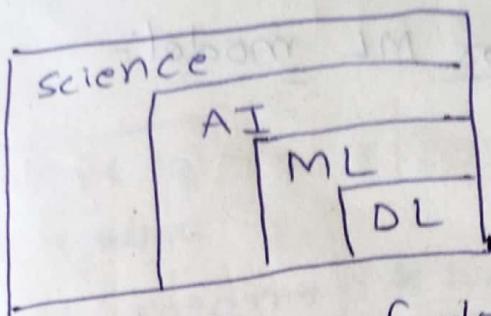


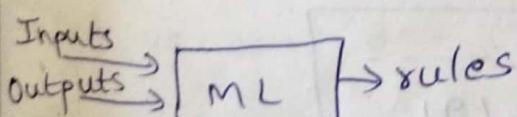
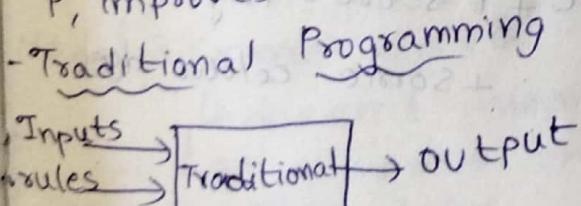
29/12/22

Deep Learning



- AI is one way of doing science.
- ML is one way of doing AI
- DL is one way of doing ML
- Deep learning - learning using neural networks
- ML - learning using other than neural networks
(except neural networks)
Ex: SVM, KNN, DTs, RFs etc..
- AI - It is not just learning, but also doing reasoning, planning and solving problems.
- Science - doing science we develop intelligence
- Science - doing science we learn from and knowledge.

ML: A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E



- Tasks
- ① Sorting & searching (we can describe these tasks using clear set of rules.)
 - ② Playing chess (Traditional algorithms)
 - ③ Recognizing speaker listening to voice

④ Understanding what is being spoken

⑤ Recognizing objects.

- Algorithmic components for ML model:-

- Dataset

- Model

- Optimization Algorithm

- Cost function

- Training strategies.

- Understanding Feature space:-

Q. Define input data of a typical ML algorithm

Ans Input data is a collection of features

that have been quantitatively measured from some object or event that we want the ML system to process.

Linear Regression

Task To prepare a smoothie and serve

Goal To predict the calories

Experience ?

Performance ?

Smoothie = 1 cup of yogurt + 1 secret ingredient
+ some colour dyes

Calories

300

400

15 smoothies

265, 281, 207, 148, 235, 72
-----, 191

smoothie

calories

268

281

207

148

235

72

$$\text{Mean} = 236.9$$

$$\text{Median} = 240$$

	calories	smoothie	weight	calories
1	79	79	26.5	26.5
2	111	111	22.1	22.1
3	115	115	15	15
4	141	141	16	16

weight of
the ingredient

$$\text{calories} = \text{mean} + p * \text{weight}$$

\hookrightarrow Intercept \hookrightarrow slope

Task or Goal - Predict calories

Model 1 :- Wild guess

\hookrightarrow very bad model

\hookrightarrow we are not having any experience

Model 2 :- we use experience

if I previously made smoothies

26.8, -----, 111 \Rightarrow mean = 236.9

26.8, -----, 111 \Rightarrow median = 240

\hookrightarrow No features are used

\hookrightarrow Prediction for 16th smoothie is 238

but actual calories = 445

\hookrightarrow ingredient quich of 88 gm

\hookrightarrow under estimation by 20%

Model 3 :- more data

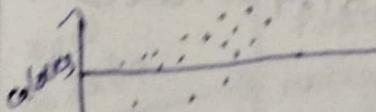
\hookrightarrow doesn't mean moderate points

\hookrightarrow It means more information about the data.

- Let us consider weight of ingredient

Data points

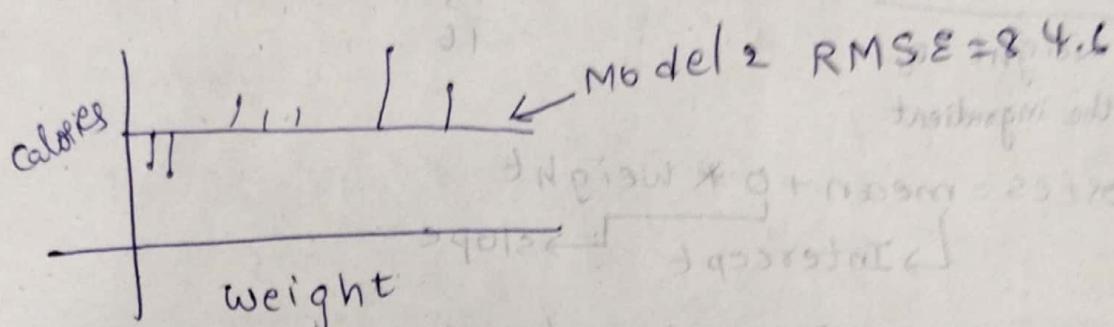
weights	calories
26.5	26.5
22.1	22.1
15	15
16	445



\downarrow weight
represents
features
 \downarrow
represents
output

Model 2 calorie $c = 236.9 + 0 * \text{weight}$

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



for model 2

$$\text{Model 3: } c = 236.9 + 0.5 * \text{weight}$$

$$\text{Model 4: } c = 236.9 + 1 * \text{weight}$$

$$\text{Model 5: } c = 236.9 + 200 + 1 * \text{weight}$$

$$\text{RMS E} = 69.3$$

$$= 69.2$$

$$= 52.8$$

$$\text{Finally: } c = 175 + 1.35 * \text{weight}$$

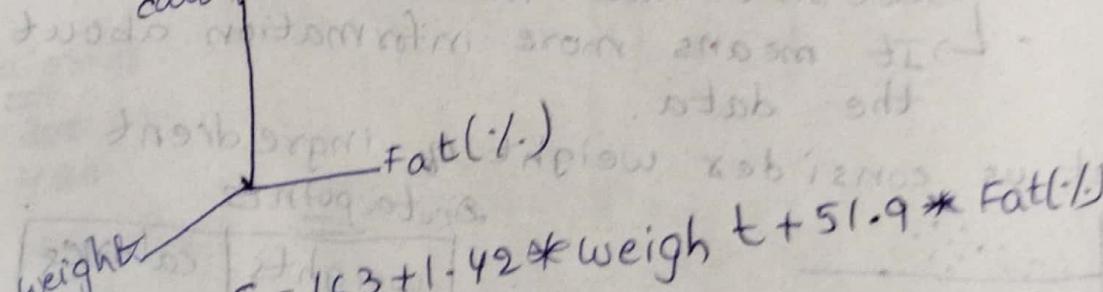
$$\text{RMS E} = 48.8$$

To make it much better, Model 6:

For 16th smoothie, the prediction = 293.8

↳ underestimate by $445 - 293.8 = 151.2$

Model 7: $c = w_0 + w_1 * \text{weight} + w_2 * \text{Fat}(\%)$



$$\text{RMS E} = 46.9$$

For 16th smoothie pred underestimated by
an amount 143

So, consider $Fat(\gamma)$, Carbohydrates, Protein
(9-4-4)

model 1: $c = w_0 + w_1 * Fat(\gamma) + w_2 * Carbohydrates + w_3 * Protein$

$$c = 143 + 8.9 * Fat + 3.9 * Carbohydrates + 4.3 * Protein$$

RMSE < 10

For smooth pred = 431

Classification - Computer Vision

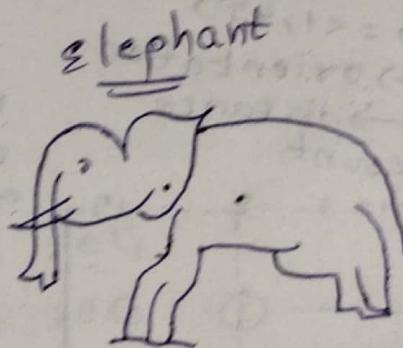
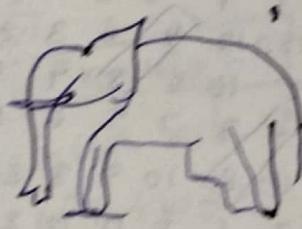
Task 1: Lion or elephant



Shape features are sufficient

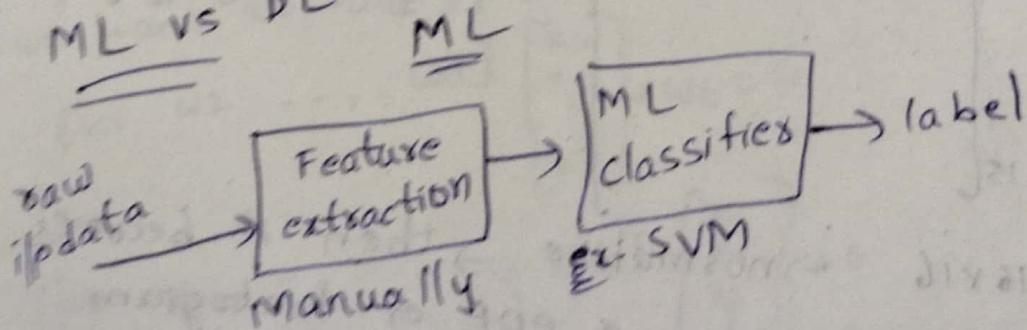
Task 2: Horse vs Zebra

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

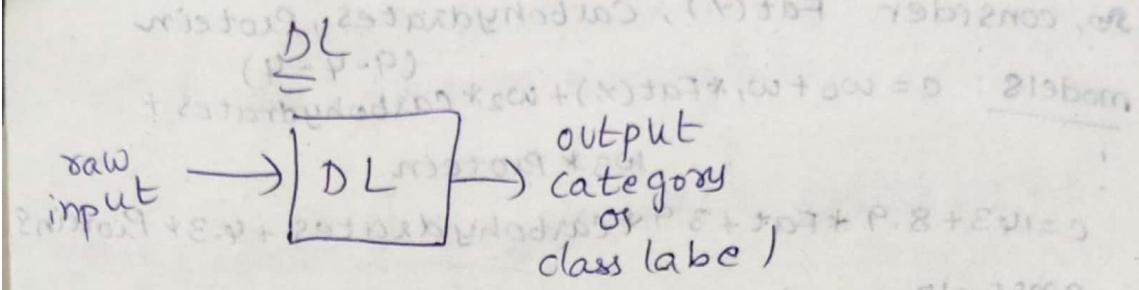


elephant

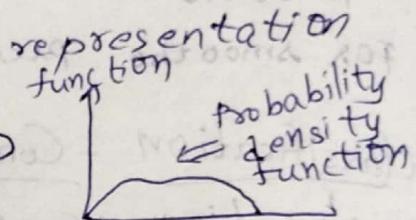
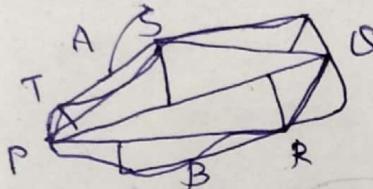
ML vs DL



