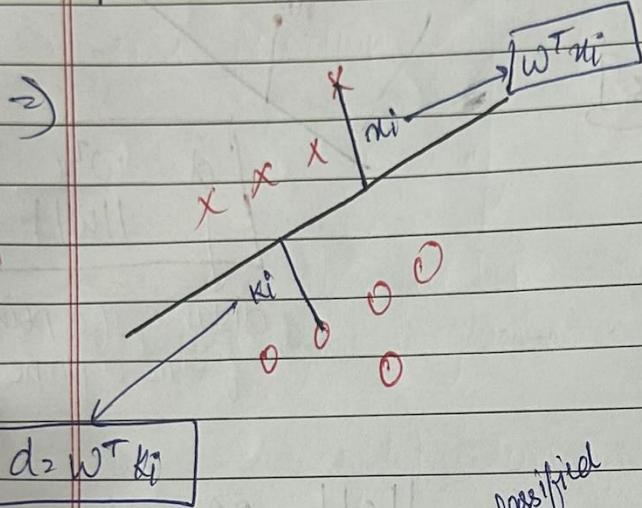
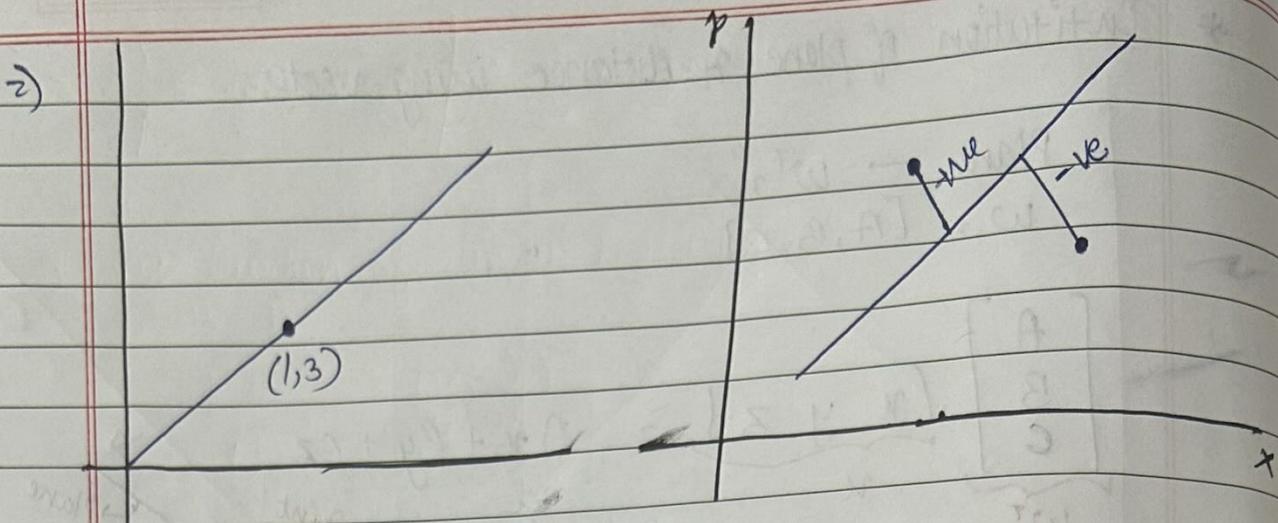
$$\Rightarrow \frac{3u - y}{\|w\|} = 0$$

$$d = \frac{w^T u}{\|w\|}$$

$w = [A, B]$
coefficients / weights / vector

$$\|w\| = \sqrt{3^2 + 1^2} = \sqrt{10}$$

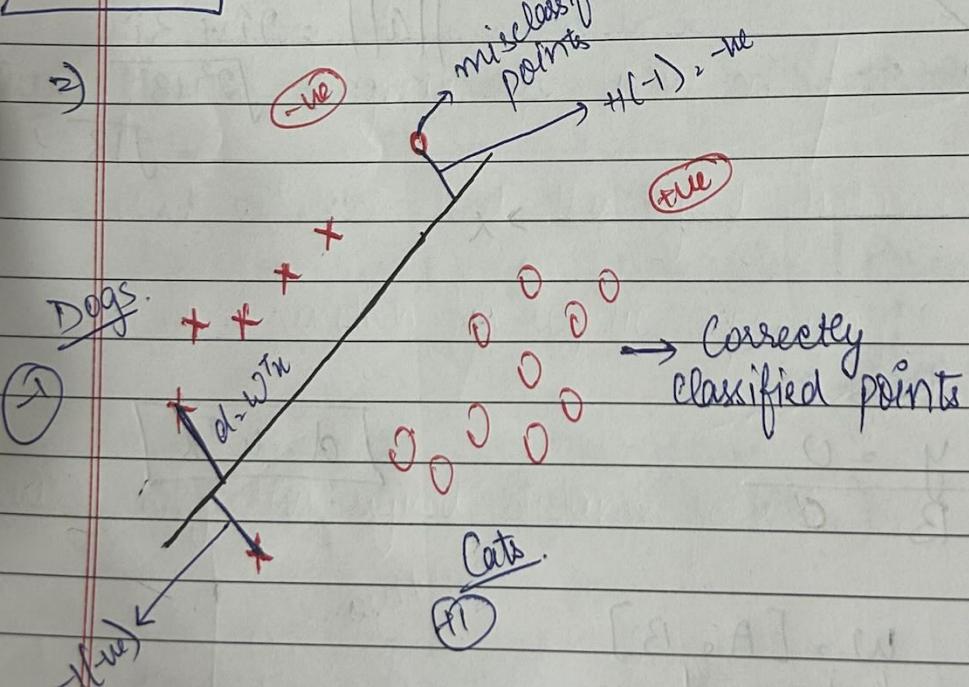
(absolute value of vector)



Distance will be true or
ve based on whether
points are below or
above.

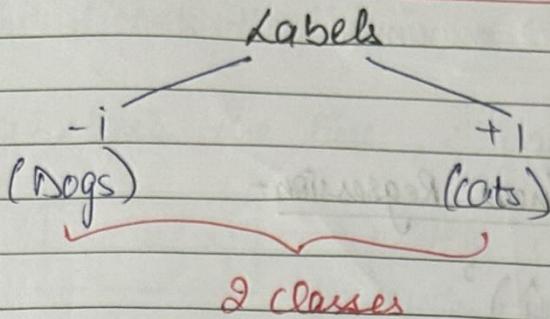
$$d = w^T x$$

$$\|w\| \rightarrow 1$$



Logistic regression is binary classification.

