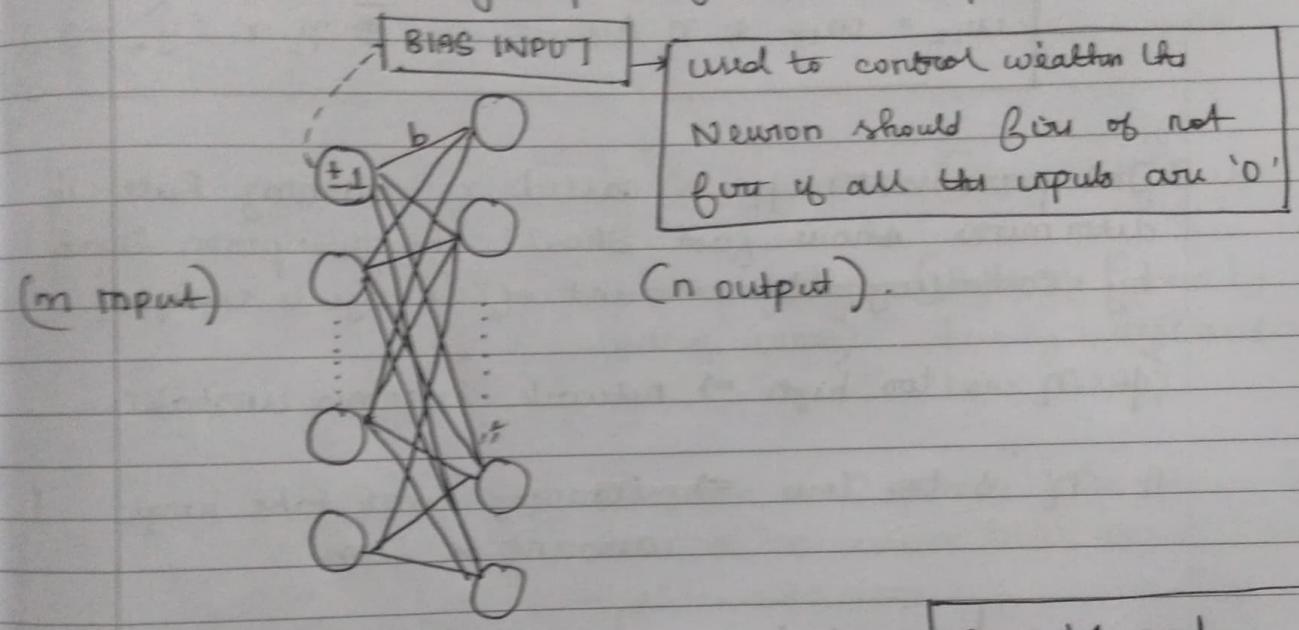


PERCEPTRON: (SINGLE LAYER PERCEPTRON)

- A linear model used for binary classification
- Step Function for activation i.e $f(x) = \begin{cases} 0, & x \leq \alpha \\ 1, & x > \alpha \end{cases}$
 $\alpha \Rightarrow$ THRESHOLD VALUE.
- Contains only input layer and output layer.



At each output neuron k we have,

$$\hat{y}_k = f\left(\sum_{i=0}^m x_i w_{ik}\right) \Rightarrow \text{Output of } k\text{th neuron.}$$

$x_0 = \pm 1$ and
 $w_{0k} = b$ which is weight of bias

⇒ The weights of the connections are updated through perception learning algorithm.

At the k^{th} output neuron the error E_k is Given by

$$E_k = \hat{y}_k - y_k \quad \text{where } \hat{y}_k = \text{predicted value}$$

y_k = actual value.

if $E_k > 0 \Rightarrow$ weights should be reduced

if $E_k < 0 \Rightarrow$ weights should be increased

$$\Rightarrow \Delta w_{ik} = -\eta \cdot x_i (\hat{y}_k - y_k)$$

Therefore the weight updation is given by.

$$w_{ik} = w_{ik} + \eta \Delta w_{ik} = w_{ik} - \eta \cdot x_i (\hat{y}_k - y_k)$$

The parameter η is called as learning Rate. It determines how fast should the perceptron learn by controlling the amount of change in weight.

