

⇒ Understanding feature space

(Q) Define Slp data of a typical ML alg?

Ans. Slp data is a collection of features that have been quantitatively measured from some object or event that we want the machine learning system to process.

Linear Regression

Task is to prepare a smoothie and serve.

Goal is To predict the calories

case I

Experience: 15 smoothies are given

Performance: by primarily guessing

Only we know our sample points

Smoothie = 1 cup of yogurt + secret ingredient
+ some color dyes.

Calories

avg mass of 150 smoothies is 236.9

case II Experience: 15 smoothies are given

148, 240, 154, 265, 281, 204,

Smoothie calories

priya

265

2

281

first Mean

= 236.9

Median = 240

value	14	15	191
	↓	↓	↓

If 16 - 234
actually 400

case III you can weight the (smoothie) ingredient

⇒ 2D features

Calorie	wt of yogurt	wt of secret ingredient
148	100 g	100 g
240	100 g	100 g
154	100 g	100 g

s.no	weight	calorie
1	88.9	265
2	156.9	281

15	110 g	191
100 g	100 g	100 g

$$\text{Calories} = \text{mean} + O + \text{weight}$$

↳ Intercept → slope

for another model we need

Task or Goal : Predict

Model 1 : wild guess

↳ very bad

↳ we are no

Model 2 : we use

268, ..., 191, ...

↳ NO

↳ Prec

but actual

↳ Pr

↳ Unc

Model 3 :-

⇒ Let us co

Calories

actual * Mod

diff * Mod

Model 3 :-

(Root Mean S

not have
or even
been
that we

for another model $c = w_0 + w_1 * \text{weight}$

so we need to find w_0, w_1 from the data

Task or Goal : Predict calories

Model 1: Wild guess

↳ very bad model

↳ we are not having any experience

Model 2: we use experience of 15 previously made smoothies

268, 271, ..., 191 \Rightarrow mean = 286.9

Median = 240

↳ No features are used

↳ Prediction for 16th smoothie is 288

but actual calories = 445

$$\frac{445}{288} = 1.5$$

↳ Under estimation by 208

↳ Under estimation by 208

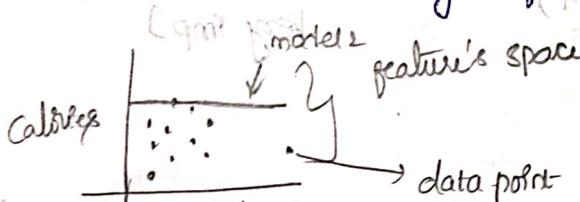
Model 3: more data

$$P_{\text{RMSE}} = 32.08$$

↳ doesn't mean more data points

SHI \Rightarrow it means more information about the data.

\Rightarrow Let us consider weight of the ingredient



weight	calorie
1	200
2	210
3	220
4	230
5	240
6	250
7	260
8	270
9	280
10	290
11	300
12	310
13	320
14	330
15	340

Model 2: calories = 286.9 + 0 * weight

↳ Not a good model.

Model 3: calories = $w_0 + w_1 * \text{weight}$

RMSE = 84.63 (for model 2)

(Root Mean Square error)

Model 3 - $C = 236.9 + 0.5 \times \text{weight}$ RMSE = 69.3

41 - $C = 236.9 + 1 \times \text{weight}$ RMSE = 69.2

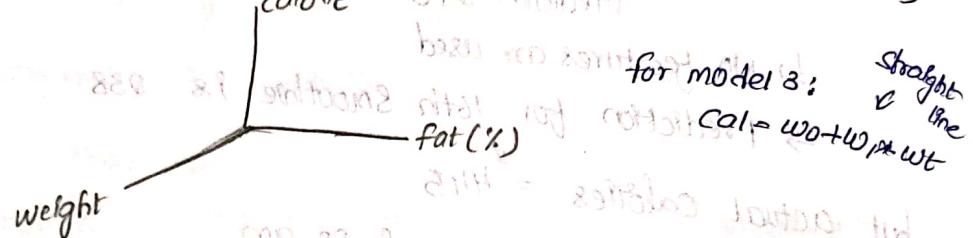
51 - $C = 200 + 1 \times \text{weight}$ RMSE = 59.8

Finally! - $C = 145 + 1.35 \times \text{weight}$ RMSE = 48.8

(to make it much better)

Model 6 - \rightarrow for 16th smoothie the prediction = 293.8
↳ underestimated by $455 - 293.8 = 151.2$

Model 4 - $C = w_0 + w_1 \times \text{weight} + w_2 \times \text{Fat (\%)} \quad \text{Plane (3D)}$



$C = 163 + 1.42 \times \text{weight} + 51.9 \times \text{Fat (\%)}$

These are the best parameters that give good performance

RMSE = 46.9

16th smoothie prediction is underestimated by an amount

so consider the other features

- Fat (%) (domain knowledge is very imp)

- carbohydrates

- protein

Model 8 - $C = w_0 + w_1 \times \text{Fat (\%)} + w_2 \times \text{Carbohydrates} + w_3 \times \text{protein}$

$C = 143 + 8.9 \times \text{fat (\%)} + 3.9 \times \text{carbohydrates} + 4.3 \times \text{protein}$

↳ this is going to be good model

The calories present in 1 cup RMSE 210
of yogurt 16th smoothie Model = $43 + 48 = 91$

Fat	Carbohydrates	Protein	Calories
8.9	3.9	4.3	143

Summary - Feature selection is very imp

Classification -

Task 1 : Lion

Shape feature

Task 2 : horse

Shape feature

but reg

ML

raw

IP data

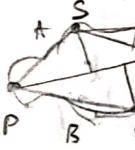
interpolation

grow

input

configuration

Shape De



rotate

SO,

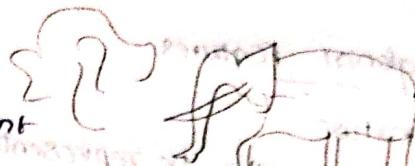
Adaboost

Unsupervised

Classification - computer vision

Task 1: Lion or Elephant

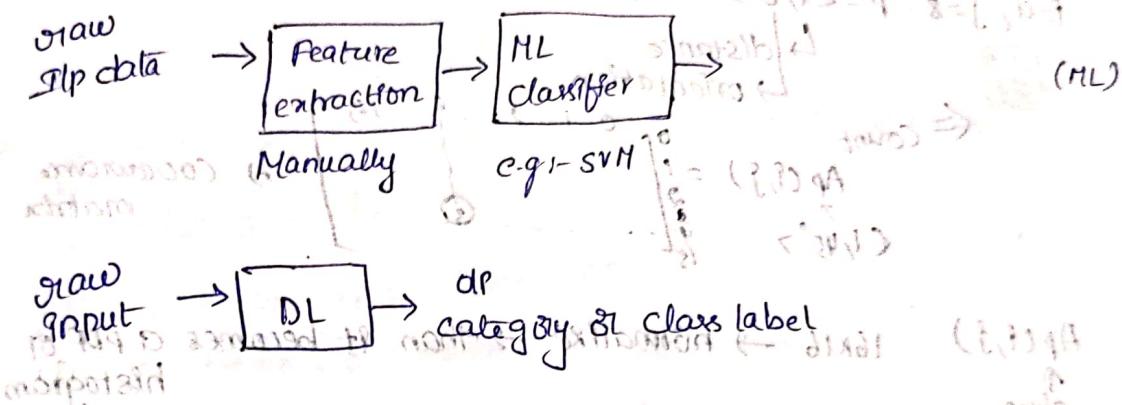
Shape features are sufficient



Task 2: horse vs zebra

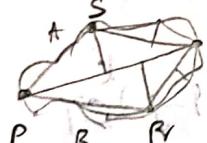
Shape features are not useful
but regional (texture) features are useful

ML vs DL



Shape Descriptors : Polygonal representation

The borders of an object



→ by joining it we get a polygonal shape (perfect)

If we want better shape → repeat the process.