→



→

$$y_i \cdot w^T x \rightarrow \text{Signed distance}$$

↓
Label Distance

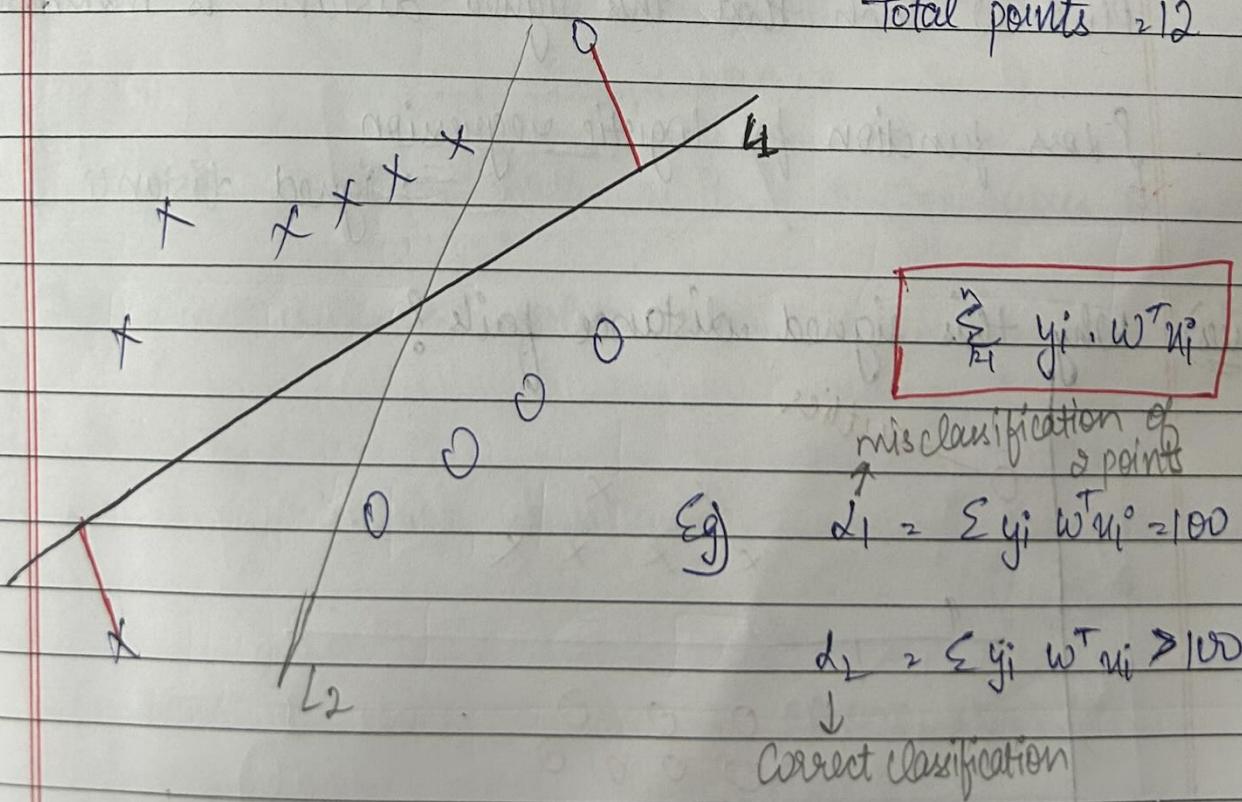
→ Dogs $\rightarrow -1$. (-ve distance) = +ve
label

→ Cats $\rightarrow +1$. (+ve distance) = +ve
label

* Product of label y_i & distance (signed distance) should be greater than zero for correct classification

→

Total points = 12



⇒ $\sum y_i \cdot w^T x_i \rightarrow \text{maximise}$

⇒ Loss function of Linear Regression -

$$\frac{1}{N} \sum_{i=1}^n (y - \hat{y})^2$$

⇒ Loss function for Logistic Regression -

$$\sum_{i=1}^n y_i w^T x_i \rightarrow \text{Maximise}$$

Signed loss function
for logistic regression

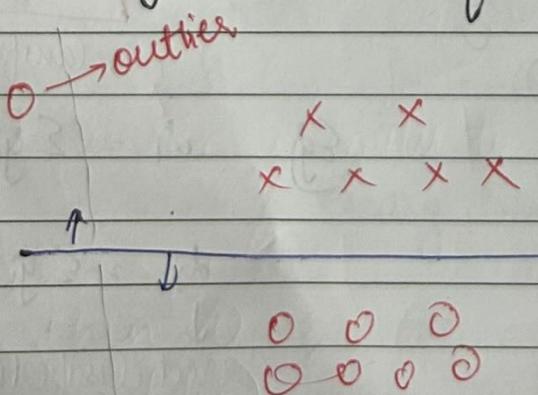
⇒ $w^T \rightarrow$ weights/vectors

$w^* = \text{argmax } \sum_{i=1}^n y_i w^T x_i$ → signed distance

Find w^T such that this signed distance is maximised.

{ Loss function for logistic regression
→ signed distance }

Ques) Why this signed distance fails?