If η is too high \Rightarrow network becomes unstable.

If η is too low \Rightarrow learning will take large amount of time.

$$0.1 < \eta < 0.4 \Rightarrow \text{optimal Range}$$

MULTILAYER PERCEPTRON

* MULTILAYER PERCEPTRON (MULTILAYER FEED FORWARD NETWORK)

\rightarrow Neurons are arranged into groups called as layers.

- \rightarrow Contains 3 main layers:
 - ① Input layer
 - ② Hidden layers
 - ③ Output layer

- ⇒ Once a given layer all the Neurons will have Same activation function
- ⇒ For input layer, input will be the actual data vector. For other layers, input will be activation of previous layer
- ⇒ Number of neurons in the input layer will be equal to number of features in the input.
- ⇒ We have two main Steps:

① Forward Propagation (Predicting Output)

② Backward Propagation (Reducing Error by (Backpropagation of error) tuning weight and bias value).

⇒ Each neuron can be given a Separate bias value (Or) each layer can have a common bias value, This is based on the implementation

⇒ BACKPROPAGATION: it is the method of calculating and propagation error from output layer to all hidden layers.

To reduce the error we using "optimisation algorithm"



Note is calculated using "loss function" ⇒ The goal of optimiser to to ~~to~~ minimise the loss function

linear Non-linear

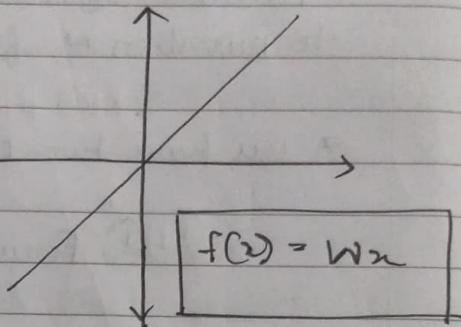
DATE

⇒ ACTIVATION FUNCTION : used to introduce non-linearity in the decision boundary (Because not all problems are linearly separable)

without ~~non-linear~~ activation functions the neural network will just perform linear weight sum (similar to linear regression).

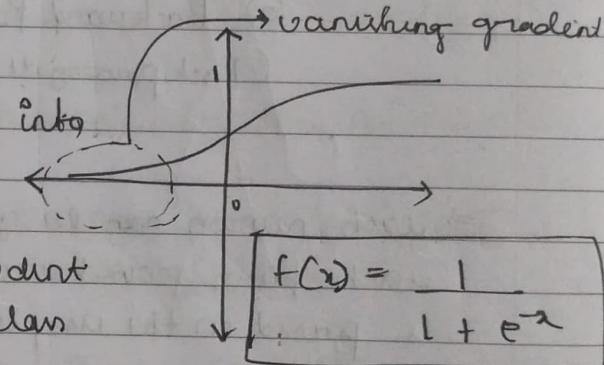
① LINEAR: $f(x) = Wx$.

⇒ identity function
⇒ used in input layer?



② SIGMOID:

⇒ Reduce the input into the Range $0 \rightarrow 1$.



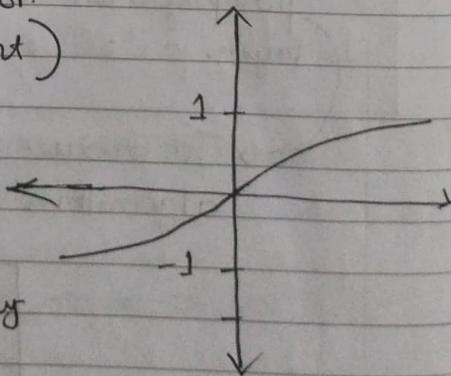
⇒ Outputs independent probability for each class

~~non-linear~~: → used in binary classification output layer.

③ TANH: (Hyperbolic Tangent)

⇒ Reduce input to the Range $-1 \rightarrow 1$

⇒ it can deal more easily with negative values.