→ start treating it as some function

rotate it → whenever we call it has probability

So, It has

$$\text{mean}(\mu) =$$

$$\text{variance}(\sigma^2) =$$

$$\text{skewness} =$$

density fun?

for Boundary Segment 1 (B_1)

$$\text{Find } N = \frac{\text{Area}}{\sigma^2}$$

$$B_2 \Rightarrow$$

$$\text{for } 8 \times 3 \Rightarrow 24 \text{ features}$$

30/12/22

Regional Features

① Texture

Texture can be represented using cooccurrence matrix

4-bit image of size 8x8

gray level i
gray level j

Some position vector

$$P=4, S=8 \quad P=\langle p, \theta \rangle = \langle 1, 45^\circ \rangle$$



↳ distance

↳ orientation

↳ count

$$Ap(p, \theta) = \begin{cases} 0 & \text{if } (i-j) \neq p \\ 1 & \text{if } (i-j) = p \end{cases}$$

10	9	1	9	5	8	11	9
6	5	15	12	6	4	3	2
9	3	2	10	6	8	4	5
2	2	10	3	7	5	6	1
21	0	11	0	10	9	8	2
8	14	4	1	6	0	7	2
2	3	0	9	11	6	3	9
7	2	8	2	6	12	6	7

↳ cooccurrence matrix

$$Ap(p, \theta)$$

$16 \times 16 \rightarrow$ normalize \rightarrow then it becomes a PDF or histogram
given $\langle 1, 45^\circ \rangle$ and denote it has $c_{ij} 16 \times 16$

① Max prob : $\max_{p, j} c_{ij}$ if uniform

② Elementwise difference moment : $\sum_i \sum_j (i-j)^k c_{ij}$

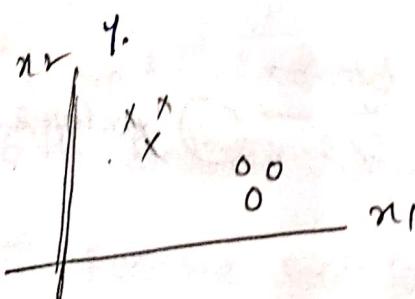
③ Inverse elementwise difference : $\sum_i \sum_j c_{ij} / (i-j)^k$

④ Uniformity : $\sum_i \sum_j c_{ij}^2$

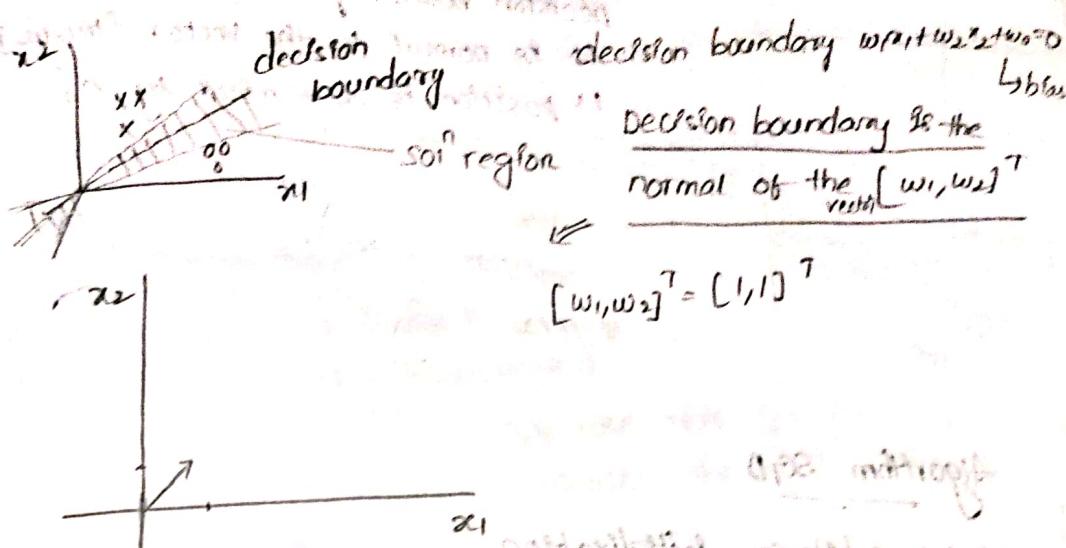
⑤ Entropy : $-\sum_i \sum_j c_{ij} \log(c_{ij})$

⇒ Dataset = ① ② ③ ④ ⑤ Label

horse/zebra



Linear classifications



$\Rightarrow \sum_{i=1}^d w_i x_i + w_0 = 0$, the decision boundary $w^T x + w_0 = 0$

Decision rule

If $w^T x + w_0 > 0$; then $x \in C_1$

If $w^T x + w_0 < 0$; then $x \in C_2$

We want to unify them into one rule.

$$\text{and if } [w_1 \ w_2 \dots \ w_d] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \\ 1 \end{bmatrix} + w_0 = 0$$

$$[w_1 \ w_2 \dots \ w_d \ w_0] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \\ 1 \end{bmatrix} = 0$$

Decision rule

$$a^T y > 0, y \in C_1$$

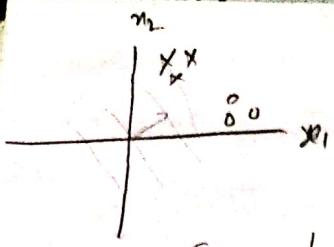
$$a^T y < 0, y \in C_2$$

Now, negate y of class 2

Unified decision rule:

$$a^T y \geq 0$$

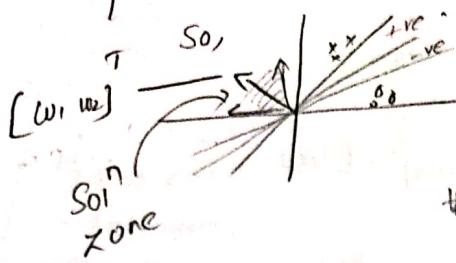
positive sign
with y



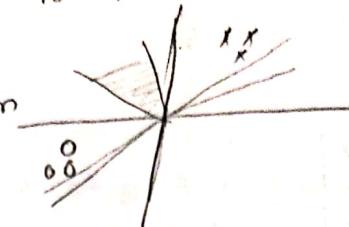
$$w_1x_1 + w_2x_2 + w_0 = 0$$

Decision boundary

- ↳ w is normal to the vector $[w_1, w_2]^T$
- ↳ position is determined by w_0
for $a^T y > 0$



After
negation
 $y \in C_2$



Algorithm SGD

- $a(0) \leftarrow$ arbitrary initialization
- $a(k-1) \rightarrow a(k')$

2/1/23.