→ Worst line

→ Outlier distorts the line

→ Outlier effects the line, \therefore affects the distance

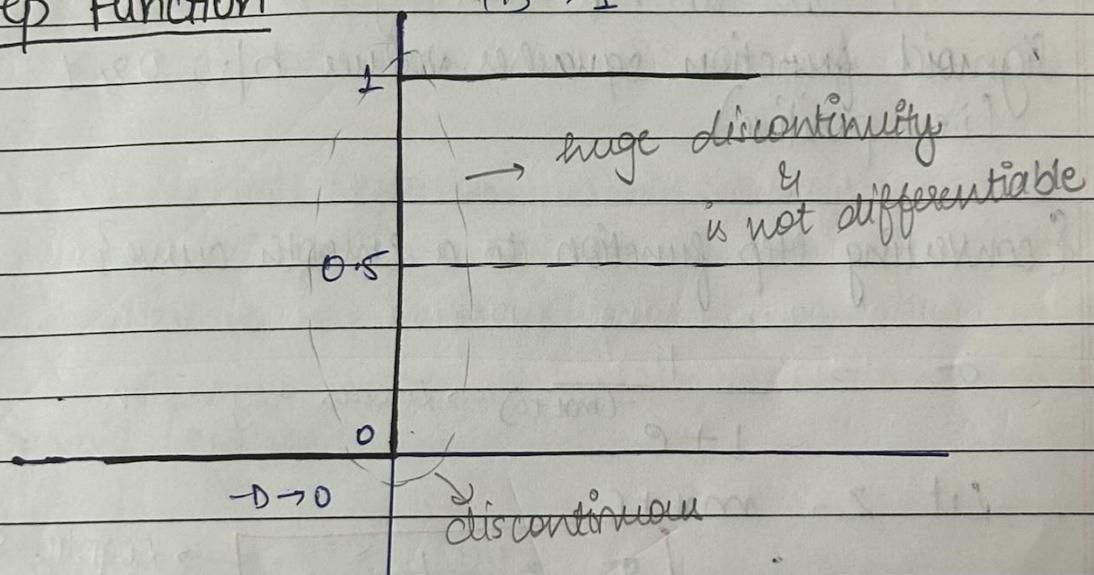
$$\boxed{y_i w^T \gamma_i}$$

\rightarrow outlier can lead the distance towards one or can minimise the distance

→ Outlier impact signed distance, \therefore this is the drawback of signed distance

→ Solution to signed distance \rightarrow STEP FUNCTION

Step function -



\therefore Step function fails as it is discontinuous on y-axis

Step function bound values b/w 0 & 1 $[0, 1]$

Step output $\rightarrow 1$

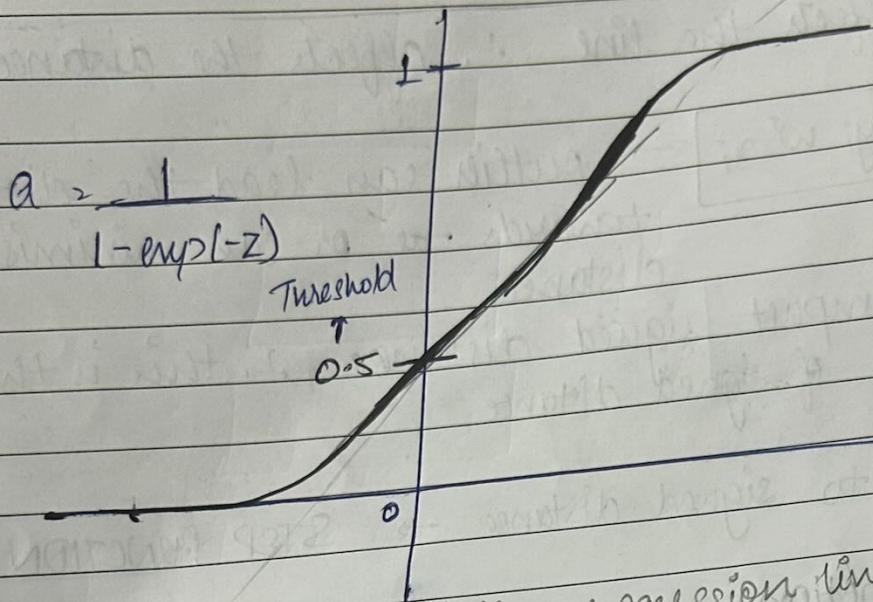
$\hookrightarrow 0$

$$[-\infty, \infty]$$

i. if $f > 0$
0 if $f < 0$

→ Solution to step function \rightarrow Sigmoid Function

SIGMOID FUNCTION



→ kind of linear regression line
 Sigmoid function squashes values b/w 0 & 1
 (i.e. positive)

Converting step function to a smooth curve?

$$\sigma \rightarrow \frac{1}{1 + e^{-(mx+c)}} \rightarrow \text{linear regression}$$

$$\text{Let } z = mx + c$$

$$\therefore \sigma = \frac{1}{1 + e^{-z}}$$

Logistic sigmoid function

Exponential -ve values to the σ data hai
 (Why we did exponential?)

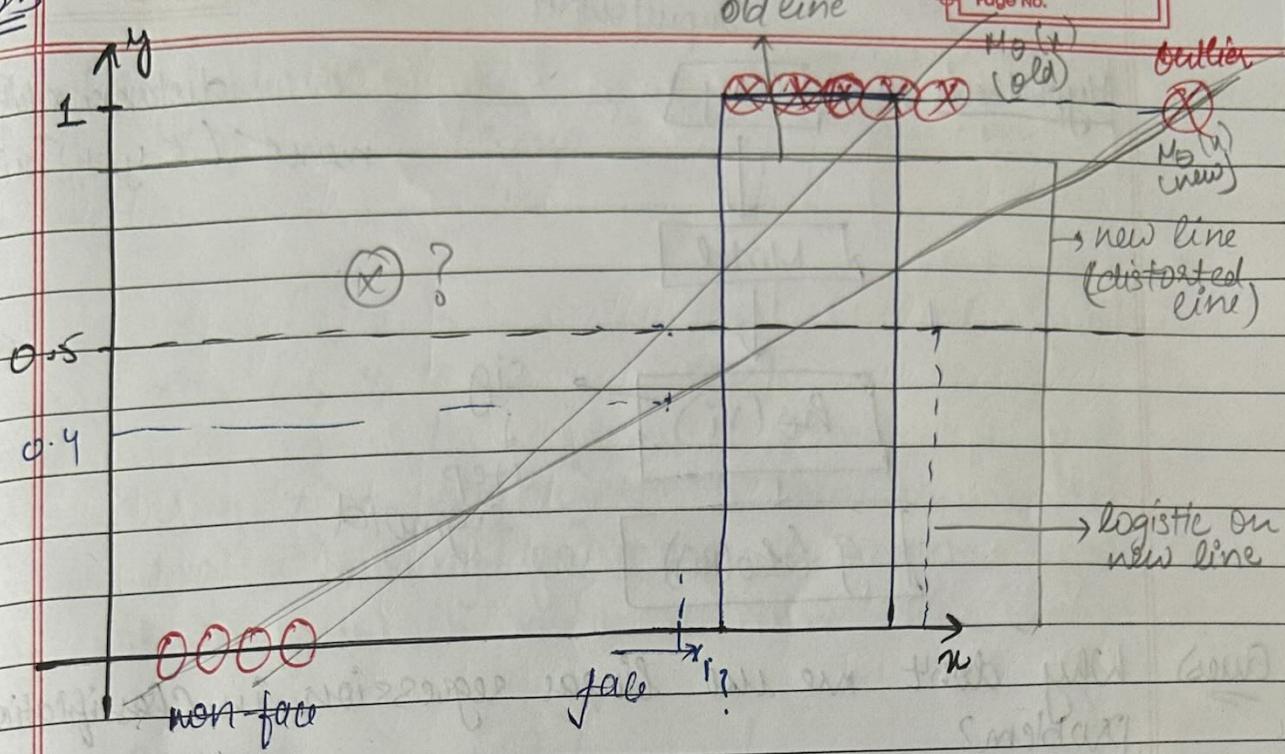
Sigmoid function is similar to linear regression

If $w^T < 0.5$ then 0 (dogs)
 Then 1 (cats)

$$w^T > 0.5$$

20/1/23

Date / /
Page No.