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

Both Sigmoid and tanh can cause vanishing Gradient Problem (For extreme high / low value) \Rightarrow So we use ReLU.

DATE

But it can cause exploding Gradient
↳ need to scale up.

④ SOFTMAX:

- provides the probability distribution of each class.
- used in output layer for multi-class classifications

$$f(x_i) = \frac{e^{x_i}}{\sum_i e^{x_i}} \Rightarrow \text{if output layer contains } M \text{ neurons and } x \text{ be the output value/activation}$$

$$\Rightarrow f(x_i) = \frac{e^{x_i}}{\sum_i^M e^{x_i}}$$

(Sum of all probability will be 1)

(Class with maximum probability will be the result).

⑤ RECTIFIED LINEAR UNIT (ReLU):

\Rightarrow overcomes vanishing/exploding gradient problem.

\Rightarrow ReLU is highly used in hidden layer

\Rightarrow Some variant of ReLU are:

$$(i) \text{ Leaky ReLU} = f(x) = \max(0.01x, x).$$

$$(ii) \text{ Parametric ReLU} = f(x) = \max(\alpha x, x)$$

$$(iii) \text{ Exponential ReLU} = f(x) = \begin{cases} x, & \text{if } x \geq 0 \\ \alpha(e^x - 1), & \text{if } x < 0 \\ \downarrow \\ \max(\alpha(e^x - 1), x). \end{cases}$$

LOSS FUNCTION: It quantifies how close the predicted value is to the Real Value.

For Given weights w and Bias b we have $L(w, b)$ as loss function

$$\hat{y}_i = h_{w,b}(x) \Rightarrow \text{predicted value for } i^{\text{th}} \text{ Sample}$$

$$y_i \Rightarrow \text{Real value for } i^{\text{th}} \text{ Sample}$$

LOSS FUNCTION FOR REGRESSION

① Mean Squared Error

$$L(w, b) = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

if more than one feature.

→ der in output

$$\Rightarrow \frac{1}{N} \sum_{i=1}^N \left(\sqrt{\sum_{j=1}^M (\hat{y}_{ij} - y_{ij})^2} \right)^2$$

$$\Rightarrow \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M (\hat{y}_{ij} - y_{ij})^2 //$$

② Mean Absolute Error

$$L(w, b) = \frac{1}{N} \sum_{i=1}^N |\hat{y}_i - y_i|$$

$$\Rightarrow \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M |\hat{y}_{ij} - y_{ij}|$$

③ Root Mean Squared

$$L(w, b) = \sqrt{MSE}$$

LOSS FUNCTION FOR CLASSIFICATION

Cross entropy loss (CE)

Negative Log Likelihood

- ↳ Binary
- ↳ Categorical

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

$$CCE = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} \log (\hat{y}_{ij})$$

Hinge loss:

↳ designed for SVM

Result should be $-1 / 1$.

$$L(w, b) = \frac{1}{N} \sum_{i=1}^N \max(0, 1 - y_i \cdot \hat{y}_i)$$

$$\frac{1}{N} \sum_{i=1}^N \max(0, 1 - y_i \cdot \hat{y}_i)$$

$$\sum_{i=1}^N \sum_{j=1}^M \max(0, y_{ij} \times \hat{y}_{ij}) \Rightarrow \text{only used for binary classification}$$

MSE and MAE are prone to outliers due to which we can use **HUBER LOSS** → For Regression

$$\hookrightarrow L(w, b) = \begin{cases} \frac{1}{2} (y_i - \hat{y}_i)^2, & |y_i - \hat{y}_i| \leq \delta \\ \delta |y_i - \hat{y}_i| - \frac{1}{2} \delta^2, & \text{else.} \end{cases}$$

Here δ is a hyperparameter which depicts the tolerance to outliers.