Regression

- ↳ fitting the hyperplane

Classification

- ↳ finding the DB which separates

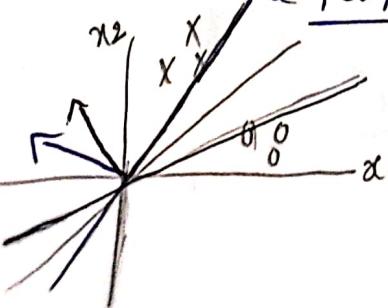
we can achieve above two classes using loss fun.,

In Regression $\text{calculus} = w_0 + w_1 * f_1(\%) + w_2 * f_2(\%) + w_3 * f_3(\%)$
Prediction = $w_0 + w_1 * \text{feature}_1 + w_2 * f_2 + w_3 * f_3$

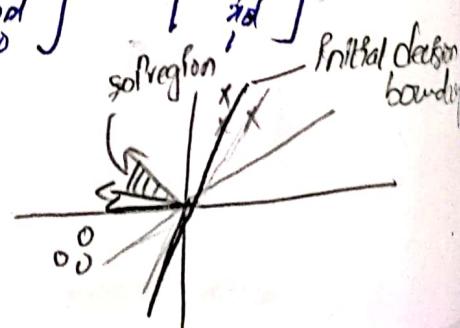
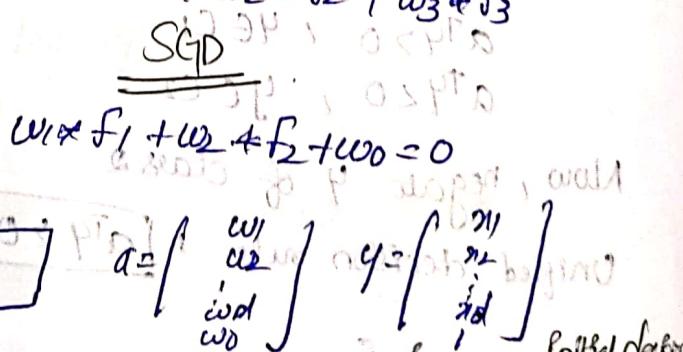
In classification

decision boundary is a hyperplane

Decision rule $|a^T y > 0|$



After negating
 $y \in \text{class } 2$



Training Alg

- (i) $a(0) \leftarrow$ arbitrary
- (ii) $a(k) \leftarrow a(k-1)$

how much the corres

$$a^T y > 0$$

$$w_1 x_1 + w_2 x_2 + w_0$$

(b) Error $J =$

$$\boxed{\nabla J =}$$

(c) $a(k) = a(k')$

η

$$\boxed{a(k) = a(k')}$$

Importance

$\eta = 0.1, 0.2$ (use

Inference

In Regression

In Classification

⇒ Inference means

Performing the

farly the model.

In the traini

$$\boxed{MSE = }$$

Training Alg

- (i) $a(0) \leftarrow$ arbitrary (number) initialize
(ii) $a(k) \leftarrow a(k-1)$, in such a way that it min some performance index
 how much the corresponding wt is useful in deciding the boundary

$$a^T y > 0$$

$$w_1 x_1 + w_2 x_2 + w_0 > 0$$

$$(b) \text{ Error } J = -\sum a^T y$$

$\forall y$ misclassified

$$\boxed{\nabla J = -\sum y}$$

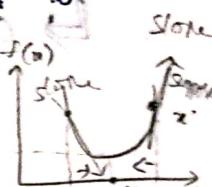
we use this gradient to update the weights.

$$(c) a(k) = a(k-1) - \eta \nabla J$$

η - learning rate

$$\boxed{a(k) = a(k-1) + \eta \sum y}$$

$\forall y$ misclassified



$$x^* = \arg \min_x f(x)$$

here x is wt

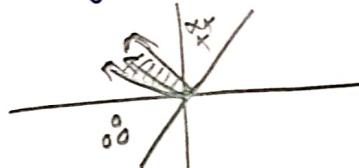
$$a(k) = a(k-1) - \eta \frac{\partial f(x)}{\partial x}$$

$$\text{if } \eta = 1 \Rightarrow a(k) = a(k-1) + \sum y$$

$\forall y$ misclassified

+ve slope \rightarrow decrease η
 -ve slope \rightarrow increase η

$$\eta = 0.1, 0.2 \text{ (usually)}$$



Inference

In Regression $w_1 x_1 + w_2 x_2 + w_0 = 0$

In classification $a^T y = 0$ & $w^T x = 0$ Ag by I where $w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$ ($d+1 \times 1$)

\Rightarrow Inference means (trying to do that) performing the tasks on an unseen data by so $x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$.

farby the model. (Test data)

In reg task, have training cur $= \frac{1}{m} \| \hat{y} - y \|_2^2$

In the training phase, have training cur $= \frac{1}{m} \| \hat{y} - y \|_2^2$

2nd Norm \downarrow

$$\boxed{MSE = \frac{1}{m} \| \hat{y} - y \|_2^2}$$

* Tasks :-

Now we have

$$(train) \quad y^{(train)} \quad x^{(train)} \quad y^{(test)} \quad x^{(test)} \quad & w$$

$$\text{Training Error } MSE = \frac{1}{m} \| \vec{y}^{(train)} - \vec{y}^{(train)} \|_2^2$$

$$= \frac{1}{m} \| w^T x^{(train)} - y^{(train)} \|_2^2$$

Similarly we have Test error on Test set

$$\text{Test error} = \| \vec{y}^{(test)} + y^{(test)} \|_2^2 \xrightarrow{\text{Optimization vs ML}}$$

↳ Called as Generalization Error. Generalize the model well on the test data

$$P \cdot P + (1-P) \cdot 0 = (1-P)$$

for training

$$P \cdot P + (1-P) \cdot 0 = (1-P) \leftarrow 1-P$$

$$(1-P) \cdot 0 = 0$$

for test

$$P \cdot P + (1-P) \cdot 0 = (1-P) \leftarrow 1-P$$

Conclusion

$$\begin{cases} P \\ 1-P \end{cases} = 0 \Rightarrow P = 0.5 \quad \text{and} \quad 1-P = 0.5$$

(both ob of prob) same weight
means no ob weight gain

$$\| P - \vec{P} \| \rightarrow 0 \Rightarrow \text{no gain}$$

↳ weight gain

$$\| P - \vec{P} \| \rightarrow \frac{1}{m} = \text{ERM}$$

we have

$$x^{(train)}, y^{(train)}, x^{(test)}, y^{(test)} \in \mathbb{R}^n$$

Training Error MSE =

$$\frac{1}{m} \sum_{i=1}^m \|w^T x^{(train)} - y^{(train)}\|_2^2$$

Similarly we have Test error on test set

$$\text{Test error} = \frac{1}{n} \sum_{i=1}^n \|g(x^{(test)}) - y^{(test)}\|_2^2$$

\Leftrightarrow called as Generalization Error.

3.1.1.2
Generalization:

It is the ability of machine learning system to perform well on the new unseen data.

\Rightarrow It means to say we want the testing error also to stay as well

\Rightarrow Capacity, underfitting & overfitting

$$\text{Prediction} = w_0 + w_1 x_1 + w_2 x_2 + \dots$$

The different weight vector values gives us different functions.

\rightarrow The range of functions which the model can learn is called the capacity of that model.

\rightarrow The goal then becomes selecting one among those functions which will best fit the data (S) which will well separate the data corresponding to two categories respectively.

For regression
 $y = w_0 + w_1 x_1 + \dots + w_n x_n$

\Leftrightarrow If its capacity is it only be only a straight line.