→ If $\geq 0.5 \rightarrow \text{face}$
 $< 0.5 \rightarrow \text{non-face}$

→ Logistic Regression
Geometric Probabilistic

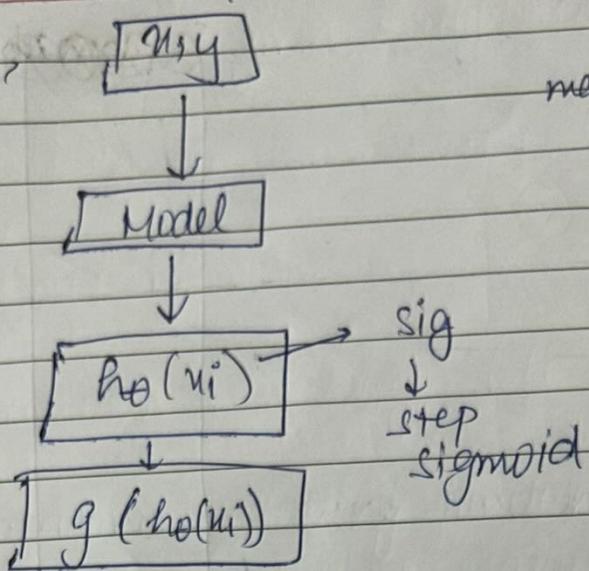
→ Sigmoid (signed)
values, 0.7 | 0.6 | 0.1

→ Step function, bound in 0, 1
Sigmoid function, bound b/w 0 & 1

→ Sigmoid $\rightarrow \frac{1}{1+e^{-z}}$ $\{z, w^T x\}$

Input | Output

2) Hypothesis,



$$\text{margin} \leftarrow (\epsilon - y_i w^T u_i)$$

Ques) Why don't we use linear regression in Classification problem?

→ line distort hui (because of outlier), toh threshold value change hogi

2) Linear regression issliye use nahi hota, kyuki, outlier kei wajah se line distort ho jati hai, jis wajah se threshold value change ho jati hai & woh wrong classification kar deتا hai.

23/1/23

$$\textcircled{1} \text{ Argmax } \sum_{i=1}^n y_i w^T u_i$$

(W)

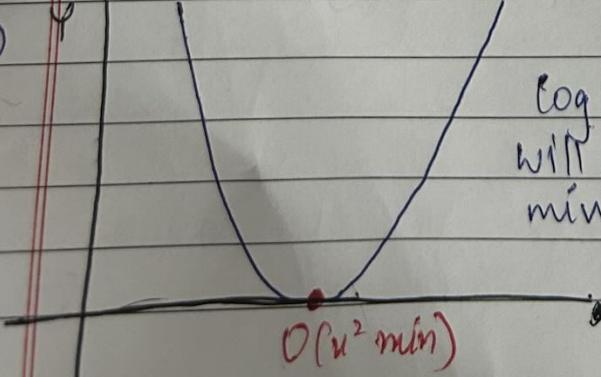
↓
find this w

$$\textcircled{2} \quad f(n) = n^2$$

Transforming to

$$g(n) = \log(f(n)) \\ = \log(n^2)$$

(3)



$\log(n^2)$
will also be
minimum

$$\textcircled{4} \quad \text{argmax } \sum \frac{1}{1 + e^{-(y_i w^T u_i)}}$$

$\log(n^2)$
min value ↑

(i) argmax Σ $\frac{1}{1 + \exp^{-y_i w^T x_i}}$ \rightarrow monotonic increasing function
 (ii)

2) $x_1 = [1, 2, 3, \dots, 100]$
 $x^2 = [1, 4, 9, \dots, 100]$

$$\log(u) = \log(10) \rightarrow \text{maximum}$$

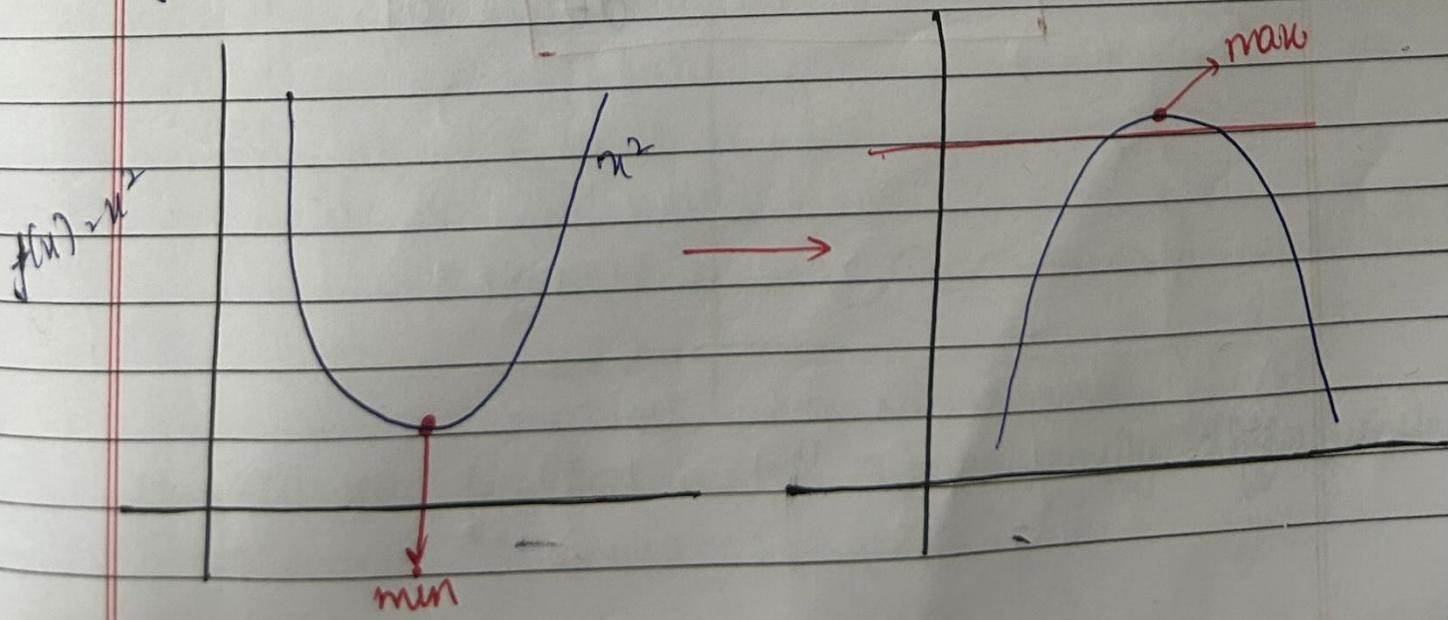
$$\log(u^2) = \log(100) \rightarrow \text{maximum}$$