Model 2 calorie = $236.9 + 0 * \text{weight}$



Shape descriptors: polygonal



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

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

$$\text{skewness} =$$

Regional features:

① Texture

Texture can be represented using the cooccurrence matrix $A_p(i,j)$

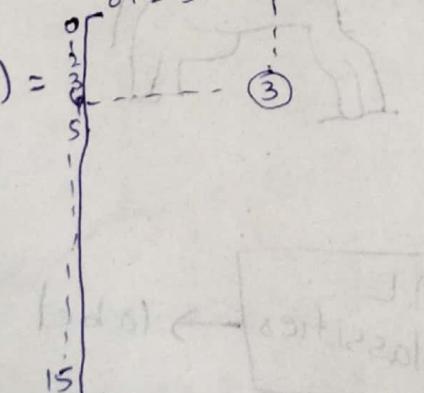
4-bit image of size 8×8

$$i=4; j=8$$

$P = \langle l, \theta \rangle = \langle 1, 45^\circ \rangle$

↳ orientation
↳ distance

4	X
X	8



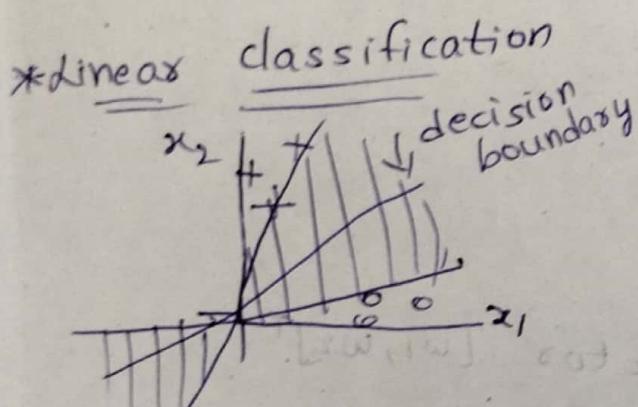
Raw image

10	9	7	9	5	8	11	9
65	5	15	12	4	6	3	2
9	3	2	10	6	8	4	5
8	2	4	3	7	5	6	1
8	0	11	8	10	9	8	2
8	7	1	6	0	7	2	
2	4	8	9	11	6	3	9
2	3	8	8	6	12	6	7
7	2	8	8	6	12	6	7

$A_p(i,j) 16 \times 16$
given
 $\langle l, 45^\circ \rangle$

\rightarrow normalize \rightarrow then it becomes a PDF or histogram denoted as $C_{ij} 16 \times 16$

- ① Max prob: $\max_{i,j} c_{ij}$
- ② Elementwise difference moment : $\sum_i \sum_j (i-j)^k c_{ij}$
- ③ Inverse elementwise diff : $\sum_i \sum_j c_{ij} / (i-j)^k$
- ④ Uniformity : $\sum_i \sum_j c_{ij}^2$
- ⑤ Entropy : $-\sum_i \sum_j c_{ij} \log(c_{ij})$



decision boundary $w_1 x_1 + w_2 x_2 + w_0 = 0$
↳ bias

Decision boundary is the normal of the vector
 $[w_1, w_2]^T$

$$\sum_{i=1}^d w_i x_i + w_0 = 0$$

$$[w^T x + w_0 = 0]$$

Decision rule is

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

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

becomes

program

$\underbrace{[w_1 \ w_2 \ \dots \ w_d]}_{\text{on}} \underbrace{w_0}_{y} = 0$

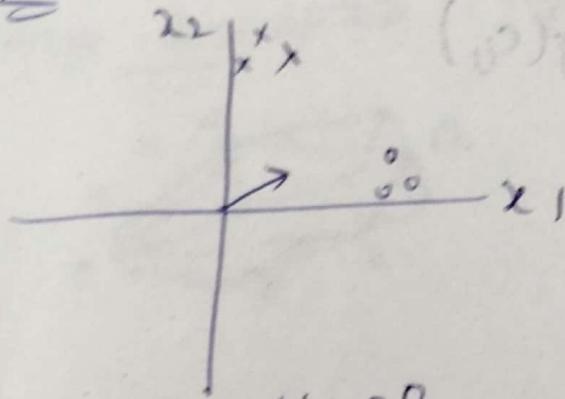
Decision rules

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

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

Now, negate y of class 2

$$\underline{DR} = a^T y > 0$$



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

Decision boundary

↳ is normal to the vector $[w_1, w_2]^T$

↳ Position is determined by w_0

AI/ML

→ Input data

→ Feature space

→ Model Architecture

→ Training strategy

→ Optimization algorithm

→ Loss function (RMSE)

Finally performance on test set.

Regression → Fitting the hyperplane

Classification → Finding the decision boundary

which separates two classes

Regression

$$\text{calories} = w_0 + w_1 * \text{fat(1)} + w_2 * \text{carbo} + w_3 * \text{protein}$$

$$\text{prediction} = w_0 + w_1 * \text{feat1} + w_2 * \text{feat2} + w_3 * \text{features}$$

classification

decision boundary
is also a hyperplane

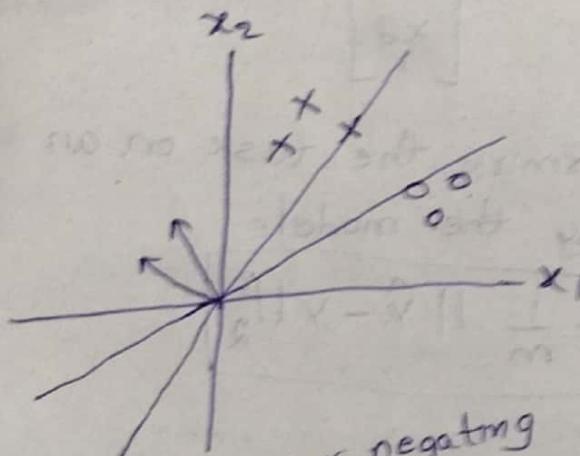
$$w_1 * f_1 + w_2 * f_2 + w_0 = 0$$

Decision rule

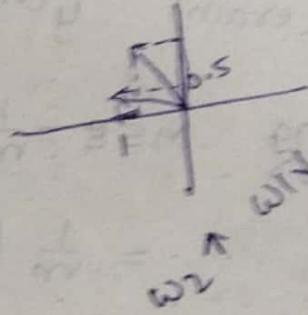
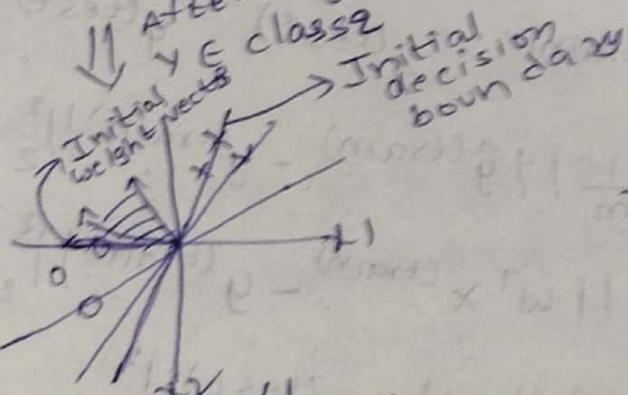
$$a^T y > 0$$

$$a = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_0 \end{bmatrix}$$

$$y = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ 1 \end{bmatrix}$$



↓ After negating
Initial weight vector



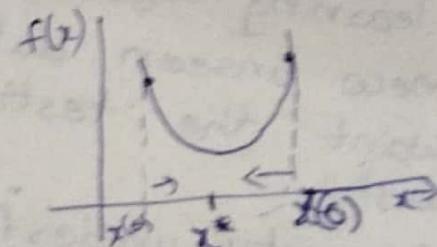
Training Algorithm

1) $a(0) \leftarrow$ Arbitrarily initialize

2) $a(k) \leftarrow a(k-1)$ in such a way that
 $a(k)$ minimizes some
performance index.

2(b) Error $J = -\sum_{Y \text{ misclassified}} a^T y$

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



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

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

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

$$a(k) = a(k-1) + \eta \sum_{Y \text{ misclassified}} y$$

* Inference:-

Regression $w_1x_1 + w_2x_2 + w_0 = 0$

Classification $\hat{y} = 0$
 $w^T x = 0$

where $W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \\ w_0 \end{bmatrix}$ $(d+1) \times 1$

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \\ 1 \end{bmatrix}$$

- Inference means performing the task on an unseen data so far by the model.