* Challenges of ML

① Data size of data

Suppose we have

+ 2nd feature
 $10^2 \Rightarrow 100$ regions

+ 3rd feature
 $10^3 \Rightarrow 1000$ regions

② Local convergence

③ Non-parallel

\Leftrightarrow Suppose

ML, Appro

Row \rightarrow Feature
Column \rightarrow Classification

DR Approach

Row
Column

- Reasons - the model complexity and distribution are almost identical
-
- \Rightarrow Model order is too high \Rightarrow Overfitting
 - \Rightarrow For the given data, it is best fitting
 - \Rightarrow Many varieties of curves \Rightarrow Model complexity
 - \Rightarrow Overfitting
 - All of the above will also use Linear Regression since parameters are zeroed out here, training error is zero.

* Challenges of ML during us to DL - 3 challenges

① the curse of (data) dimensionality

Suppose we have only 1 feature \Rightarrow need only 100 regions \Rightarrow 100 decision boundaries

+ 2nd feature $10^2 \Rightarrow 100$ regions \Rightarrow we need to separate with 100 decision boundaries

+ 3rd feature

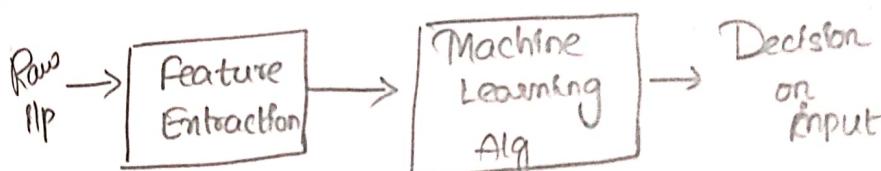
$10^3 \Rightarrow 1000$ regions \Rightarrow Local constancy and smoothness regularization.

② Local constancy and smoothness regularization.

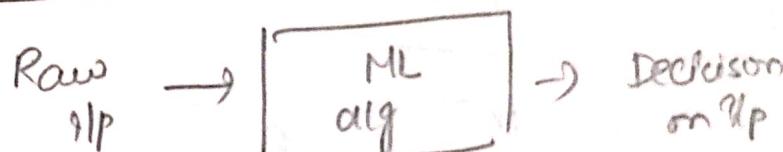
③ Manifold learning.

\Rightarrow Suppose take an e.g. \Rightarrow differentiating dog & cat

ML Approach



DL Approach

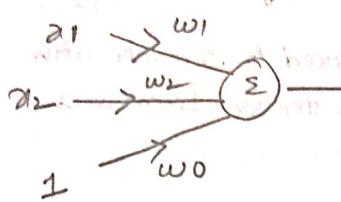


In traditional ML the features are decided by the user accordingly we need feature extraction algorithm. In case of deep learning we will not tell what features are to be analyzed. Instead during the training phase the model learns features from the training samples that are useful in taking decisions.

→ During inference similar features are extracted by the model from the so far unseen data in order to make a decision.

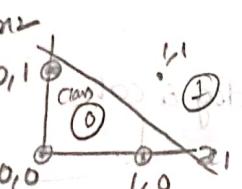
Perception :-

$$y = w_1x_1 + w_2x_2 + w_0$$

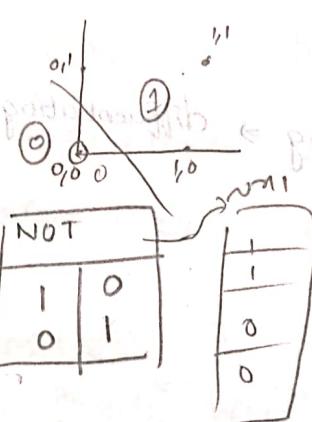


$$y = \begin{cases} 1 & \text{if } w^T x > 0 \\ 0 & \text{if } w^T x \leq 0 \end{cases}$$

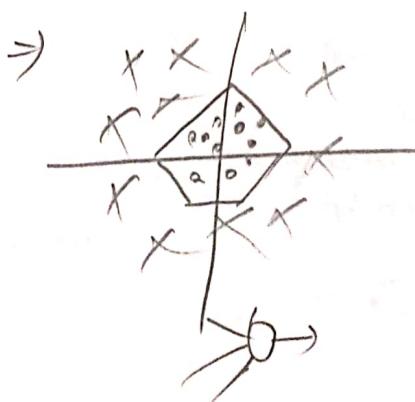
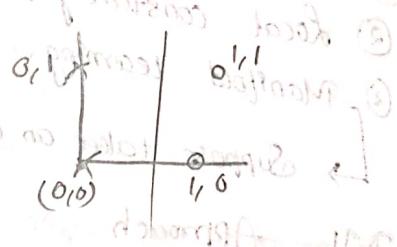
AND



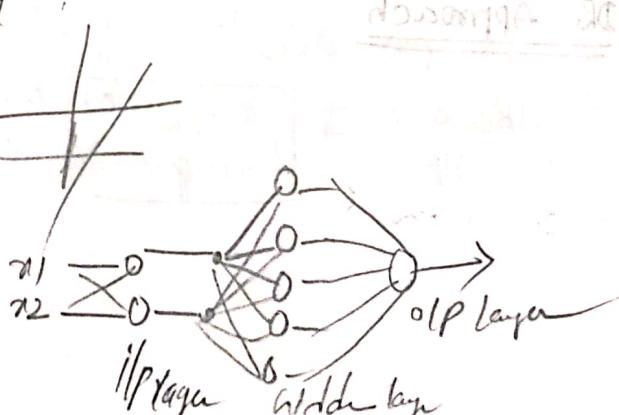
OR



NOT



for



$$a/1/23. \quad y = w^T x$$

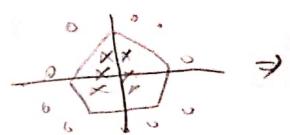
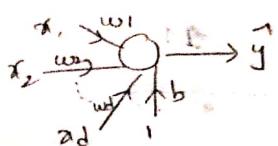


Fig (a)

$$\Rightarrow x_1 - \frac{w_{11}}{w_{11} + w_{12}} \rightarrow y \\ x_2 - \frac{w_{21}}{w_{21} + w_{22}} \rightarrow y$$

$$w = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix}, \quad G = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix}$$

we tried

$$\text{ie } w \in R^{2 \times 2}$$

$$w \in R^{2 \times 1}$$

$$G \in R^2$$

$$b \in R^1$$

→ for above fig

More clearly,

$$\Rightarrow \hat{y} = f^{(m)}(f)$$

→ for fig (a),

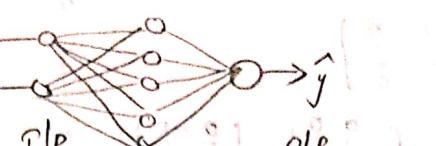
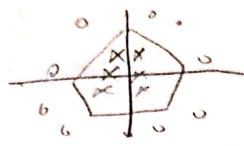
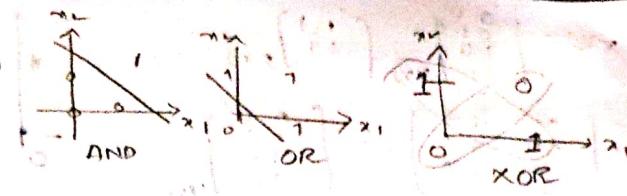
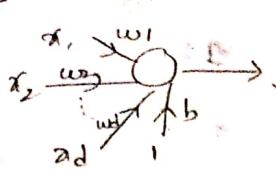
$$w = \begin{bmatrix} \dots \end{bmatrix}$$