→ argmax Σ $\log \left(\frac{1}{1 + \exp(-y_i w^T x_i)} \right)$

2) argmax $\Sigma \log \left(\frac{1}{1 + \exp(-y_i w^T x_i)} \right)$
 maximise

$$\left\{ \begin{array}{l} \log u^n = n \log u \\ \log \left(\frac{1}{n} \right) = \log n^{-1} = -\log n \end{array} \right\}$$

3) $f(u) \rightarrow \text{max}$



→ Loss function = $\operatorname{argmin}_w \sum_{i=1}^N \log(1 + \exp(-y_i w^T x_i))$

* $\left\{ \begin{array}{l} w^2 = \text{max} \\ \log(w^2) \rightarrow \text{max} \end{array} \right\}$

→ log optimisation mein easy sehta hai, issiliye log lete hai

log → monotonic increasing function

→ optimisation mein monotonic increasing function use hota hai

→ log ki values 0-1 ke beech mein $-\infty$ hoti hai & fir increase kerti hai

→ exponential \rightarrow monotonic increasing function

$(-\infty \text{ to } 0) \rightarrow y = 0$
Then starts increasing

→ Loss function ko minimise kerte hai using gradient descent
then,

$$\theta = \theta - n \cdot \frac{dL}{d\theta}$$

Probabilistic Approach

- Probabilistic approach uses MLE (Maximum Likelihood Estimation)
- Probability & likelihood are proportional.
- Probabilistic Approach - finding the best parameters for the best fit line using partial derivatives.

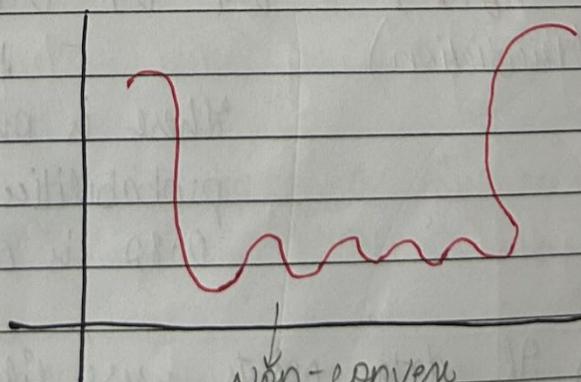
$$\Rightarrow h_{\theta}(x) = \frac{1}{1 + e^{-(\theta^T x)}} \rightarrow \begin{array}{l} 100 \rightarrow 1 \\ 90 \rightarrow 0.9 \\ 80 \rightarrow 0.8 \end{array}$$

$$\boxed{\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2}$$

↓
Loss ↓
Squared Loss

$$\left\{ \theta = \theta - \frac{\partial L}{\partial \theta} \right\}$$

$$\Rightarrow \left((1 - 0.8)^2 + (1 - 0.2)^2 + (1 - 0.4)^2 \right) \rightarrow \text{Curvature}$$



Non-convex function is not always easily differentiable.

→ Probabilistic approach loss function:

Log Loss | Binary Cross Entropy

→ Multiclassifications:-

Categorical cross Entropy pheli class

$$\rightarrow \text{Avg Loss} = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log (1-h_\theta(x^{(i)})) \right]$$

descri. class

Actually minimise the diff
b/w the probabilities distribution

2) Dog ki agar 0.9 {pheli class x},
for cat ki 0.1 hogi {then second class 1-x}

Ques) Let us see what this function does?

→ Actually minimize diff b/w 2 probabilities distributions

Eg) $y = 1 \quad 0 \quad 1 \quad 1 \quad 0$ {5 data points}

$$h_\theta(x^{(i)}) = 0.9 \quad 0.1 \quad 0.99 \quad 0.82 \quad 0.1$$

(Prediction)

There is only little bit difference in
probabilities, like 0.1 is close to 0 &
0.82 is close to 1

Eg) If data points were like

$$1 \quad 0 \quad 1 \quad 1 \quad 0$$

$$0.5 \quad 0.5 \quad 0.6 \quad 0.4 \quad 0.3$$

→ This difference is large

- Logistic Regression is a model for binary classification predictive modeling. The parameters of a logistic regression model can be estimated by the probabilistic framework called maximum likelihood estimation.
- Geometric mean → signed distance
Probabilistic mein → Maximum Likelihood Estimation
- The parameters of logistic regression equation are estimated via maximum likelihood, considering the conditional response variable as Bernoulli.

Maximum Likelihood Estimation

- i) Sometimes the probability density functions are not known.
- ii) So, these probability density functions are to be estimated from available data.
- iii) Maximum Likelihood method views the parameters as quantities whose values are fixed but unknown.
- iv) The best estimate of their value is defined to be the one that minimises the probability of obtaining the samples actually observed.