* GRADIENT DESCENT:

Huber loss is like a combination of MSE and MAE

→ An optimisation algorithm for ~~minimising~~ minimising the cost / loss function

- (i) Batch Gradient Descent
- (ii) Stochastic Gradient Descent
- (iii) Mini-Batch Gradient Descent

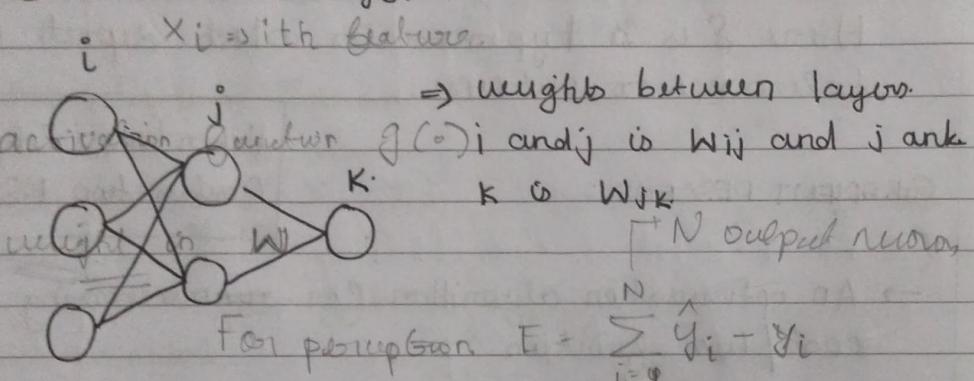
- ⇒ In (i) The complete training set is used for every iteration / epoch.
- ⇒ In (ii) Only Single Sample is used for one iteration (i.e. Number of iterations / epochs = Number of Samples in training Set → maximum).
- ⇒ In (iii) A Batch of Samples used for each iteration from the training dataset.
- ⇒ Some other optimisers are: RMSProp, ADADelta, ADAGrad, Adam.

* BACKPROPAGATION DERIVATION:

⇒ Backpropagation is that method of sample calculating and propagating error from output layer to the entire network algorithm (Stochastic Gradient Descent).

⇒ It updates weight of the network by using a Gradient Descent Algorithm. as for Singil input (x_i) → M feature Sample

⇒ Consider a 3 layer MLP.



(INPUT) (HIDDEN) (OUTPUT)

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

we will discuss using Sigmoid activation function for the hidden and output layers.

let i = input layer if L Layer \exists

 $O_K = \frac{\sum_{j=1}^L O_j \cdot w_{kj}}{1 + e^{-\sum_{j=1}^L O_j \cdot w_{kj}}} \quad | \quad i=1 \quad O_K = \text{output of neuron in layer } K$

$I_K = \sum O_j \cdot w_{kj}, \quad I_j = \sum x_i \cdot w_{ji} \quad | \quad L-1$

$E = \frac{1}{2} \sum_{i=1}^N \left(\sum_{k=1}^M g(w_{ik} \cdot a_j) \right)$

O_j = output of neuron in layer K

Why I_k and I_j are inputs

According to gradient descent we have

$$N \Rightarrow \Delta W \propto -\frac{\partial E}{\partial W} \quad (\Delta W \text{ is change in weights})$$

$N_i \Rightarrow$ Number of neurons in i th layer ($\frac{\partial E}{\partial W}$ is Gradient of loss function).

$$N_o \Rightarrow \Delta W = -\eta \frac{\partial E}{\partial W} \quad (\eta \text{ is learning rate})$$

$N_{Hj} \Rightarrow$ Number of neurons in j th hidden layer.

* On Output layer:

$$E = \frac{1}{2} \sum_{k=0}^{N_o} (Y_k - \hat{Y}_k)^2 = \frac{1}{2} \sum_{k=0}^{N_o} \left(g\left(\sum_{j=0}^{N_{Hj}} w_{kj} o_j\right) - Y_k \right)^2$$

$$\Delta w_{jk} = -\eta \frac{\partial E}{\partial w_{jk}}$$

Consider a perceptron \Rightarrow No Hidden layers

$$\Rightarrow E = \sum_{k=0}^{N_o} \left(\frac{1}{2} \left(\sum_{i=0}^{N_i} w_{ik} o_i + b_k \right)^2 \right) \quad \begin{array}{l} \text{Perceptron activation} \\ \text{and Hidden layers} \end{array}$$