Feed forward

by the
algorithm.
features
phase the
that are

by the
note a

9/11/23. $y = w^T x$



$$\Rightarrow x_1 \xrightarrow{w_{11} w_{12}} c_1 \\ x_2 \xrightarrow{w_{21} w_{22}} c_2$$

$$w = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix}, \quad c = \begin{bmatrix} c_1 \\ c_2 \end{bmatrix}$$

$$y = f^*(c)$$

$$y = f^*(m)$$

↳ the "fun" representing XOR

$G = [c_1 \dots c_2]_{2 \times 1}, [b]_{1 \times 1}$ ↳ we want to learn this fun'

we tried $\hat{y} = f(m; w, b)$ ↳ linear classifier

↳ this is failed to
separate XOR

$$\text{ie } w \in R^{2 \times 2}$$

$$w \in R^{2 \times 1} \text{ where}$$

R is set of real nos.

$$G \in R^2$$

$$b \in R^1$$

→ for above fig (a) $\Rightarrow \hat{y} = f^2(f^1(m))$ ↳ IIP layer
↳ fun for 1st layer (S1)
↳ second layer (S2) OLP layer

More clearly, $\hat{y} = f(x; w, G, b)$

The
Parameters (S1)
Parameters (S2)

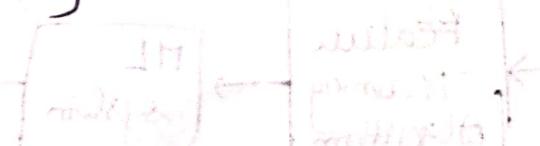
Learnable S1 trainable
parameters

$$\Rightarrow \hat{y} = f^{(m)}(f^{(m-1)}(\dots(f'(x))\dots))$$

→ for fig (a), $X = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$ & $y = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$ for XOR gate

$$w = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \quad w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, \quad G = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \quad b = 0$$

feed forward Neural N/W



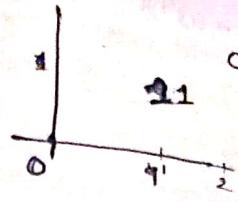
and maximizing it

2nd maximization
with respect to w

$$xw = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}_{4 \times 2} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}_{2 \times 2}$$

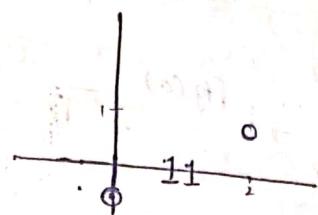
$$= \begin{bmatrix} 0+0 & 0+0 \\ 0+1 & 0+1 \\ 1+0 & 1+0 \\ 1+1 & 1+1 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}_{4 \times 2}$$

$$y = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 2 \end{bmatrix}$$



$$xw + c = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} + \begin{bmatrix} 0 & -1 \\ 0 & -1 \\ 0 & -1 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 2 & -1 \end{bmatrix} \text{ & } y = \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix}$$

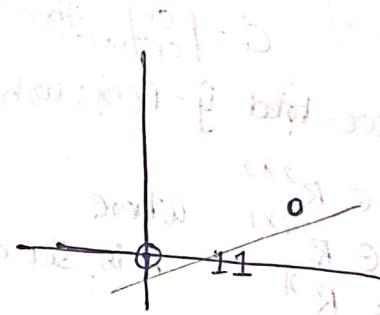


$w_1x_1 + w_2x_2 + \dots + w_0$ \rightarrow is always linear

Consider $\text{ReLU}(xw + c)$

$$\text{ReLU}(z) = \max(0, z)$$

$$\text{ReLU}(xw + c) = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 2 & 0 \end{bmatrix} \text{ & } y = \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix}$$



$$= \begin{bmatrix} xw + c \\ xw + c \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 2 & 0 \end{bmatrix} \text{ & } y = \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix}$$

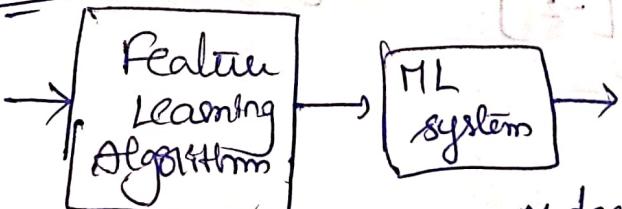
$$\text{here } y = \hat{y}, \text{ Err} = 0$$

$$\Rightarrow f(a) = f^2(f'(a))$$

\downarrow Linear classifier feature learning layer

\downarrow feature learning layers Linear classifier

ML vs DL



\downarrow Designer has to consider which features to be included

Random forests how to get them

Activation Function	Equation	Derivation	MLP	DL
Step function	$\phi(z) = \begin{cases} 0 & z < 0 \\ 1 & z \geq 0 \end{cases}$	$\phi'(z) = 0$	$\phi'(z) = 0$	$\phi'(z) = 0$
Binary step (g)	$\phi(z) = \begin{cases} 0 & z < 0 \\ 0.5 & z = 0 \\ 1 & z > 0 \end{cases}$			
Graph	$y = \phi(z)$	$y = \phi(z)$	$y = \phi(z)$	$y = \phi(z)$
Graph	$y = \phi(z)$	$y = \phi(z)$	$y = \phi(z)$	$y = \phi(z)$
Graph	$y = \phi(z)$	$y = \phi(z)$	$y = \phi(z)$	$y = \phi(z)$

Activation Function

Equation

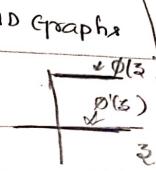
Derivation

1D Graphs

① Binary step (δ_1)
Unit step (δ_1) Heaviside

$$\phi(z) = \begin{cases} 0 & z < 0 \\ 0.5 & z = 0 \\ 1 & z > 0 \end{cases}$$

$$\phi'(z) = 0$$

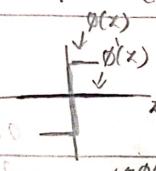


MLP

② Signum

$$\phi(z) = \begin{cases} -1 & z < 0 \\ 0 & z = 0 \\ +1 & z > 0 \end{cases}$$

$$\phi'(z) = 0$$

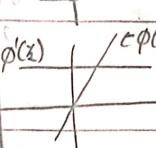


MLP

③ Linear

$$\phi(z) = z$$

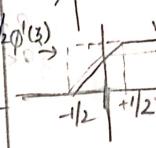
$$\phi'(z) = 1$$



④ Piecewise Linear

$$\phi(z) = \begin{cases} 1 & z \geq 1/2 \\ z + 1/2 & -1/2 < z < 1/2 \\ 0 & z \leq -1/2 \end{cases}$$

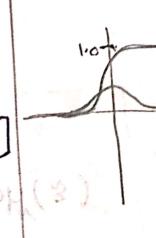
$$\phi'(z) = \begin{cases} 1 & z \geq 1/2 \\ 1/2 & -1/2 < z < 1/2 \\ 0 & \text{else} \end{cases}$$



⑤ Sigmoid

$$\phi(z) = \frac{1-e^{-z}}{1+e^{-z}}$$

$$\phi'(z) = \sigma(z)[1-\sigma(z)]$$

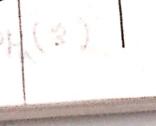


Logistic Regression
Binary Classification

⑥ Tanh

$$\phi(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\phi(z) = 1 - \tanh(z)$$