~~24/11/23~~

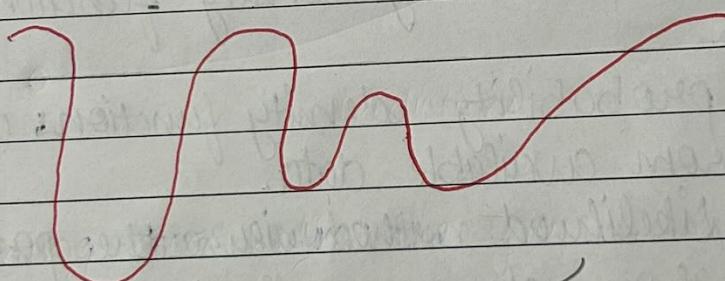
$$\text{h}_0(x) = \frac{1}{1 + e^{-\theta^T x}}$$

This is our hypothesis function,
we know this function lies in range [0, 1]
It represents line separating two values

$$\text{Loss} = \frac{1}{n} \sum_i \epsilon(y^{(i)} - h_\theta(x^{(i)})^2)$$

Squared error loss function
this is non-convex

- if class is +ve, $(1-0.8)^2$
↓
confidence for pt(a) will be high
- if class is -ve, confidence will be low
 $\Rightarrow (1-0.2)^2 + (0-0.6)^2$



multiple local minima

- For multiclass \rightarrow Categorical cross entropy
- Log Loss | Binary Cross Entropy \rightarrow for binary classification

$$\text{Avg Loss} = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log (1-h_\theta(x^{(i)}))$$

means?

confidence of point
being +ve

confidence of point
being -ve

$$h_0(x) = 0.6$$

$1 - 0.6 \rightarrow 0.4 \rightarrow$ it says that there is 60% certainty that it belongs to class 1, & 40%, that it belongs to class 0.

* (Jiska confidence jyada hai, woh model jyada accha hai)

→ I want my predictions to be as close as the real probabilities.

→ And this loss function will help to do this task.

→ Let's visualise,

Case 1

$$y^{(i)} = 1 \quad 1 \quad 1 \quad 1 \quad 1$$

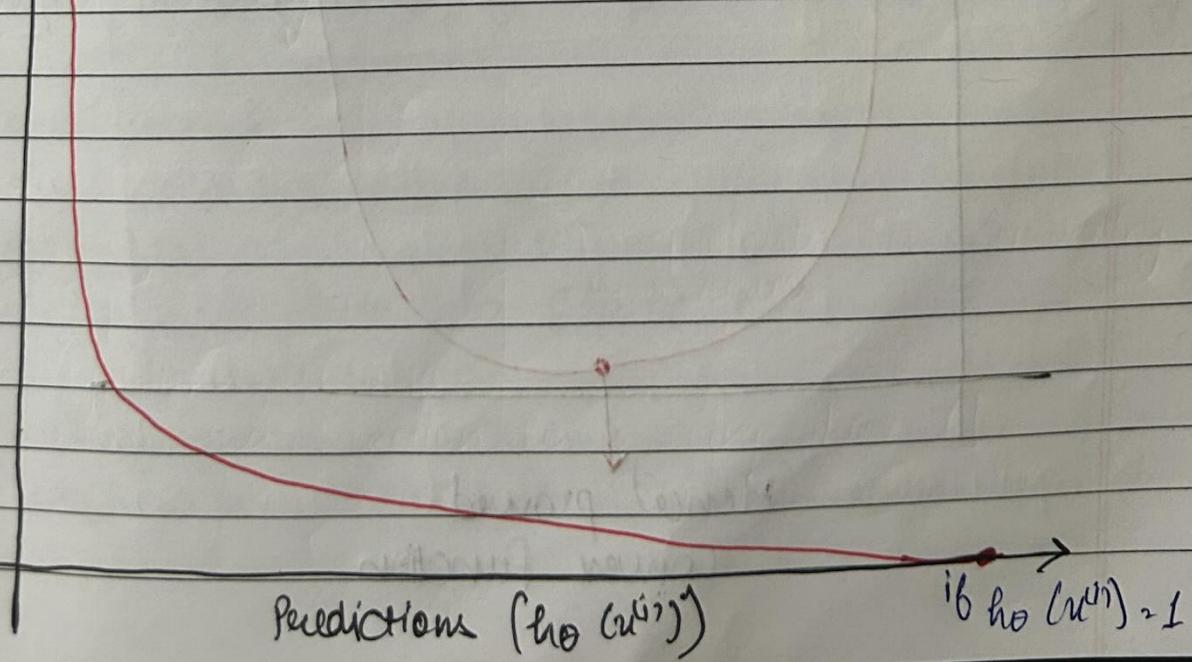
$$\text{Loss function} = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_0(x^{(i)})$$

(Second term of loss function will be zero)

$$\therefore (1 - y^{(i)}) \log (1 - h_0(x^{(i)})) \rightarrow 0$$

$$\uparrow +\infty \quad h_0(x^{(i)}) = 0$$

Loss



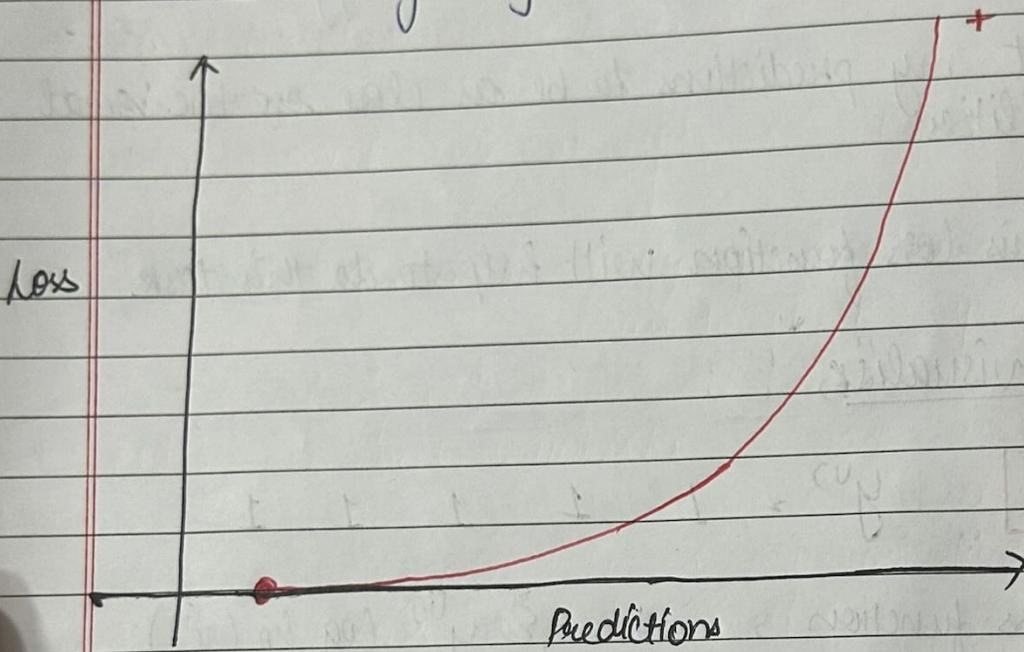
Case 2.