By chain rule we have

$$\Delta w_{jk} = -\eta \left[\frac{\partial E}{\partial o_k} \cdot \frac{\partial o_k}{\partial I_k} \cdot \frac{\partial I_k}{\partial w_{jk}} \right]$$

To reduce for all k we have $\frac{\partial o_k}{\partial I_k} = \frac{\partial}{\partial I_k} \frac{1}{1+e^{-I_k}} = \frac{-e^{-I_k}}{(1+e^{-I_k})^2}$

$$\frac{\partial E}{\partial o_k} = \frac{1}{2} \sum_{k=0}^{N_o} (Y_k - \hat{Y}_k) \cdot 2 \hat{Y}_k$$

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

$$\frac{\partial E}{\partial I_k} = - (1+e^{-I_k})^{-2} (-e^{-I_k})$$

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

$$\left| \begin{array}{l} \frac{\partial O_k}{\partial I_k} = O_k(1-O_k) \\ \hline \end{array} \right|$$

$$\frac{\partial I_k}{\partial W_{jk}} = \frac{\partial}{\partial W_{jk}} \sum O_j W_{jk} = O_j$$

$$\Rightarrow \boxed{\Delta W_{jk} = -n \cdot O_k \cdot (1-O_k) (O_T - O_k) (O_j)}.$$

$$\Rightarrow (\text{new}) W_{jk} = (\text{old}) W_{jk} + \Delta W_{jk}$$

For Hidden Layer,

$$\Delta W_{ij} = -n \frac{\partial E}{\partial W_{ij}}$$

$$= \frac{\partial E}{\partial O_k} \cdot \frac{\partial O_k}{\partial I_k} \cdot \frac{\partial I_k}{\partial O_j} \cdot \frac{\partial O_j}{\partial I_k} \cdot \frac{\partial I_k}{\partial W_{ij}} \quad (\text{By chain rule})$$

$$\frac{\partial I_k}{\partial O_j} = \cancel{O_{jk}} W_{jk} \quad \frac{\partial O_j}{\partial I_k} = O_j(1-O_j)$$

$$\frac{\partial I_j}{\partial W_{ij}} = x_i$$

$$\boxed{\Delta W_{ij} = O_k \cdot x_i \cdot O_k \cdot O_j \cdot (O_T - O_k) (1 - O_k) (1 - O_j) W_{jk}}$$

From ΔW_{jk} and ΔW_{ij} we have

$$\left. \begin{array}{l} E_k = O_k(1-O_k)(O_T - O_k) \\ E_j = O_j(1-O_j) \sum_{k=1}^N W_{jk} E_k \end{array} \right| \Rightarrow \begin{array}{l} W_{ij} = W_{ij} + n O_i E_j \\ B_i = B_i + n (\pm 1) (E_i) \end{array}$$

* HYPERPARAMETERS :

⇒ Configuration options that can be varied or altered to improve the performance of model.

- ① No. of Hidden Layers
- ② No. of Neurons in each layer
- ③ Activation Function
- ④ Loss Function
- ⑤ Optimizer Function
- ⑥ No. of epochs/iterations
- ⑦ Type of Regularization
- ⑧ Weight initialisation
- ⑨ Learning Rate & Momentum.

* BIAS AND VARIANCE : (Trade off)

⇒ When a model is trained it captures the patterns in the training data and tries to generalise it.

⇒ BIAS is the measure of our model's inability to learn patterns from data set.

⇒ When model has HIGH BIAS ⇒ did not train well in training Set

When model has LOW BIAS ⇒ it has identified new patterns from Training Set properly.

⇒ VARIANCE is the measure of model's fluctuation in performance when there is a change in data set.

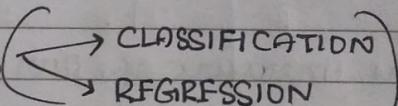
⇒ UNDERFITTING: when a model does not perform well in training data.
(HIGH BIAS).