Training error

$$MSE = \frac{1}{m} \sum_{i=1}^m ||\hat{y}_i - y_i||^2$$

m samples
 $x^{(train)}$ $y^{(train)}$ $x^{(test)}$ $y^{(test)}$

$$\text{Training error } MSE = \frac{1}{m} \sum_{i=1}^m ||\hat{y}^{(train)} - y^{(train)}||^2$$

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

$$\text{Test error} = ||\hat{y}^{(test)} - y^{(test)}||^2$$

↳ Generalization error

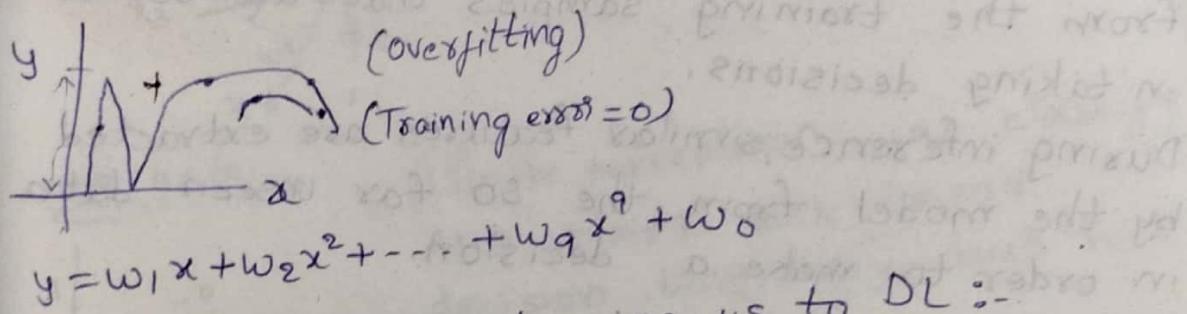
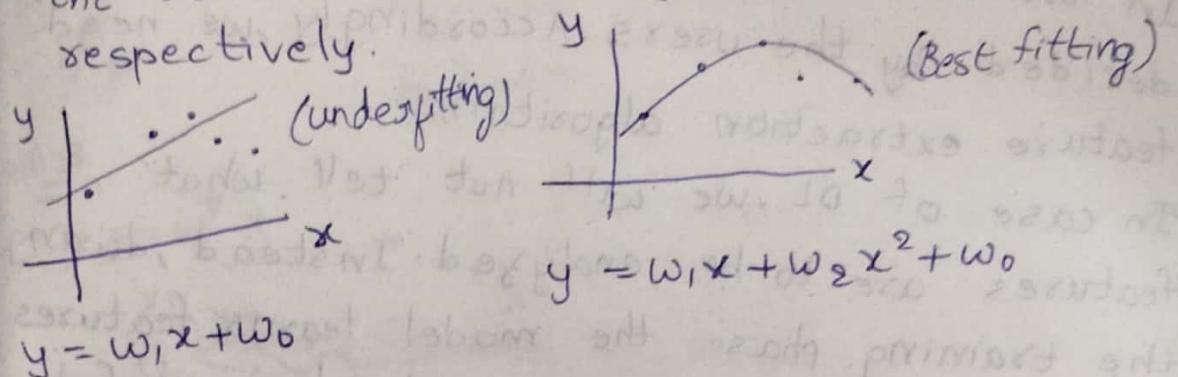
* Generalization:- It is the ability of machine learning systems to perform well on the new unseen data. That means to say we want the testing error also too small as well.

R - Capacity, underfitting and overfitting:-

$$pred = w_1x_1 + w_2x_2 + w_0$$

Different weight vector values gives us different functions.

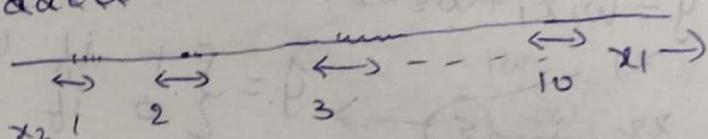
Range of functions which the model can learn is called the capacity of that model. The goal then becomes selecting one among those functions which will best fit the data (or) which will well-separate the data corresponding to two categories respectively.



- challenges of ML driving us to DL :-

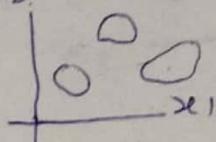
① The curse of data dimensionality

one feature



+ second feature

$$10^2 = 100$$

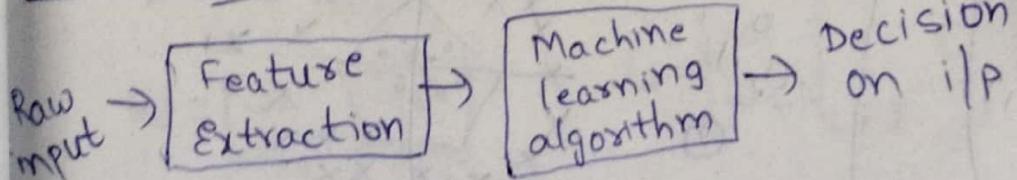


+ third feature

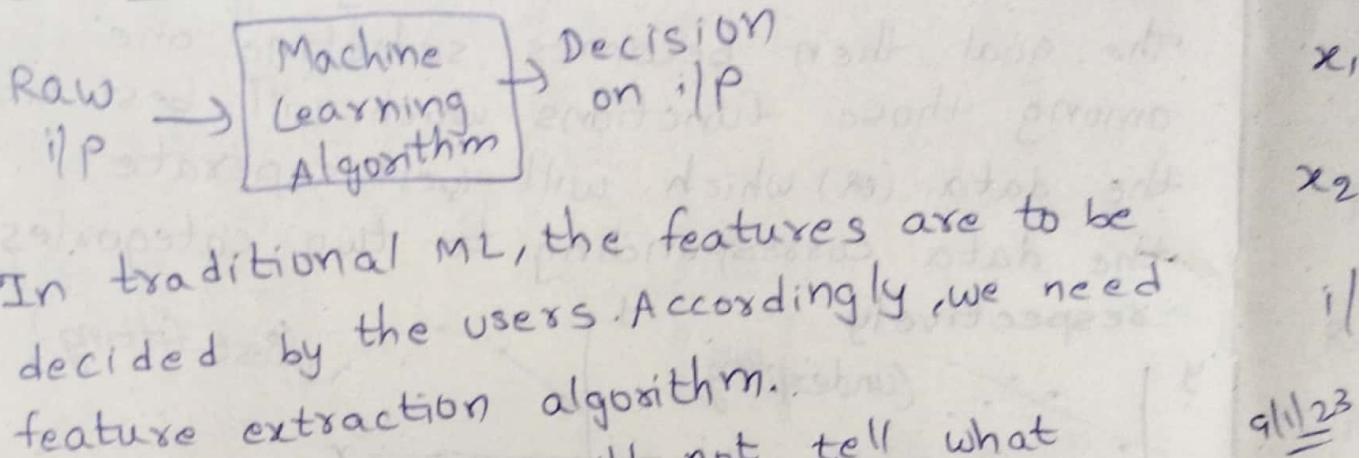
② local constancy and smoothness regularization

③ Manifold learning.

ML Approach



DL Approach



- In traditional ML, the features are to be decided by the users. Accordingly, we need feature extraction algorithm.
- In case of DL, 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.

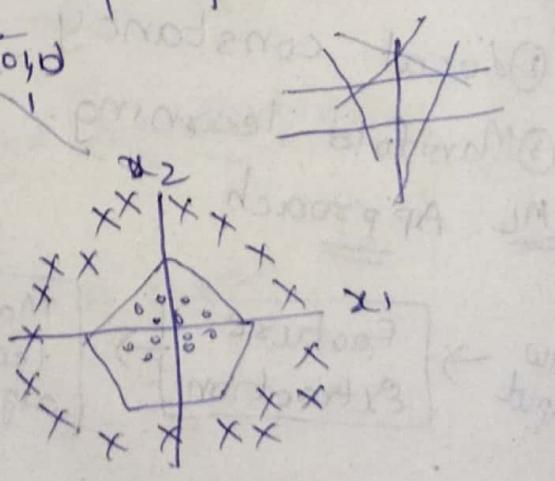
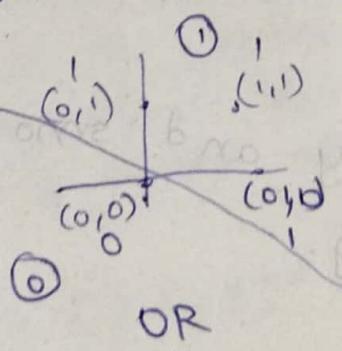
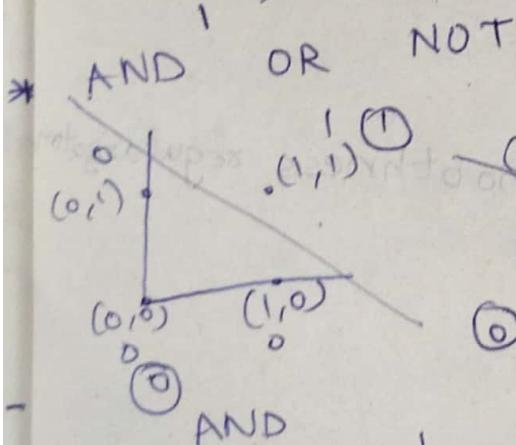
*Perceptron:-

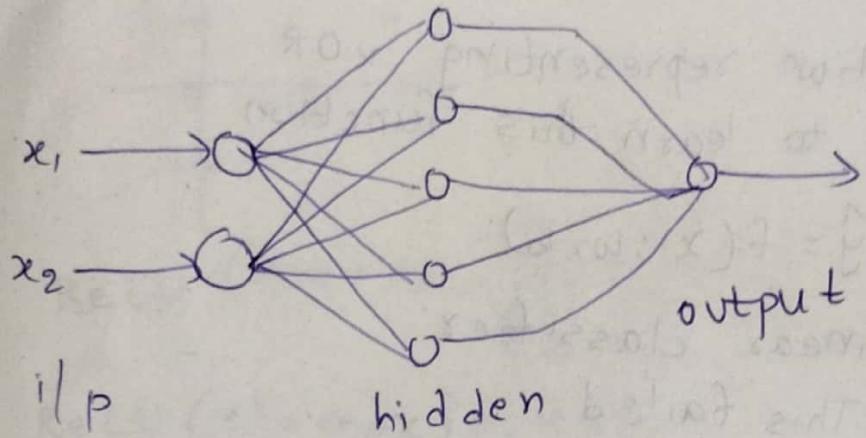
$$y = w_1x_1 + w_2x_2 + w_0$$

$$\begin{matrix} x_1 \\ x_2 \end{matrix} \rightarrow \sum_{i=1}^2 w_i x_i + w_0$$

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

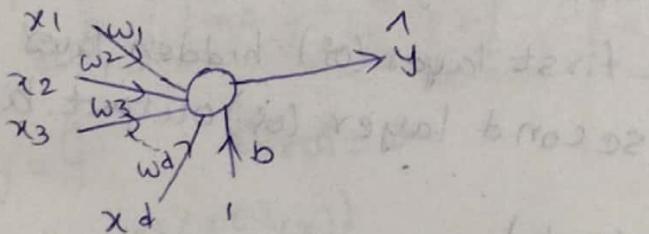
		AND	OR	NOT
x ₁	x ₂	0	0	1
0	0	0	1	1
0	1	0	1	0
1	0	1	1	0
1	1	1	1	0