⑦ ReLU

$$\phi(z) = \max(0, z)$$

$$\phi'(z) = \begin{cases} 1 & z > 0 \\ 0 & z \leq 0 \end{cases}$$

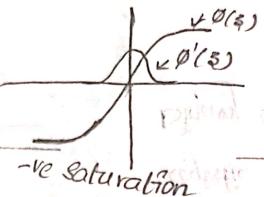


DL
Feature Processing
+ classifier

6) Tanh

$$\phi(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

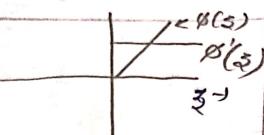
$$\phi'(z) = 1 - \tanh^2(z)$$



7) Softplus

$$\phi(z) = \log(1 + e^z)$$

$$\phi'(z) = \frac{e^z}{1 + e^z}$$

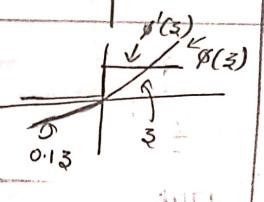


8) ReLU

$$\phi(z) = \max(0, z)$$

$$(or) \phi(z) = \begin{cases} z & z \geq 0 \\ 0 & z < 0 \end{cases}$$

$$\phi'(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$

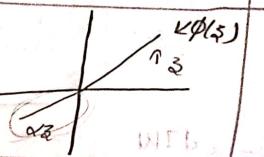


9) Leaky ReLU

$$\phi(z) = \max(0.1z, z)$$

$$(or) \phi(z) = \begin{cases} z & z \geq 0 \\ 0.1z & z < 0 \end{cases}$$

$$\phi'(z) = \begin{cases} 1 & z \geq 0 \\ 0.1 & z < 0 \end{cases}$$



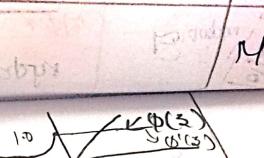
10) Parametric

Relu

$$\phi(z) = \max(\alpha z, z)$$

$$(or) \phi(z) = \begin{cases} z & z \geq 0 \\ \alpha z & z < 0 \end{cases}$$

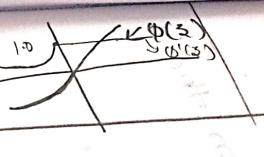
$$\phi'(z) = \begin{cases} 1 & z \geq 0 \\ \alpha & z < 0 \end{cases}$$



11) Softmax

$$\phi(v_i) = \frac{\exp(v_i)}{\sum \exp(v_j)}$$

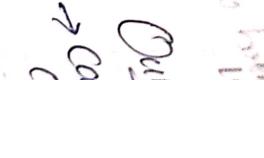
$$\phi'(z) = \begin{cases} 1 & z \geq 0 \\ \alpha e^z & z < 0 \end{cases}$$

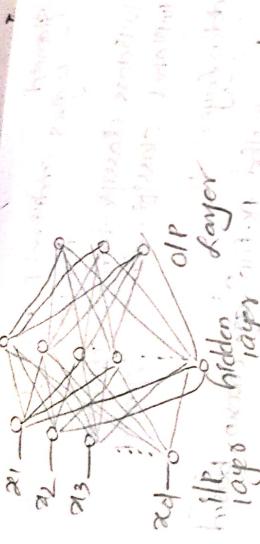


12) Exponential / Linear unit

$$\phi(z) = \begin{cases} z & z \geq 0 \\ \alpha(e^z - 1) & z < 0 \end{cases}$$

$$\phi'(z) = \begin{cases} 1 & z \geq 0 \\ \alpha e^z & z < 0 \end{cases}$$





Output Layer

→ Can see the labels or targets be np data it can see.

→ during training phase; parameters are updated such that the o/p becomes close to the target.

→ $y = f(x; w, b)$
 \Rightarrow hidden layer of model may

Hidden Layers → Not involved in it

→ They don't have access to the labels & targets for training datasets.

→ Obviously; the hidden layers won't get any direct guidance from the dataset in updating its

Parameters: → So it is back propagation

→ Thus they are called hidden layer

→ Output layer

for the activation function chosen based on the type of task.

e.g.: ① Regression -
② Binary Classification -
③ Multi-class classification -
④ ...

$$\text{linear unit}$$

$$f(x) = \sum_{i=1}^n w_i x_i + b$$

$$\text{sigmoid}$$

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

$$\text{tanh}$$

$$f(x) = \tanh(x)$$

$$\text{ReLU}$$

$$f(x) = \max(0, x)$$

$$\text{softmax}$$

$$f(x) = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

Sigmoid gives conditional probability
 $p(y=1/x)$

- ③ Multiclass classification - softmax
- ④ Multilabel classification - sigmoid.

Hidden layers

For hidden layers the "act. fun" are chosen based on the type of neural netw architecture.

- ① CNN - ReLU
- ② RNN - Tanh (or) Sigmoid

20/1/23. Cost/Loss function

→ Cost-fun is a method of evaluating how well your algorithm is modelling your dataset.

→ It's a mathematical fun' of the parameters of the model denoted by $J(\theta)$

where θ represents the parameter vector of the model.

Loss → represents error corresponding to one training sample.

Cost → represents error over the entire dataset.

1) Linear

Regression Loss

- ① MSE (Mean Square error)
- ② MAE (Mean Absolute error)
- ③ Huber loss

Classification Loss

- ① Binary Cross-Entropy
- ② Categorical cross-entropy

Auto encoders

- ① KL divergence

GAN

- ① Discriminative
- ② Minmax GAN

Object Detection

- ① Focal loss

word Embeddings

- ① Triplet loss

Problem Setting

$$x, y \sim P_{\text{data}}$$

$$\hat{y} \sim P_{\text{model}}$$

$$\hat{y} = f^*(x)$$

hidden

w, c

$y =$

Auto encoders

- ① KL divergence loss

GAN

- ① Discriminative loss

- ② Minmax GAN loss

Object Detection

- ① Focal loss

word Embeddings

- ① Triplet loss

Problem Setting

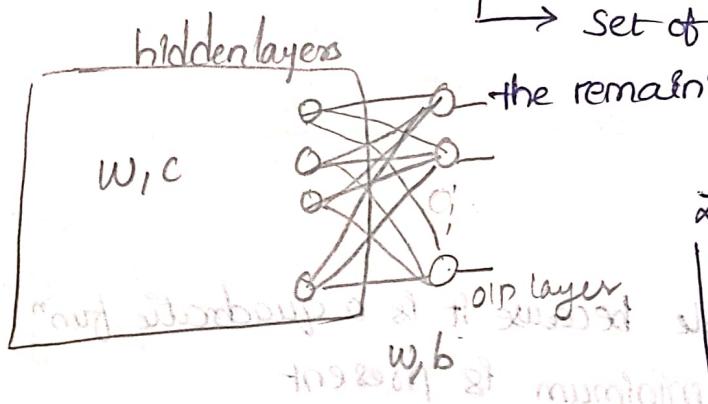
$x, y \sim P_{\text{data}}$ \rightarrow the entire dataset coming from distribution P_{data}

$\hat{y} \sim P_{\text{model}}$ \Rightarrow predictions produced by the model according to its own learned distribution

$P_{\text{model}} \rightarrow$ complete set of model parameters.

$$\hat{y} = f^*(x; \Theta) \quad \begin{matrix} \text{set of parameters corresponding} \\ \text{to OLP layers} \end{matrix}$$

\rightarrow set of parameters of all the remaining layers of DL model



$$z = w^T h + b \quad \begin{matrix} \text{further vector} \\ \text{obtained from the last hidden} \\ \text{layer.} \end{matrix}$$

\rightarrow Linear OLP of the OLP layer.

$$y = \phi(z) = \phi(w^T h + b)$$

\rightarrow Activation function
This decides what is the type of OLP units.

e.g.: If you choose ϕ to be Sigmoid then the OLP unit is said to be of Sigmoid type.