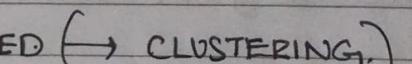
⇒ OVERRFITTING: when a model performs well in train data but it lacks performance in testing data
(LOW BIAS, HIGH VARIANCE)

⇒ A Good model should have LOW BIAS and LOW VARIANCE.

MACHINE LEARNING: A Algorithm is said to be learning from experience E with respect to some task T and performance measure P if the performance P increases with increase in experience E.

TYPES OF MACHINE LEARNING:-

① SUPERVISED 
CLASSIFICATION
REGRESSION

② UNSUPERVISED 
CLUSTERING.

③ SEMI-SUPERVISED

④ REINFORCEMENT

LINEAR SEPARABILITY:

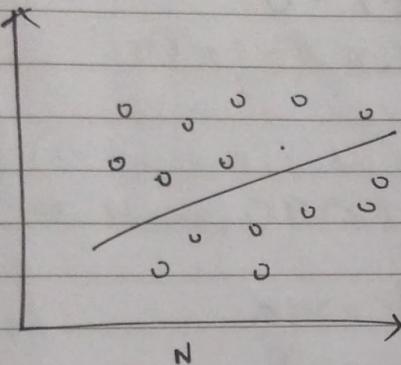
A data is Said to be linearly Separable if only a single line / Plane / Hyperplane can partition the data into 2 classes.

Eg: Perceptron can only work on a data set that's linearly separable.

→ One of the main reason for MLP improvement

LINEAR REGRESSION:

Consider the following equation $y = \theta_0 x + \theta_1$



Linear regression
Gives to fit a line of the form on the given dataset.

$$E = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2 \Rightarrow \text{where } \hat{y}_i \text{ is the predicted value and } y_i \text{ is the actual value}$$

The error E is called as MEAN SQUARED ERROR.
we have to Reduce this error then it is called LEAST SQUARE REGRESSION.

$$E = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2 = \frac{1}{N} \sum_{i=1}^N (\theta_0 x_i + \theta_1 - y_i)^2$$

$$\frac{\partial E}{\partial \theta_0} = 2 \times \frac{1}{N} \sum_{i=1}^N (\theta_0 x_i + \theta_1 - y_i)(x_i + 0 - 0) = 0$$

$$= \theta_0 \bar{x^2} + \theta_1 \bar{x} - \bar{x} \cdot \bar{y} = 0 \quad \dots \text{---(1)}$$

$$\frac{\partial E}{\partial \theta_1} = 2 \times \frac{1}{N} \sum_{i=1}^N (\theta_0 x_i + \theta_1 - y_i)(0 + 1 - 0) = 0$$

$$= \theta_0 \bar{x} + \theta_1 - \bar{y} = D \quad (1)^2$$

$$\Rightarrow \theta_1 = \bar{y} - \theta_0 \bar{x} \quad (2)$$

From ① and ② we have.

$$\theta_0 \bar{x}^2 + (\bar{y} - \theta_0 \bar{x}) \cdot \bar{x} - \bar{x} \cdot \bar{y} = D$$

$$\theta_0 \bar{x}^2 + \bar{x} \cdot \bar{y} - \theta_0 \bar{x}^2 - \bar{x} \cdot \bar{y} = D$$

$$\theta_0 = \frac{\bar{x} \cdot \bar{y} - \bar{x} \cdot \bar{y}}{\bar{x}^2 - \bar{x}^2}$$

\Rightarrow

$$\boxed{\theta_0 = \frac{\bar{x} \cdot \bar{y} - \bar{x} \cdot \bar{y}}{\bar{x}^2 - \bar{x}^2}}$$

(X)

$$\boxed{\theta_1 = \bar{y} - \theta_0 \bar{x}}$$

\Rightarrow why Huber Loss for Regression?

MSE is Highly Sensitive to outliers because we square the other all error values. MAE just take the absolute value of error it gives same weightage to all errors due to which outlier can be completely Neglected. Huber Loss provides Balance (8).