$$y^{(i)} = 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0$$

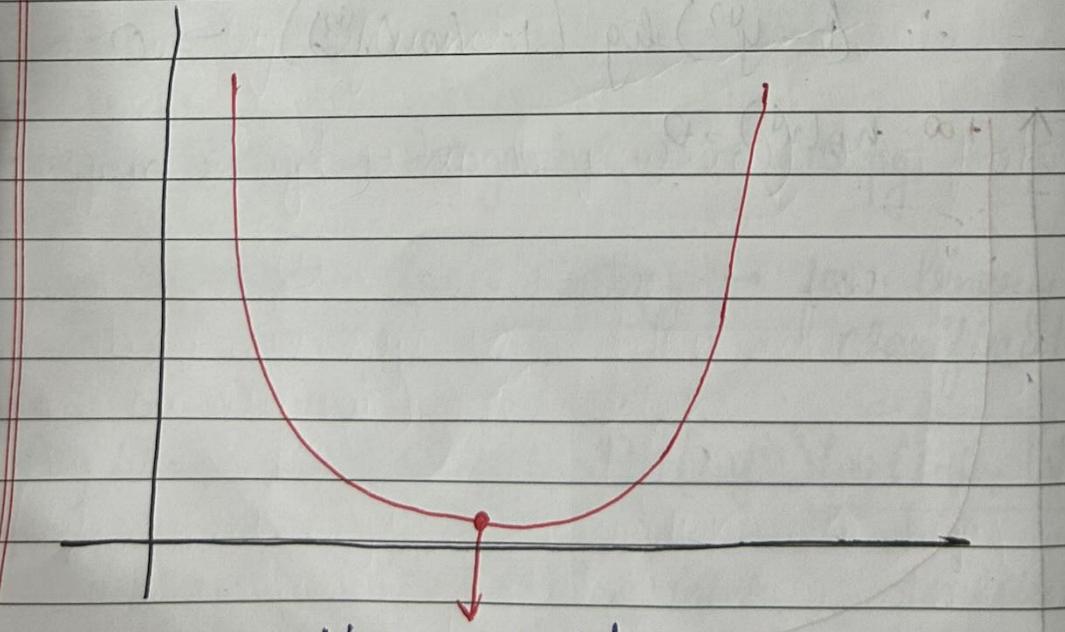
$$\text{Loss} = -\frac{1}{m} \sum_{i=1}^m (1-y^{(i)}) \log (1-\hat{y}_0(x^{(i)}))$$

(first term will be zero)

$$\therefore y^{(i)} \log (\hat{y}_0(x^{(i)})) \rightarrow 0$$



⇒ Combining both the graphs.



Hence proved
Convex Function

25/1/23

Date / /
Page No.

2) Loss = argmin $\sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$

$$J(\theta) = \sum_{i=1}^m y^{(i)} (\log h_\theta(x^{(i)}) + (1-y^{(i)}) \log (1 - h_\theta(x^{(i)})))$$

$$= - \sum_{i=1}^m [y^{(i)} - h_\theta(x^{(i)})] x^{(i)}_j \quad \begin{matrix} \rightarrow \text{Sample} \\ \downarrow \rightarrow \text{feature} \\ (\text{After doing derivation}) \end{matrix}$$

2) Final update,

$$\theta_j' = \theta_j + n \sum_{i=1}^m [y^{(i)} - h_\theta(x^{(i)})] \quad \begin{matrix} \downarrow \\ \text{Learning rate} \end{matrix}$$

2) Learning Rate

Ridge | Lasso Regression
(L2 loss) (L1 loss)

$$\boxed{\text{Loss} = \text{argmin}_{(w)} \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))}$$

2) Ridge | Lasso \Rightarrow is used to control the overfitting
(Regularisation)

Techniques

Ridge + Lasso \rightarrow Elastic Net

Ridge loss $\rightarrow \lambda_2$

Lasso loss $\rightarrow \lambda_1$

Ridge Regression & Lasso Regression

$$\text{Loss} \rightarrow \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i}) + d \|\mathbf{w}\|^2$$

Penalty term
(L2 loss)
Hyperparameters
or
L2 norm

→ Loss function define kerna hai,
↓

who optimise kerna hai, (G.D) ke through)

→ $n, d, c \rightarrow$ Hyperparameters
(have to check experimentally)

* Ridge has some drawbacks ∵ we use Lasso

① Loss → minimise (loss + regularization)
= Loss + $\lambda \|\mathbf{w}\|^2$

→ $d = 0 \Rightarrow$ Loss = minimize \Rightarrow Overfitting

* λ has big influence.

$$\Rightarrow \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i)) + d \|\mathbf{w}\|^2$$

(1) (2)

(Jab w mare hoga, tab function minimize hoga kyuki
(-) lga hua hai) → first section

(Jab w min hogta, tabhi function min hogta) → second section

→ Equilibrium to find optimal values of w

- 2) $\lambda \uparrow \rightarrow \text{underfit}$
 $\lambda \downarrow \rightarrow \text{overfit}$

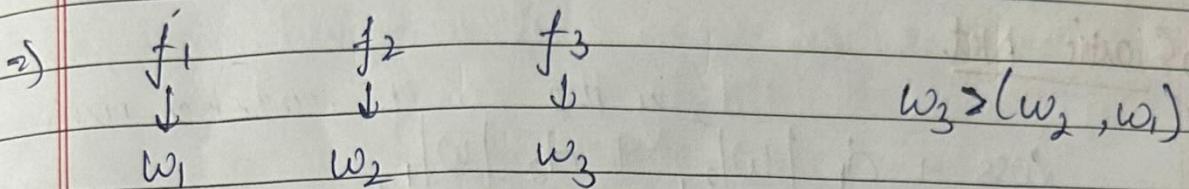
L-1 regularization $|w|_1$

$$\downarrow \\ |w_1, w_2, w_3|$$

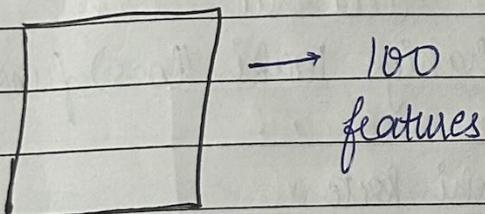
$$|w_1| + |w_2| + |w_3|$$

If we add this above penalty term,

$$\text{argmin } \sum_i \log(1 + \exp(-y_i w^T u_i)) + \lambda |w|_1$$



- 2) f_3 is most important feature
('Jiska w jyada, woh jyada imp')