\rightarrow Goal we want \hat{y}_i as close to y_i as possible \leftarrow prediction target

\rightarrow If they are not equal, then there is error given by $y_i - \hat{y}_i$ \leftarrow +ve or -ve

\rightarrow we define cost function $J(\theta)$ which is a function of model parameters & we try to update θ such that $J(\theta)$ gets minimized.

Linear Regression Loss

which $\phi(x)$ is to be used?

$\phi(x)$ must be linear function

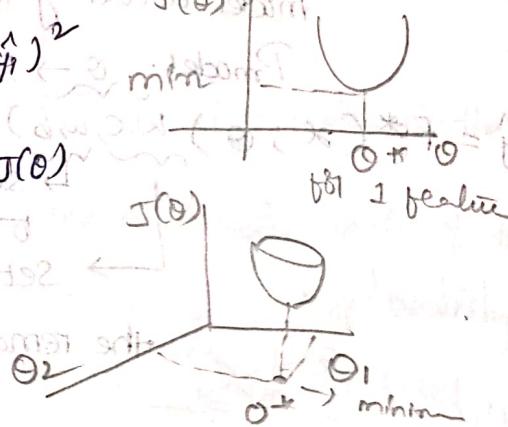
$$\phi(x) = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\rightarrow \hat{y} = \phi(w^T h + b) = w^T h + b$$

① MSE cost function of $J(\theta)$

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$$\theta^* = \text{arg min}_{\theta} J(\theta)$$



Advantages

① Easy to interpret

② Always differentiable because it is a quadratic function

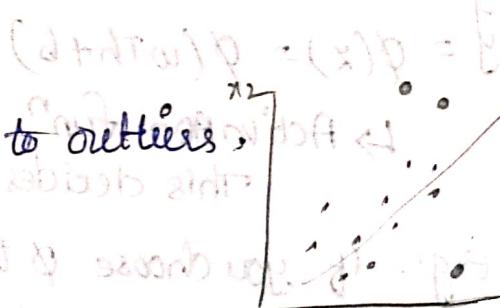
③ Only one local minimum is present

Disadvantages

① MSE is not robust to outliers

outliers bring all the search up

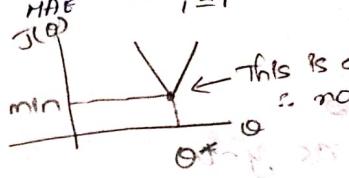
outliers to do a lot of work



② MAE

$$\hat{y} = \phi(w^T h + b) =$$

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$



Advantage

① Intuitive & easy

② Robust to outliers

Disadvantages

① Not differentiable used directly in algorithms.

③ Huber Loss

$$J_{\text{Huber}}(\theta) =$$

δ - hyper parameter

The hyperparameter δ controls the transition to linear.

Note: δ !

① It

② It

①

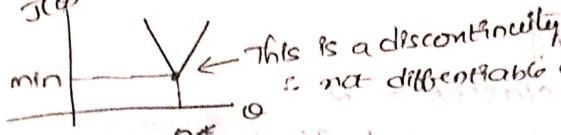
will also need requirement

② MAE

$$\hat{y} = \phi(w^T h + b) = w^T h + b$$

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$

(Error is true)



Advantage

- ① Intuitive & easy to implement
- ② Robust to outliers

Disadvantages

- ① Not differentiable so gradient descent cannot be used directly but we can use Sub gradient descent algorithm.
- ③ Huber Loss

$$J_{\text{Huber}}(\theta) = \begin{cases} \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (y_i - \hat{y}_i)^2 & \text{if } |y_i - \hat{y}_i| \leq \delta \\ \text{linear} & \text{if } |y_i - \hat{y}_i| > \delta \\ \frac{1}{N} \sum_{i=1}^N \delta |y_i - \hat{y}_i| - \frac{1}{2} \delta^2 & \end{cases}$$

δ - hyper parameter

The hyperparameter (δ) defines the point where the transition cost funⁿ $J(\theta)$ transitions from quadratic to linear.

NOTE :- It's Performance

① It's robust to outliers.

② It's Performance lies b/w MAE & MSE

Disadv

- ① Associated Complexity: The hyperparameter (δ) will also need to be optimized which is an additional requirement on the training algorithm.

Note :- For linear regression the OLP layer contains only one OLP unit and with linear activation.

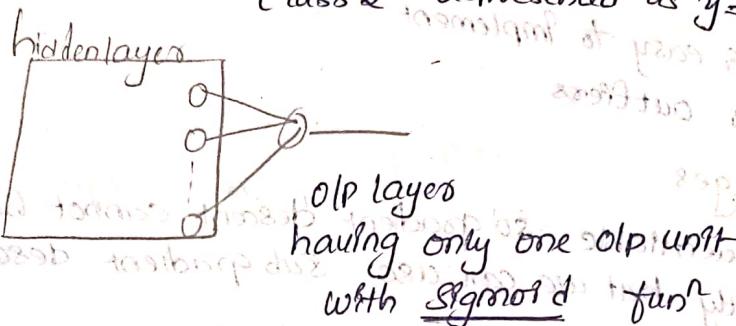
Classification

① Binary-cross entropy

If no. of categories 2 classes = 2

Let us say class 1 represented as $y=0$

Class 2 represented as $y=1$



$$\text{Sigmoid } \phi(x) = \sigma(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{1+e^x}$$

$$z = \mathbf{w}^T \mathbf{x} + b \quad \hat{y} = \sigma(z) = \frac{1}{1+e^{-(\mathbf{w}^T \mathbf{x} + b)}}$$

0.5

So when $\hat{y} < 0.5$ then $y=0$ and when $\hat{y} > 0.5$ then $y=1$

- ① $-\infty < z < \infty$
- ② min value of $\sigma(z) = 0$
- ③ max value of $\sigma(z) = 1$
- ④ If $0 \leq \hat{y} \leq 1 \Rightarrow$ can be interpreted as probability

$\sigma(z)$ & converting the ' \hat{y} ' into probability.

\Rightarrow we know that ' y ' has one of the two possible values

$y=0$ or $y=1$

$y=0$ If $x \in \text{class 1}$

$y=1$ If $x \in \text{class 2}$

\Rightarrow But \hat{y} can take any value
↳ represents

$\Rightarrow \hat{y} = P_{\text{model}}(y=1|x)$
Now \hat{y}

i.e. $\hat{y} = P_{\text{model}}(y=1|x)$

If $P_{\text{model}}(y=1|x) =$

If $P_{\text{model}}(y=0|x) =$

$$J(\hat{y}) = -\sum_{i=1}^N y_i \log$$

Case 1 :- Assume

$y_i = 1$ and

$$\text{Loss} = -1 \cdot \log$$

Case 2 :- $y_i = 0$ and

$$\text{Loss} =$$

Case 3 :- $y_i = 0.5$

respectively

\Rightarrow Then there will be) and adding

$$\text{us } J(\hat{y}).$$

Advantages

The cost function

Properties

① The loss

② This is

\Rightarrow But \hat{y} can take any real no. in b/w 0 & 1 (including)
 \Downarrow represents the probability

$$\Rightarrow \hat{y} = P_{\text{model}}(y=1|x) \quad (\text{or}) \quad \hat{y} = P_{\text{model}}(y=0|x)$$