* BAYESIAN STATISTICS:

↳ (methods that used Bayes Theorem)

⇒ BAYES THEOREM: Finds the probability of an event occurring, given the probability of another event that has already occurred

↳ (called as conditional probability)

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

⇒ Here, $\Rightarrow P(B)$ is called as evidence.

⇒ $P(A)$ is called as prior probability
(i.e., probability occurrence of event A without considering event B)

⇒ $P(A|B)$ is called as posterior Probability
(i.e., probability of occurrence of event A with considering event B)

⇒ we can $P(Y|x)$ using above method.

$Y \Rightarrow$ Target Class and $X =$ Feature Vector.

$$\Rightarrow P(Y_j | x_1, x_2, \dots, x_m) = \frac{P(x_1, x_2, \dots, x_m | Y_j) \cdot P(Y_j)}{P(x_1) \cdot P(x_2) \cdots P(x_m)}$$

$$= \frac{P(Y_j | x_1) \cdot P(Y_j | x_2) \cdots P(Y_j | x_m) \cdot P(Y_j)}{P(x_1) \cdot P(x_2) \cdots P(x_m)}$$

→ since prior evidence is same.

for all target we can ignore
thus

\Rightarrow Result will be class j having max $P(y_j | x)$.

\Rightarrow There are various variation of Naive Bayes algorithms like

(i) Gaussian / Normal NB

(ii) Multinomial NB

(iii) Bernoulli NB

;

\Rightarrow Important assumption for Naive Bayes algorithm is all the features are independent.

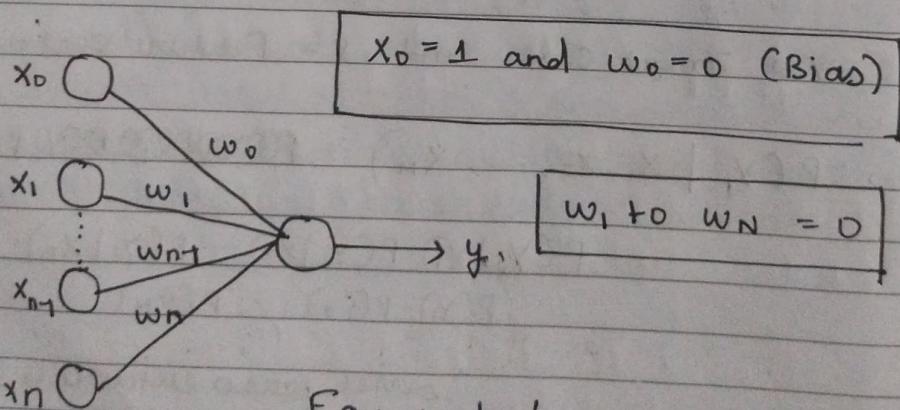
$$\text{i.e } P(x_i | x_j) = P(x_i)$$

$\overbrace{\hspace{15em}}$

* HEBBIAN LEARNING RULE (HEBB NETWORK)

\rightarrow inputs can only be 1 and -1

\rightarrow initialize all weights and bias to 0.
in Starting



For each training sample we do the following

$$\Rightarrow \Delta w_i = x_i y \text{ and } w_i = w_i + \Delta w_i$$

— X —

PROBLEMS:

① AND GATE WITH HERB NETWORK:

(B_x)(ΔB_w)

X ₀	X ₁	X ₂	Y	Δw ₀	Δw ₁	Δw ₂	Δw ₀	Δw ₁	Δw ₂
1	-1	-1	-1	-1	1	1	1	1	1
1	-1	1	-1	-1	-1	1	-1	-1	-1
1	1	-1	-1	-1	-1	-1	1	1	1
1	1	1	1	1	1	1	1	1	1
				-2	2	2			

Initial value of $w_0 = 0$ (B_w) $x_0 = 1$ (B_x).