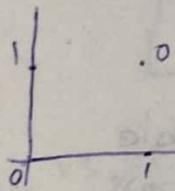


g1/123

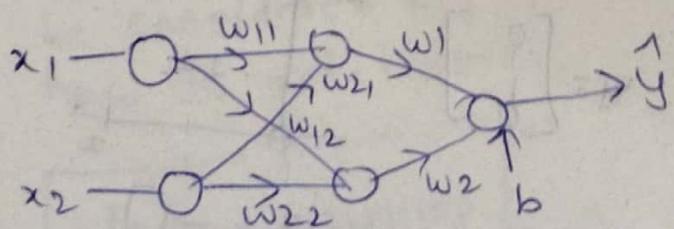
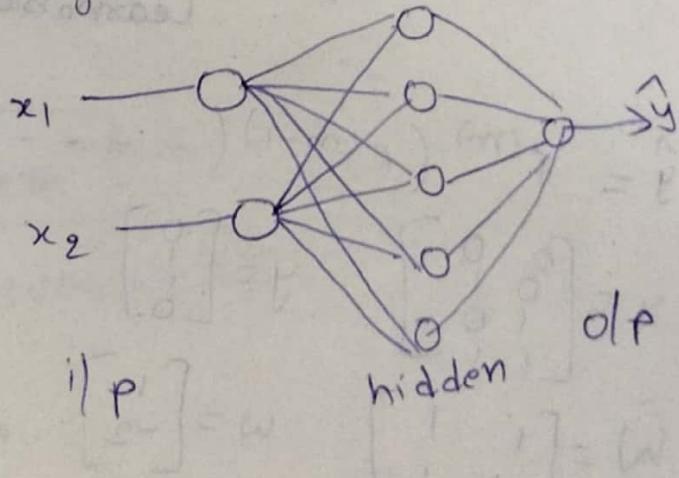
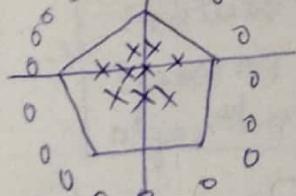
$$y = w^T x$$



XOR gate



Single perceptron can classify AND, OR, NOT with linear decision boundary but not XOR gate.



$$W = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix}_{2 \times 2}$$

$$w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}_{2 \times 1}$$

$$C = \begin{bmatrix} c_1 \\ c_2 \end{bmatrix}_{2 \times 1}$$

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

The function representing $x \text{ OR}$
we want to learn this function

we tried $\hat{y} = f(x; w, b)$

Linear classifier

This failed

$$\hat{y} = f^{(2)}(f^{(1)}(x)) = f(x)$$

Input layer

first layer (or) hidden layer

second layer (or) output layer

$$\hat{y} = f(x; w, c, w, b)$$

θ (or)

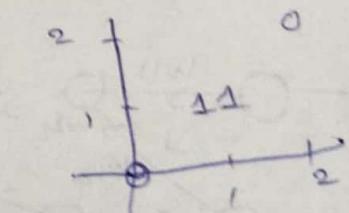
Θ
parameters (Θ)

learnable (Θ) trainable
parameters

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

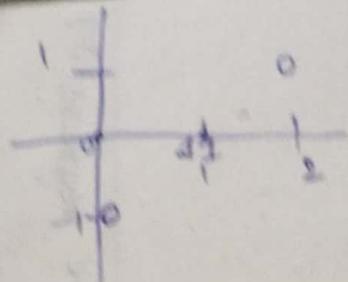
$$x = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \quad y = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$$w = \begin{bmatrix} 1 & 1 \end{bmatrix} \quad b = 0.$$



$$xw = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} \quad u = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

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

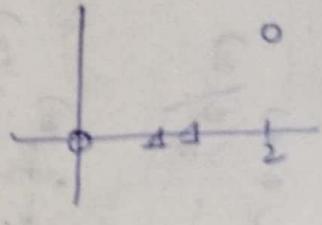


$$\text{ReLU}(xw + c) = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$$

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

$$\text{ReLU}(xw + c) \cdot w + b$$

$$= \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -2 \end{bmatrix} + \begin{bmatrix} 0 \\ 8 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \vec{y}$$

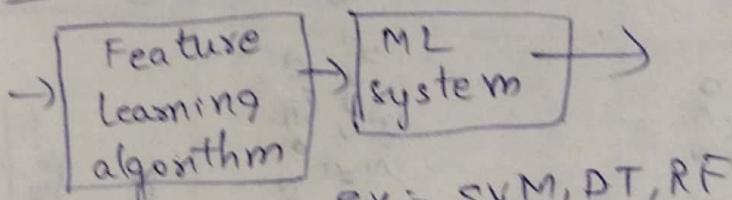


yes

$$y = \vec{y} \quad \text{error} = 0$$

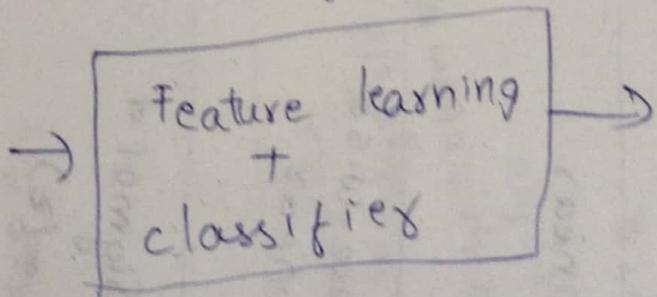
$$f(x) = \begin{cases} f^{(2)}(x) & f^{(1)}(x) \\ \text{feature learning layer} \\ \text{Linear classifier} \end{cases}$$

ML vs DL



ex: SVM, DT, RF

↳ Designer has to decide which features to be considered & how to get them.



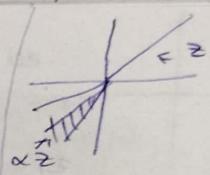
Types of Activation functions:-				
Activation function	Equation	Derivation	ID Graphs	Examples
① Binary step (or) unit step (or) Heaviside	$\phi(z) = \begin{cases} 0 & z < 0 \\ 1 & z \geq 0 \end{cases}$	$\phi'(z) = 0$		MLP
② signum	$\phi(z) = \begin{cases} -1 & z < 0 \\ 0 & z = 0 \\ 1 & z > 0 \end{cases}$	$\phi'(z) = 0$		MLP
③ linear	$\phi(z) = z$	$\phi'(z) = 1$		
④ Piece-wise linear	$\phi(z) = \begin{cases} 1 & z \geq 1/2 \\ z + 1/2 & -1/2 < z < 1/2 \\ 0 & z \leq -1/2 \end{cases}$	$\phi'(z) = \begin{cases} 1 & -1/2 < z < 1/2 \\ 0 & \text{otherwise} \end{cases}$		
⑤ Sigmoid or(z)	$\phi(z) = \frac{1}{1+e^{-z}} = \frac{e^z}{1+e^z}$	$\phi'(z) = \sigma(z) [1 - \sigma(z)]$		logistic regression Binary-classification
⑥ tanh	$\phi(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$	$\phi'(z) = 1 - \tanh^2(z)$		
⑦ Soft plus Note :- Sigmoid = differentiation of soft plus	$\phi(z) = \log(1+e^z)$	$\phi'(z) = \frac{e^z}{1+e^z}$		
⑧ ReLU	$\phi(z) = \max(0, z)$ $\phi(z) = \begin{cases} z & ; z \geq 0 \\ 0 & ; z < 0 \end{cases}$	$\phi'(z) = \begin{cases} 1 & ; z \geq 0 \\ 0 & ; z < 0 \end{cases}$		RNN, CNN, Hidden layers
⑨ Leaky ReLU	$\phi(z) = \max(0, \alpha z, z)$ $\phi(z) = \begin{cases} z & ; z \geq 0 \\ 0.1z & ; z \leq 0 \end{cases}$	$\phi'(z) = \begin{cases} 1 & ; z \geq 0 \\ 0.1 & ; z < 0 \end{cases}$		

⑩ Parametric ReLU

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

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

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



⑪ Softmax

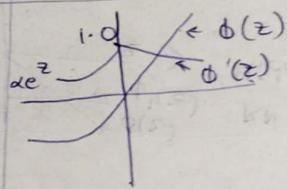
$$\phi(z_i) = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

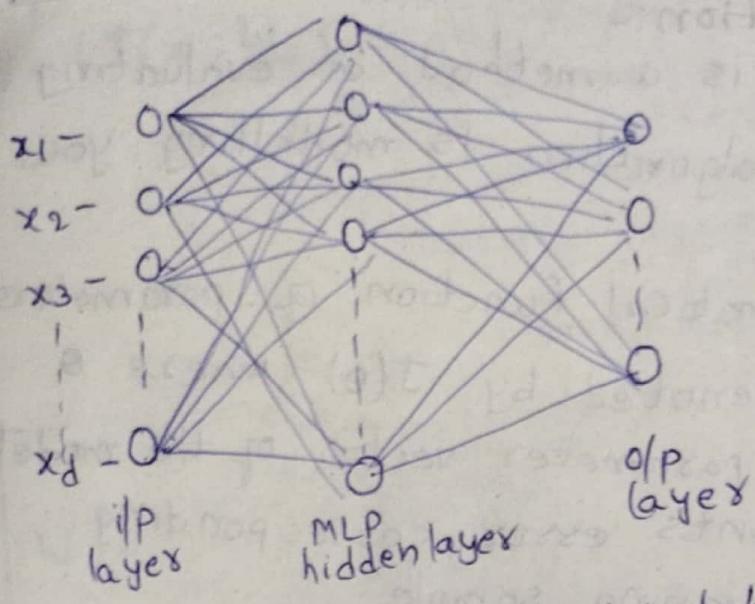
Multiclass classification

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 layers :- It can see labels/targets i.e. input data it can see.

During training phase, parameters are updated such that the dp becomes close to targets.

Hidden layers :- they don't have access to labels or targets i.e., training datasets.

obviously, the hidden layers won't get any direct guidance from dataset in updating its parameters.

↳ Solution is backpropagation.

→ Thus, they are called hidden layers.

For output layers, the activation function is decided chosen based on type of tasks. The tasks can be

① Regression - Linear Activation function

② Binary classification - sigmoid activation function
(gives conditional probability)

③ Multi class classification - softmax

④ Multi label classification - sigmoid

Hidden layers - The activation functions are chosen based on the type of neural network architecture

① CNN - ReLU

② RNN - Tanh (or) sigmoid

Cost / Loss function:-

① Cost function is a method of evaluating how well your algorithm is modelling your dataset.