Now \rightarrow

$$\text{i.e. } \hat{y} = P_{\text{model}}(y=1|x) \quad \boxed{\hat{y}} = P_{\text{model}}(y=1|x)$$

if $P_{\text{model}}(y=1|x) = \hat{y} \geq 0.5$ then x class 2

if $P_{\text{model}}(y=1|x) = \hat{y} < 0.5$ then x class 1

$$J(\theta) = -\sum_{i=1}^N y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i)$$

Case 1 - Assume $y_i = 1$ (100%) LOSS = $-y_i \log \hat{y}_i - (1-y_i) \log (1-\hat{y}_i)$

$y_i = 1$ and the model prediction also $\hat{y}_i = 1$

$$\text{Loss} = -1 \cdot \log 1 - (1-1) \log (1-1) = 0$$

Case 2 - $y_i = 0$ and $\hat{y}_i = 0$ LOSS = $0 \cdot \log 0 + (1-0) \log (1-0) = 0$

Case 3 - \hat{y}_i is other than 0 or 1 when $y_i \neq 0, 1$ respectively. Example: $\hat{y}_i = 0.21, 0.47, 0.78, 0.68, \dots$

\Rightarrow Then there will be non-zero loss (that leads to be) and adding loss corresponding to all training samples

is $J(\theta)$.

Advantages

The cost function is

differentiable

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

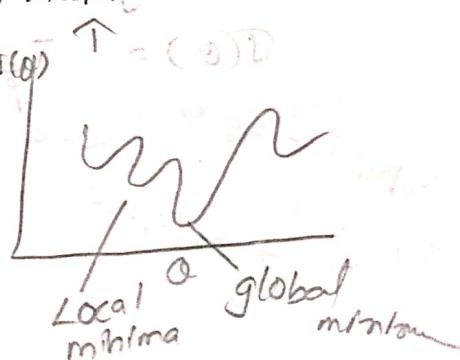
$\frac{\partial J(\theta)}{\partial \theta_0}$

Disadvantage

① The loss landscape will have multiple local minima.

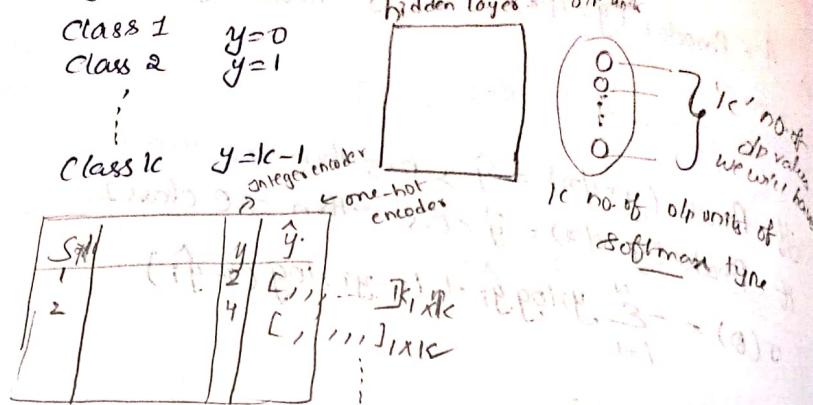
② This is not intuitive

③ This is not differentiable at $\hat{y}_i = 0, 1$



② Multiclass classification 1088

Categories = # classes = k



$$\text{Softmax} \quad \phi(\tilde{z}_i) = \frac{\phi(z_i)}{\sum_{j=1}^k \phi(z_j)}$$

e.g. Iris dataset

$k=3$ integer encoder One-hot Encoder

Setosa	0	[1, 0, 0]
Versicolor	1	[0, 1, 0]
Virginica	2	[0, 0, 1]

\tilde{z}_i	y_i	loss	Model accuracy
1	1	0.475	after 10 epochs
2	0	0.155	$\frac{\sum \hat{y}_i}{N} = 0.938$
3	0	0.070	$\frac{\sum \hat{y}_i}{N} = 0.925$

$$\text{Loss} = -\sum_{j=1}^k y_j \log \hat{y}_j$$

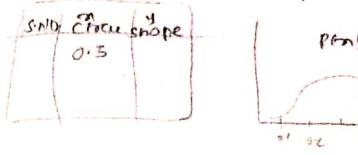
$$T(O) = -\sum_{i=1}^N \sum_{j=1}^k y_{ij} \log \hat{y}_{ij}$$

MNIST - Digit recognition
Image-net
 $\rightarrow k=10,000$

Sequential fit
model compile predict
Summary evaluate
Plot-model

2/1/23

$$\text{Circularity} = \frac{4\pi A}{P^2}$$



$P(m)$ = probability distribution

$$P(y = \text{Triangle}) = 0.5$$

$$P(y = \text{Rect}) = 0.5$$

$$P(y = T|x) = \frac{1}{2}$$

$$P(y = R|x) = \frac{1}{2}$$

$P(\alpha|y=T) \rightarrow P(m)$ given
 $P(\alpha|y=R) \rightarrow$ Conditioned on

$$P(y=T|x) = P(m)$$

$$[P(y=T|x)P(x)] = P(m)$$

$$P(y=R|x) = P(m)$$

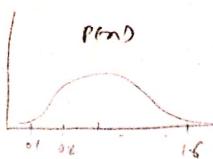
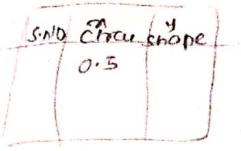
$$P(y=T|x) > P(y=R|x)$$

$$\Rightarrow P(\alpha|y=T) P(y=T|x) > P(\alpha|y=R) P(y=R|x)$$

$$P(\alpha|y=)$$

2/1/23

$$\text{Circularity} = \frac{4\pi A}{P^2} \quad \begin{matrix} \text{Suppose} \\ \downarrow \end{matrix} \quad \begin{matrix} A = \text{Area} \\ P = \text{Perimeter} \end{matrix}$$



$P(m)$ = probability distribution of actual measurements

$$P(y=T|\text{angle}) = 0.5 \quad \begin{matrix} \text{a priori} \\ \text{probabilities} \end{matrix}$$

$$P(y=R|\text{rect}) = 0.5 \quad \begin{matrix} \text{a priori} \\ \text{probabilities} \end{matrix}$$

$$P(y=T|x) = \quad \begin{matrix} \text{a posteriori} \\ \text{probabilities} \end{matrix}$$

$$P(y=R|x) = \quad \begin{matrix} \text{a posteriori} \\ \text{probabilities} \end{matrix}$$

↳ we want these probabilities

↳ But very different to get

$$P(x|y=T) \quad \begin{matrix} \rightarrow p(m) \text{ given that } y=T \\ \int_{m=y}^{M} \end{matrix}$$

$$P(x|y=R) \quad \begin{matrix} \rightarrow \text{Conditional} \\ \text{probability distribution} \end{matrix}$$

$$P(y=T|x) = \frac{P(x|y=T) P(y=T)}{P(x)} \quad [\text{Bayes theorem}]$$

$$[P(y=T|x) P(x) = P(x|y=T) P(y=T)]$$

$$P(y=R|x) = \frac{P(x|y=R) P(y=R)}{P(x)}$$

$$P(y=T|x) > P(y=R|x) \Rightarrow \text{Triangle}$$

$$< \Rightarrow \text{Rectangle}$$

$$\Rightarrow P(x|y=T) P(y=T) > P(x|y=R) P(y=R) \quad (\text{if same})$$

$$P(x|y=T) > P(x|y=R) \Rightarrow \text{Triangle}$$

$$< \Rightarrow \text{Rectangle}$$