* $\alpha 1$, not important features ko 0 krdeta hai

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \rightarrow \text{Sparse Matrix}$$

- 2) Feature selection $\alpha 1$ khud hi ker deta hai

→ Categorical

- ↳ Chi-square test
 ↳ Anova

Numerical

- ↳ Z-test
 ↳ t-test

→ To find the values on which the model is working very well → Grid Search CV.

$$\rightarrow d = \left[10^{-3}, 10^{-2}, 10^1, 1, 10^2, 10^3 \right] \quad \text{Grid Search CV}$$

↓

$$\left[0.001, 0.01, 1, 100, 1000 \right] \quad \text{(Random Search Induction)}$$

6 Models

→ Jis value pe, good accuracy aati hai, woh model choose kiya jaata hai.

Elastic Net

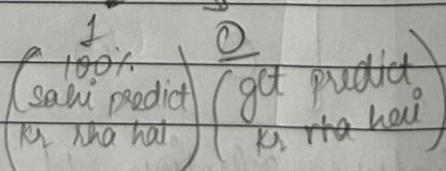
$$\text{Loss} + \underbrace{\alpha_1 |w|_2^2}_{\text{yeh slow computation}} + \underbrace{\alpha_2 |w|_1}_{\text{yeh fast keta hai, kyuki yeh unimportant ko zero karta hai}}$$

→ hume w ki min value chahiye toaki (loss) function min ho ske.

Desliye hum α use nahi kerte

6/2/23

→ Linear regression → predicted value is continuous
 Classification → predicted value is discrete



Continuous mein accuracy ka formula problem karta hai
 hum error nikalte hai

Model Predictions

Date / /
Page No.

no. of observations n = 165	Predicted NO	Predicted YES	
	(Model → NO) Real → NO)	(Model → YES) Real → NO)	
Actual NO	TN = 50	FP = 10	60 (Total actually NO obs)
Actual YES	FN = 5	TP = 100	105 (Total actual YES obs)
Ground Truth (dataset)	55 (Total predicted NO obs)	110 (Total predicted YES obs)	

~ Confusion Matrix

Ques) confusion matrix kyun aayi?
 → Kyuki accuracy score aaya.

$$\boxed{\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}}$$

Accuracy measures how often the classifier correctly predicts

⇒ Eg) 85% accurate

15% kheab

yeh kahan get hai
woh confusion matrix
bataa hai.

* accuracy (sk learn) → use

⇒ TN & TP values are important (Always +ve)
 Focus on reducing FN & FP
 (100% accuracy if FP = FN = 0)

→ Confusion Matrix

Confusion Matrix is a performance measurement for the machine learning classification problems, where the output can be two or more classes. It is a table with combinations of predicted & actual values.

* It is done on testing data.

→ Precision $\rightarrow \frac{TP}{TP + FP}$

→ Precision - Precision for a label is defined as the no. of true positives divided by the no. of predicted positives.

→ F1 Score - It gives a combined idea about precision & recall metrics. It is a maximum when precision is equal to recall.

F1 score is the harmonic mean of precision & recall

$$F1 = \frac{2 \cdot \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

		Actual Values	
		Positive (1)	Negative (0)
Predicted Value	Positive (1)	TP	FP
	Negative (0)	FN	TN

1 → Fraud credit card

0 → Not fraud

- Actual → NO } Toh fir kisi padega, kya kisi in
Predicted → YES } real mein kisi hai
- Actual → YES } Toh problem hogi
Predicted → NO
- Precision & Recall Trade off → improves the model a little bit.

7/2/23

Confusion Matrix

It is a performance measurement for machine learning classification problem where output can be two or more classes. It is a table with 4 different combinations of predicted & actual values.

Key classifications

- Accuracy
- Recall
- Precision
- F1 - Score

		Actual	
		Pos (1)	Neg (0)
Predicted	P	TP	FP
	N	FN	TN

- True positive - Predicted something yes & it came out to be Yes.
- True Negative - predicted no it comes NO
- False Positive - predicted its Yes , it came out to be NO.
- False Negative - predicted its no , it comes Yes.

→ Accuracy-

$$\frac{\text{True Positive} + \text{True Negative}}{\text{True Pos} + \text{True Neg} + \text{False Pos} + \text{False Neg}}$$

$$\Rightarrow \frac{\text{Correct Prediction}}{\text{Total Prediction}}$$

→ Precision - Ability This will give you how many predictions you made are right

→ Recall - Ability to recall to the number of events

⇒ Precision is the ratio of a no. of events you can correctly recall to the no. of all events you recall
(min of ✓ & ✗ events)

$$\text{Precision} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive.}}$$

* In general words,

~~Precision~~ → Accuracy of positive values

~~Recall~~ → Accuracy of negative values

$$\text{Precision} \uparrow \frac{\text{TP}}{\text{TP} + \text{FP} \downarrow}$$

$$\text{Recall} \uparrow \frac{\text{TP}}{\text{TP} + \text{FN} \downarrow}$$

⇒ Applications

Eg) Image Classification

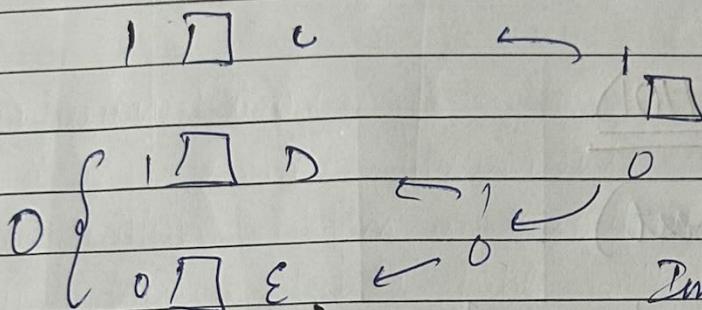
Precision & Recall → for which predictor

Biometric - Precision & Recall hi kaam keta hai

Recall

- choose this accuracy
- Precision

- 1 v/s Rest
logistic → for multiple



Imbalanced dataset
(leads to overfitting)

Code Discussion,
penalty, "none"
uses neither L1 or L2

- by default uses L2
- solver?

9/2/23

FEATURE SELECTION TECHNIQUES In MACHINE LEARNING