• It is a mathematical function of parameters of the model denoted by $J(\theta)$ where θ represents the parameter vector of the model.

• Loss \rightarrow represents error corresponding one training sample

cost \rightarrow represents error over the entire dataset.

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 loss

GAN :-

① Discriminative loss

② Minmax GAN loss

Object Detection :-

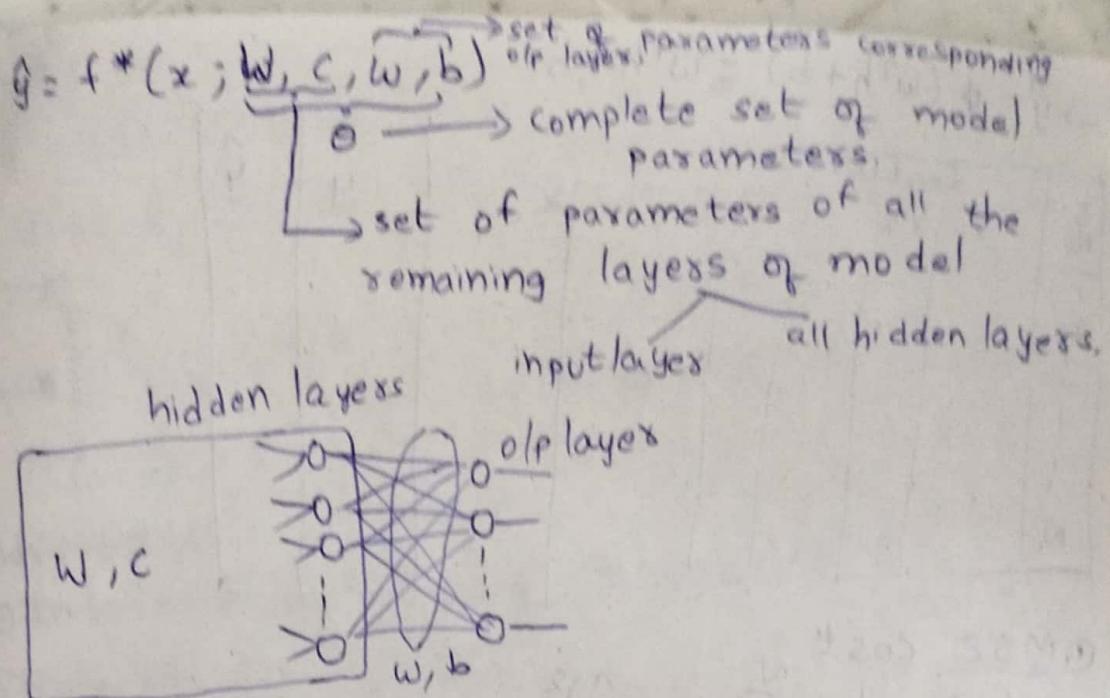
① Focal loss

Word embeddings :-

① Triplet loss

Problem setting

$x, y \sim P_{\text{data}}$ \Rightarrow the entire dataset coming from the distribution P_{data}
 $\hat{y} \sim P_{\text{model}}$ \Rightarrow predictions produced by the model (\hat{y}) according its own learned distribution P_{model}



$z = w^T h + b$
 feature vector from last hidden layer
 linear output of the output layer.

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

Activation function
 decides what is the type of units.
 e.g. if you choose ϕ to be sigmoid, then output unit is said to be of sigmoid type.

Problem setting
 We want \hat{y}_i as close to y_i as possible \rightarrow target.

for all training samples.

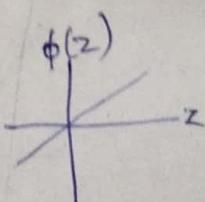
→ If they are not equal, then there is error given by $y_i - \hat{y}_i$

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

the Linear Regression Loss:

→ Which $\phi(z)$ is to be used?

→ $\phi(z)$ must be linear function $\phi(z) = z$

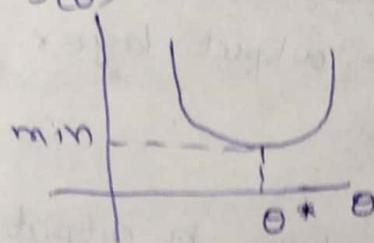


$$\rightarrow \hat{y} = \phi(w^T h + b) = w^T h + b$$

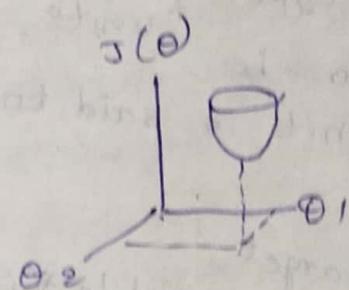
S. No	features					x_d	y	\hat{y}
1	x_1	x_2						
2								
3								
4								
5								
N								

① MSE cost

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



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



Advantages:-

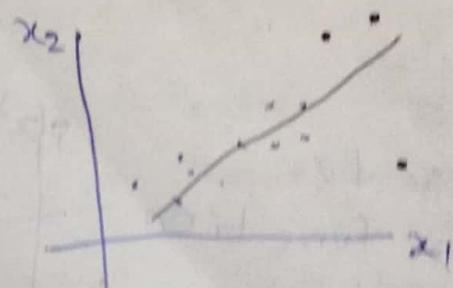
① Easy to interpret.

② Always differentiable, because it is a quadratic function.

③ Only one local minima present

Disadvantages:-

① MSE is not robust to outliers

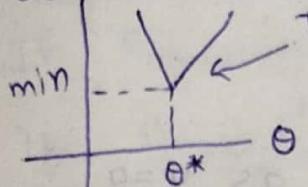


② MAE

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

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

$J(\theta)$



this is a discontinuity

∴ not differentiable at that point

Advantages:-

① Intuitive and easy to implement.

② Robust to outliers

Disadvantage:-

① Not differentiable. So gradient descent cannot be used directly but we can use subgradient descent algorithm.

③ Huber loss

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

s-hyper parameters

the hyper parameter s defines the point where the cost function $J(\theta)$ transitions from quadratic to linear.

Advantages:-

Note:- ~~Huber~~ performance

① Robust to outliers

② Its performance lies b/w MAE & MSE

Disadvantage:-

① Associated complexity. The hyper parameter s which also needed to be optimized which is an additional requirement on the training algorithm.

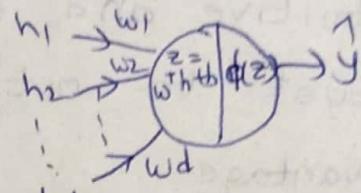
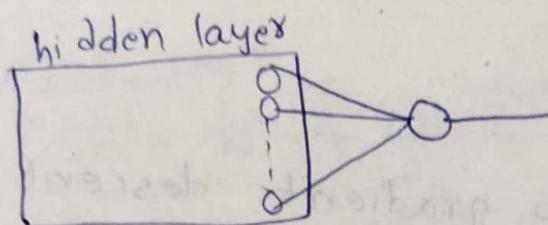
Note:- For linear regression, the output layer contains only one output unit with linear activation.

Classification losses:-

① Binary cross-entropy

#categories (OR) #classes = 2

Let us say class 1 represented as $y=0$
class 2 represented as $y=1$



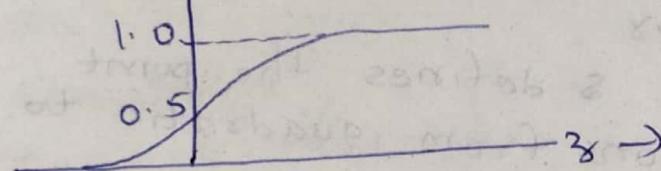
dp layers having only one o/p unit with sigmoid function

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

$$z = w^T h + b$$

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

$$\phi(z)$$



- ① $-\infty < z < \infty$
- ② Min value of $\sigma(z) = 0$
- ③ Max value of $\sigma(z) = 1$
- ④ Max value of $\sigma(z) = 1$
- ⑤ $\therefore 0 \leq \hat{y} \leq 1 \Rightarrow$ can be interpreted as probability.
- $\sigma(z)$ is converting ' z ' into probabilities.
- 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 $x \in \text{class 2}$
- But \hat{y} can take any real number in between

0 and 1 $\Rightarrow \hat{y}$ is a probability

$$\hat{y} = P_{\text{model}}(y=1/x)$$

$$\hat{y} = P_{\text{model}}(y=0/x) \times$$

If $P_{\text{model}}(y=1/x) = \hat{y} > 0.5$ then $x \in \text{class 2}$

If $P_{\text{model}}(y=1/x) = \hat{y} < 0.5$ then $x \in \text{class 1}$

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

S.No		y	\hat{y}
1			
2			
3			
N			

Suppose the accuracy = 1
 $\text{Loss} = -y_i \log \hat{y}_i - (1-y_i) \log (1-\hat{y}_i)$

Case I :- $y_i = 1$ & the model prediction also $\hat{y}_i = 1$

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

Case II :- $y_i = 0$ and $\hat{y}_i = 0$

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

Case III :- \hat{y}_i is other than 0 (or) 1

when y_i is 0 (or) 1 respectively.

$$\text{Eg } \hat{y}_i = 0.21, 0.77, 0.98$$

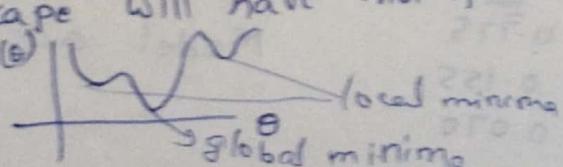
→ then there will be non-zero loss, and adding losses corresponding to all training samples gives us $J(\theta)$

Advantages:-

① The cost function is differentiable

Disadvantages:-

② The loss landscape will have multiple local minima



② It is not intuitive

③ Multiclass classification loss

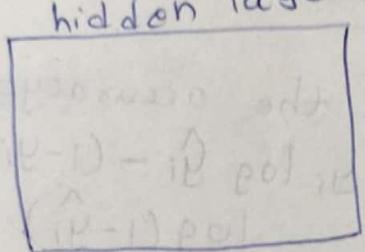
#categories = #classes = k

class 1 $y=0$

class 2 $y=1$

class k $y=k-1$

hidden layers



k number of o/p units
of softmax type

{ 'k' number of output values

S.No		y	\hat{y}
1		2	$[0, 1, 0]_{1 \times k}$
2		4	$[0, 0, 1]_{1 \times k}$
3		1	$[1, 0, 0]_{1 \times k}$
N	$\theta = (0-1)$	$\theta = (0-1)$	$\theta = (0-1)$

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

Ex Iris dataset

$k=3$ one-hot encoding

setosa 0 [1, 0, 0]

virginica 1 [0, 1, 0]

versicolor 2 [0, 0, 1]

$$\frac{y_i}{\hat{y}_i} \rightarrow \text{Loss} = 0.3675$$

$$\begin{array}{cc} 1 & 0.775 \\ 1 & 0.155 \\ 2 & 0.070 \end{array} \quad \text{Model accuracy} = 63.25\%$$

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

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

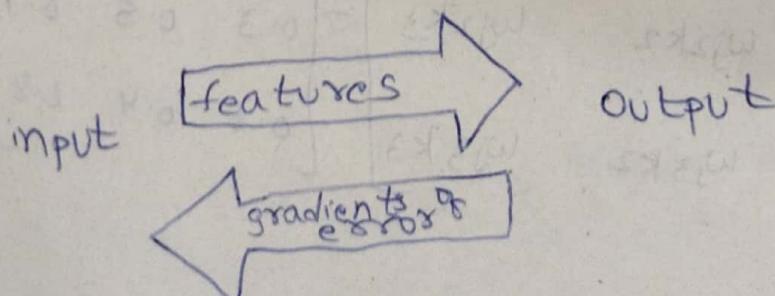
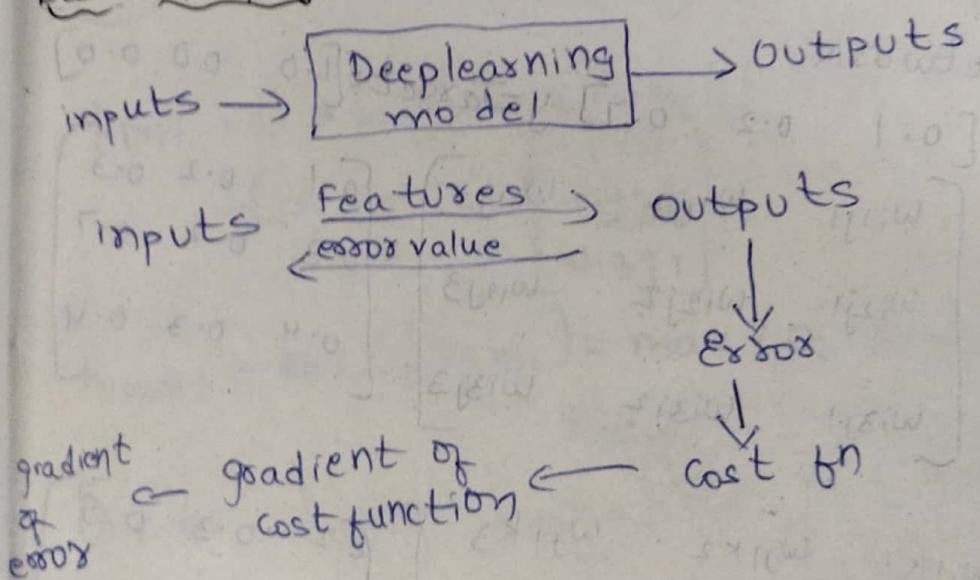
After few epochs

$\frac{g_1}{0.938}$	$\rightarrow \text{loss} = 0.095$
0.028	$\rightarrow \text{Accuracy} = 99.05\%$
0.036	

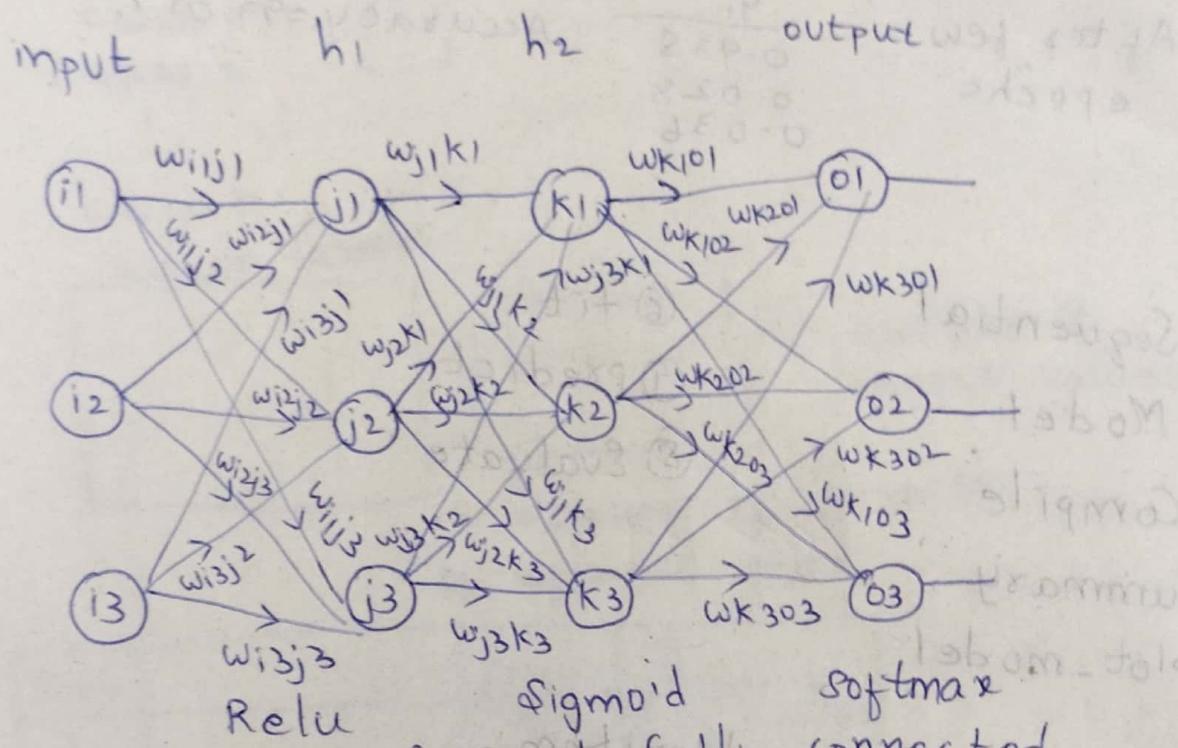
- ① Sequential
- ② Model
- ③ Compile
- ④ Summary
- ⑤ plot-model

- ⑥ fit
- ⑦ predict
- ⑧ evaluate

Back propagation algorithm:-



- 4 layers: i/p, h₁, h₂, o/p
- 3 neurons for each layer
- For simplicity all biases = 0



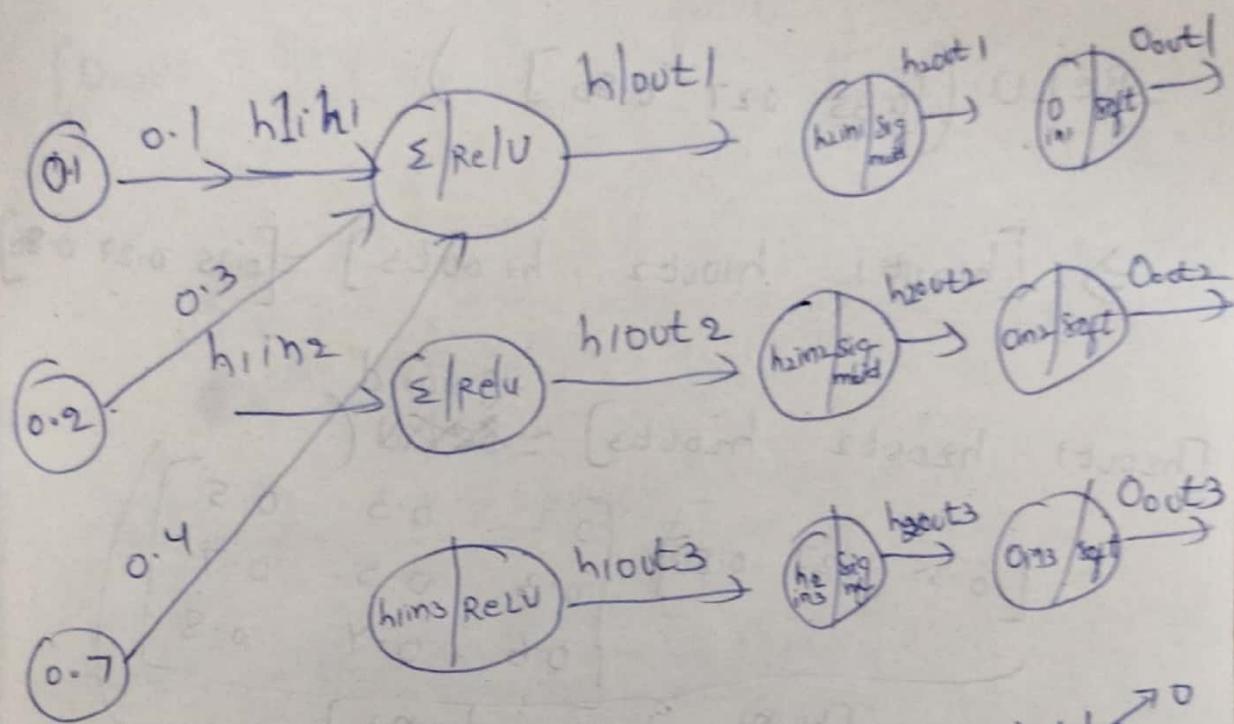
This is a feed forward fully connected neural network

$$\text{Input} = \begin{bmatrix} 0.1 & 0.2 & 0.7 \end{bmatrix} \quad y = \text{Target} = \begin{bmatrix} 1.0 & 0.0 & 0.0 \end{bmatrix}$$

$$W_{ij} = \begin{bmatrix} w_{i1j1} & w_{i1j2} & w_{i1j3} \\ w_{i2j1} & w_{i2j2} & w_{i2j3} \\ w_{i3j1} & w_{i3j2} & w_{i3j3} \end{bmatrix} = \begin{bmatrix} 0.1 & 0.2 & 0.3 \\ 0.3 & 0.2 & 0.7 \\ 0.4 & 0.3 & 0.4 \end{bmatrix}$$

$$W_{jk} = \begin{bmatrix} w_{j1k1} & w_{j1k2} & w_{j1k3} \\ w_{j2k1} & w_{j2k2} & w_{j2k3} \\ w_{j3k1} & w_{j3k2} & w_{j3k3} \end{bmatrix} = \begin{bmatrix} 0.2 & 0.3 & 0.5 \\ 0.3 & 0.5 & 0.7 \\ 0.6 & 0.4 & 0.8 \end{bmatrix}$$

$$W_{K1} = \begin{bmatrix} w_{K111} & w_{K112} & w_{K113} \\ w_{K211} & w_{K212} & w_{K213} \\ w_{K311} & w_{K312} & w_{K313} \end{bmatrix} = \begin{bmatrix} 0.1 & 0.4 & 0.8 \\ 0.3 & 0.7 & 0.2 \\ 0.5 & 0.2 & 0.9 \end{bmatrix}$$



$$\begin{aligned} h1in_1 &= i_1 * w_{ij_1} + i_2 * w_{ij_2} + i_3 * w_{ij_3} + b_{j_1} \\ &= [i_1 \ i_2 \ i_3] \begin{bmatrix} w_{ij_1} \\ w_{ij_2} \\ w_{ij_3} \end{bmatrix} + b_{j_1} \end{aligned}$$

$$h1out_1 = \text{ReLU}(h1in_1) = \max(0, h1in_1)$$

$$\begin{aligned} h1out_1 &= \text{ReLU}([i_1 \ i_2 \ i_3]) \\ &= \text{ReLU} \left(\begin{bmatrix} w_{ij_1} & w_{ij_2} & w_{ij_3} \\ w_{ij_1} & w_{ij_2} & w_{ij_3} \\ w_{ij_1} & w_{ij_2} & w_{ij_3} \end{bmatrix} \right) \end{aligned}$$

$$= \text{ReLU} \left(\begin{bmatrix} 0.1 & 0.2 & 0.3 \\ 0.3 & 0.2 & 0.7 \\ 0.4 & 0.3 & 0.9 \end{bmatrix} \right)$$

$$= \text{ReLU} \left(\begin{bmatrix} 0.35 & 0.27 & 0.80 \end{bmatrix} \right)$$

$$\Rightarrow [h_{1\text{out}1} \ h_{1\text{out}2} \ h_{1\text{out}3}] = [0.35 \ 0.27 \ 0.80]$$

sigmoid

$$[h_{2\text{out}1} \ h_{2\text{out}2} \ h_{2\text{out}3}] = \text{sigmoid} \left(\begin{bmatrix} 0.2 & 0.3 & 0.5 \\ 0.3 & 0.5 & 0.7 \\ 0.6 & 0.4 & 0.8 \end{bmatrix} \right)$$

$$h_{2\text{out}1} = \frac{1}{1 + e^{-h_{2\text{in}1}}}$$

$$[h_{2\text{out}1} \ h_{2\text{out}2} \ h_{2\text{out}3}]$$

$$= \text{sigmoid} \left(\begin{bmatrix} h_{2\text{in}1} & h_{2\text{in}2} & h_{2\text{in}3} \end{bmatrix} \right)$$

$$= \text{sigmoid} \left(\begin{bmatrix} 0.631 & 0.56 & 1.004 \end{bmatrix} \right)$$

$$= [0.6527 \ 0.6364 \ 0.7318]$$

$$[O_{\text{out}1} \ O_{\text{out}2} \ O_{\text{out}3}] = \text{Softmax} \left(\begin{bmatrix} 0.1 & 0.4 & 0.8 \\ 0.3 & 0.7 & 0.2 \\ 0.5 & 0.2 & 0.9 \end{bmatrix} \right)$$

$$= \text{Softmax} \left(\begin{bmatrix} O_{\text{in}1} & O_{\text{in}2} & O_{\text{in}3} \end{bmatrix} \right)$$

$$\text{softmax} = \text{softmax}([o_{in1} \quad o_{in2} \quad o_{in3}])$$

$$o_{out1} = \frac{e^{o_{in1}}}{e^{o_{in1}} + e^{o_{in2}} + e^{o_{in3}}}$$

$$[o_{out1} \quad o_{out2} \quad o_{out3}] = \begin{bmatrix} e^{o_{in1}} \\ e^{o_{in1}} + e^{o_{in2}} + e^{o_{in3}} \end{bmatrix}$$

$$\frac{e^{o_{in2}}}{e^{o_{in1}} + e^{o_{in2}} + e^{o_{in3}}}$$

$$\frac{e^{o_{in3}}}{e^{o_{in1}} + e^{o_{in2}} + e^{o_{in3}}}$$

$$= [0.223 \quad 0.3064 \quad 0.470]$$

$$\text{output}^T = \text{softmax}(\text{sigmoid}(\text{ReLU}(\text{Input}^T w_{jk}) w_{jk}))$$

$$\text{output} = f(x_j \theta)$$

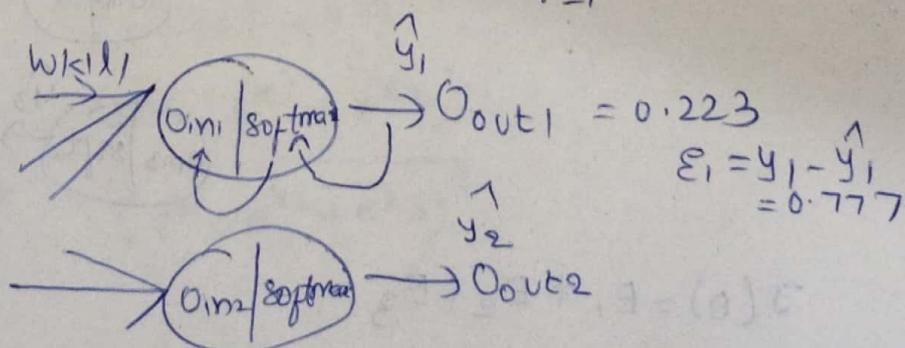
$$\text{s.GD} = \theta^{(2)} = \theta^{(1)} - \eta \nabla_{\theta} J(\theta)$$

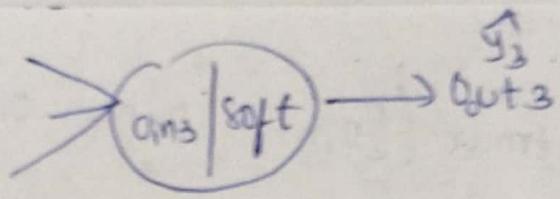
Optimized

$$\text{Gradient descent: } \theta^{(2)} = \theta^{(1)} - \eta \nabla_{\theta} J(\theta); \eta = 0.01$$

$$\text{Cost function: categorical cross-entropy } J(\theta) = - \sum_{i=1}^3 \text{Target}_i \log(o_{outi})$$

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

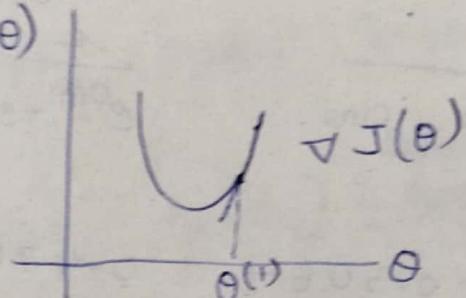




$$\frac{\delta \hat{y}_1}{\delta w_{k1l1}} = \frac{\delta \hat{y}_1}{\delta \phi} \times \frac{\delta \phi}{\delta o_{in1}} \times \frac{\delta o_{in1}}{\delta w_{k1l1}}$$

$$o_{out1} = \hat{y}_1 = \text{softmax}(h_2 o_{out1} w_{k1l1} + h_2 o_{out2} w_{k2l1} + h_2 o_{out3} w_{k3l1})$$

$J(\theta)$

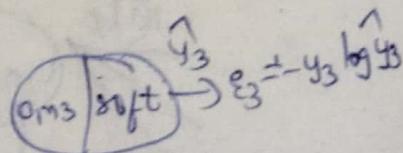
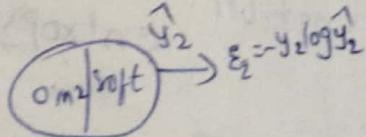
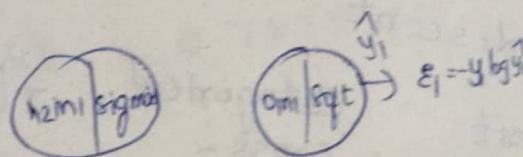
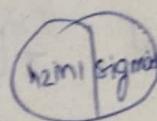
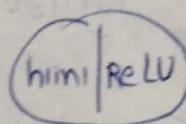


$$\theta^{(2)} = \theta^{(1)} - n \nabla J(\theta^{(1)})$$

$$\begin{aligned} \frac{\delta J(\theta)}{\delta y_1} &= \frac{\delta}{\delta \hat{y}_1} [-y_1 \log \hat{y}_1 - y_2 \log \hat{y}_2 - y_3 \log \hat{y}_3] \\ &= \left[-y_1 \left(\frac{1}{\hat{y}_1} \right) - 0 - \theta \right] = -y_1 / \hat{y}_1 \end{aligned}$$

$$\frac{\delta J(\theta)}{\delta w_{k1l1}} = \frac{\delta J(\theta)}{\delta \hat{y}_1} * \frac{\delta \hat{y}_1}{\delta \phi} * \frac{\delta \phi}{\delta o_{in1}} * \frac{\delta o_{in1}}{\delta w_{k1l1}}$$

(i)



$$J(\theta) = E_1 + E_2 + E_3$$

$$J(\theta) = E_1 + E_2 + E_3$$

$$= -y_1 \log \hat{y}_1 - y_2 \log \hat{y}_2 - y_3 \log \hat{y}_3$$

$\nabla_{\theta} J(\theta)$ shows the cost value changes when all the parameters gets updated.

$$f(x; \theta) = \text{ReLU}(\text{Input}^T w_{ik})$$

$$\text{softmax}(\text{sigmoid}(\text{ReLU}(\text{Input}^T w_{ij})w_{jk})w_{kl})$$

$$\frac{\partial J(\theta)}{\partial w_{kli}} = \frac{\partial J(\theta)}{\partial \text{softmax}} * \frac{\partial \text{softmax}}{\partial o_{inj}} * \frac{\partial o_{inj}}{\partial w_{kli}}$$

Derivative of softmax:-

$$y_1 = \frac{e^{x_1}}{e^{x_1} + e^{x_2} + e^{x_3}}$$

$$y_2 = \frac{e^{x_2}}{e^{x_1} + e^{x_2} + e^{x_3}}$$

$$y_3 = \frac{e^{x_3}}{e^{x_1} + e^{x_2} + e^{x_3}}$$

$$\frac{\partial y_1}{\partial x_1} = \frac{\partial}{\partial x_1} \left(\frac{e^{x_1}}{e^{x_1} + e^{x_2} + e^{x_3}} \right) = \frac{(e^{x_1} + e^{x_2} + e^{x_3}) e^{x_1} - e^{x_1} (e^{x_1})}{(e^{x_1} + e^{x_2} + e^{x_3})^2}$$

$$\frac{\partial y_1}{\partial x_1} = \frac{e^{x_1} (e^{x_2} + e^{x_3})}{(e^{x_1} + e^{x_2} + e^{x_3})^2}$$

$$\frac{\partial y_2}{\partial x_2} = \frac{e^{x_2} (e^{x_1} + e^{x_3})}{(e^{x_1} + e^{x_2} + e^{x_3})^2}$$

$$\frac{\partial y_3}{\partial x_3} = \frac{e^{x_3} (e^{x_1} + e^{x_2})}{(e^{x_1} + e^{x_2} + e^{x_3})^2}$$

derivative of sigmoid

$$y = \frac{1}{1+e^{-x}}$$

$$\frac{\partial}{\partial x}(f(x)) = -\frac{1}{x^2}$$

$$y = \frac{e^x}{1+e^x}$$

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

(or)

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

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

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

Derivation of ReLU

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

$$= \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases}$$

$$\frac{\partial \text{ReLU}}{\partial x} = \begin{cases} 1 & x > 0 \\ 0 & \text{elsewhere} \end{cases}$$

$$w_{kl}' = w_{kl} - \eta \delta w_{kl}$$

$$w_{jk}' = w_{jk} - \eta \delta w_{jk}$$

$$w_{ij}' = w_{ij} - \eta \delta w_{ij}$$

$$\delta w_{kl} = \begin{bmatrix} \frac{\partial \varepsilon_1}{\partial w_{k1l1}} & \frac{\partial \varepsilon_2}{\partial w_{k1l2}} & \frac{\partial \varepsilon_3}{\partial w_{k1l3}} \\ \frac{\partial \varepsilon_1}{\partial w_{k2l1}} & \frac{\partial \varepsilon_2}{\partial w_{k2l2}} & \frac{\partial \varepsilon_3}{\partial w_{k2l3}} \\ \frac{\partial \varepsilon_1}{\partial w_{k3l1}} & \frac{\partial \varepsilon_2}{\partial w_{k3l2}} & \frac{\partial \varepsilon_3}{\partial w_{k3l3}} \end{bmatrix}$$

chain-rule

$$\frac{\partial \varepsilon_1}{\partial w_{K1l1}} = \frac{\partial \varepsilon_1}{\partial y_1} * \frac{\partial \hat{y}_1}{\partial o_{in1}} * \frac{\partial o_{in1}}{\partial w_{K1l1}}$$

↓
w.r.t
activation
function

$$\varepsilon_1 = -y_1 \log \hat{y}_1$$

$$\frac{\partial \varepsilon_1}{\partial \hat{y}_1} = \frac{\partial}{\partial \hat{y}_1} (-y_1 \log \hat{y}_1) = -y_1 \left(\frac{1}{\hat{y}_1} \right)$$

$$\frac{\partial \hat{y}_1}{\partial o_{in1}} = \frac{e^{\hat{y}_1} (e^{\hat{y}_2} + e^{\hat{y}_3})}{(e^{\hat{y}_1} + e^{\hat{y}_2} + e^{\hat{y}_3})^2}$$

$$\frac{\partial o_{in1}}{\partial w_{K1l1}} = \frac{\partial}{\partial w_{K1l1}} \left(h_2 \text{out}_1 * w_{K1l1} + h_2 \text{out}_2 * w_{K2l1} + h_2 \text{out}_3 * w_{K3l1} \right)$$

$$\frac{\partial o_{in1}}{\partial w_{K1l1}} = h_2 \text{out}_1$$

$$\delta w_{jk} = \begin{cases} \frac{\partial \varepsilon}{\partial w_{jk}} & k=1 \\ \frac{\partial \varepsilon}{\partial w_{jk}} & k=2 \\ \frac{\partial \varepsilon}{\partial w_{jk}} & k=3 \end{cases}$$

$$\delta w_{kl} = \begin{cases} \frac{\partial \varepsilon}{\partial w_{kl}} & l=1 \\ \frac{\partial \varepsilon}{\partial w_{kl}} & l=2 \\ \frac{\partial \varepsilon}{\partial w_{kl}} & l=3 \end{cases}$$

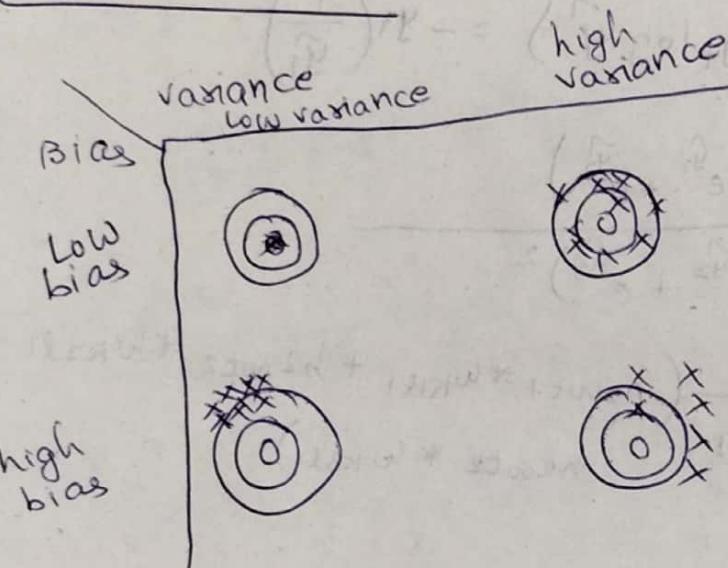
UNIT-11

*Regularization & Optimization:

Regularization

what is meant by Generalisation?

Bias and Variance



- Generalisation is the ability of model to perform well on training & testing data.

- Regularization reduce the generalisation / test error possibly at the expense of training error

- Regularization is the strategy to reduce the test error possibly at the expense of increased training error.

- It is any modification that we make to a learning algorithm that is intended to reduce its generalisation error but not its training error.

- There are mainly 3 strategies used by regularization

① Constraints based regularization

② In this we put an extra constraint on

the model such as adding restrictions on the parameter values.

② Penalties based regularization

In this we add extra terms in the objective cost function

③ Ensemble methods

In this we combine multiple hypothesis that explains the training data. Any of these regularization methods works by trading increased bias for reduced variance. An effective regularizer is one that makes a profitable trade off such as reducing variance significantly while not overly increasing the bias.

Parameter Norm Penalties

Many regularization approaches

J-cost fn
θ-parameters
wt
x-training set
y-labels

$$\tilde{J}(\theta; x, y) = J(\theta; x, y) + \alpha \cdot \Omega(\theta)$$

↓
cost function
of the
regularised
model

↳ cost function
of an
unregularized
model

hyperparameters
that weights
the contribution
of penalty term

↳ Penalty term
↳ which is a
function of
only the
parameters
↳ This term is
nothing to do
with training
data set

when $\alpha=0 \Rightarrow$ no regularization

higher the
value of α \Rightarrow more the regularization

consider an unregularized model

→ we try to calculate w^* at which $J(\theta; x, y)$ is min

In regularised model,

our goal becomes

→ find \tilde{w} at which $\tilde{J}(\theta; x, y)$ is min

↳ this may increase data term

↳ decreased by reducing penalty term.

Data & Penalty terms are independent. So

that they can ↑ or ↓

↳ depends on features ↳ depends on parameters

Ultimately our goal is to minimize $\tilde{J}(\theta; x, y)$

Different Regularisation functions:-

$$P\text{-Norm } L_p(\theta) = \left(\sum_i |\theta_i|^p \right)^{1/p} = \|\theta\|_p$$

where $p = 0.5, 1, 2, 3, \dots$ any positive number

L^2 Norm. Regularization (OR) Ridge Regularization

$$\|\theta\|_2 = \left(\sum_i |\theta_i|^2 \right)^{1/2} = \theta \theta^T$$

$$\tilde{J}(\theta; x, y) = J(\theta; x, y) + \alpha \theta \theta^T$$

For simplicity ignore the bias parameter

$$\tilde{J}(w; x, y) = J(w; x, y) + \alpha w w^T$$

$$\tilde{J}(w; x, y) = J(w; x, y) + \alpha w_2$$

GD: $w \leftarrow w - \nabla_w \tilde{J}(w; x, y)$

L^1 Regularization (OR) Lasso Regularization

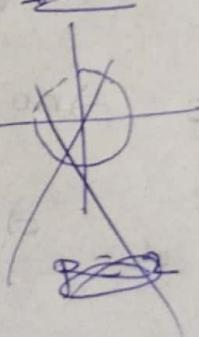
$$\frac{w_1}{w_2}$$

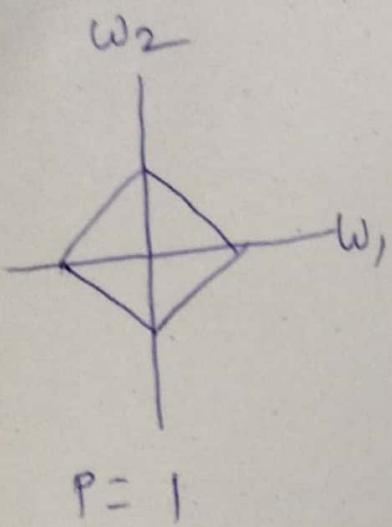
$$\frac{w_1}{w_2}$$

$$P=0.5$$

$$\frac{w_1}{w_2}$$

$$P=2$$





Lasso