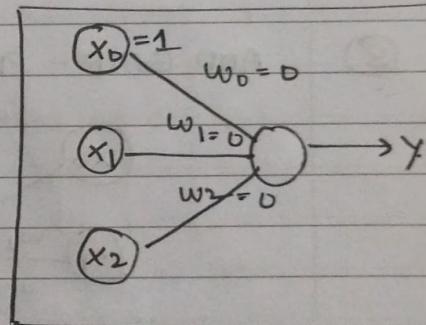
$w_1 = 0$

$w_2 = 0$

⇒ Final weights are $w_0 = -2$

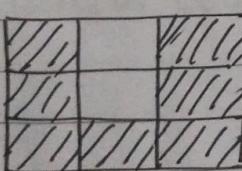
$w_1 = 2$

$w_2 = 2$



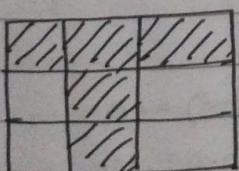
V and T.

② 3x3 IMAGE OF WITH HERB NETWORK.



1	-1	1
1	-1	1
1	1	1

(Target = 1)



1	1	1
-1	1	-1
-1	1	-1

(Target = -1)

Initial value of w_0 to w_9 is 0 and $x_0 = 1$

x_0	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	y
1	1	-1	1	1	-1	1	1	1	1	1
1	1	1	1	-1	1	-1	-1	1	-1	-1

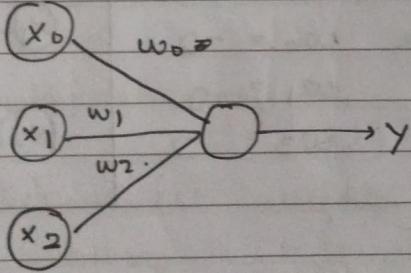
Δw_0	Δw_1	Δw_2	Δw_3	Δw_4	Δw_5	Δw_6	Δw_7	Δw_8	Δw_9
1	1	-1	1	1	-1	-1	1	1	1
-1	-1	-1	-1	1	-1	1	-1	-1	1
0	0	-2	0	2	-2	2	2	0	2

\Rightarrow Final weights are

$$w_0 \quad w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6 \quad w_7 \quad w_8 \quad w_9 \\ 0 \quad 0 \quad -2 \quad 0 \quad 2 \quad -2 \quad 2 \quad 2 \quad 0 \quad 2.$$

----- x -----

(3) AND GATE - PERCEPTRON :



$$x_0 = 1 \text{ (BIAS INPUT)}$$

$$w_0 = -0.5, \quad w_1 = 0.2, \quad w_2 = 0.3$$

(BIAS WEIGHT)

x_0	x_1	x_2	y	f	\hat{y}	E
1	0	0	0	$(1)(-0.5) + (0)(0.2) + (0)(0.3) = 0$	0	0
1	0	1	0	$(1)(-0.5) + (0)(0.2) + (1)(0.3) = 0$	0	0
1	1	0	0	$(1)(-0.5) + (1)(0.2) + (0)(0.3) = 0$	0	0
1	1	1	1	$(1)(-0.5) + (1)(0.2) + (1)(0.3) = 0$	0	-1

$$w_0 = -0.5 \quad \text{OR} \quad (0.4)(-1)(1) = -0.1$$

$$w_1 = 0.2 \quad \text{OR} \quad (0.4)(-1)(0) = 0.6$$

$$w_2 = 0.3 \quad \text{OR} \quad (0.4)(1)(1) = 0.7$$

ITERATION #2: ($w_0 = -0.1, w_1 = 0.6, w_2 = 0.7$)

x_0	x_1	x_2	y	\hat{y}	E
1	0	0	0	$(1)(-0.1) + (0)(0.6) + (0)(0.7) = 0$	0
1	0	1	0	$(1)(-0.1) + (0)(0.6) + (1)(0.7) = 1$	-1
1	1	0	0	$(1)(-0.5) + (1)(0.6) + (0)(0.3) = 1$	1
1	1	1	1	$(1)(-0.9) + (1)(0.2) + (1)(0.3) = 0$	-1

$$w_0 = -0.1 - (0.4)(0)(1) = -0.5$$

$$w_1 = 0.6 - (0.4)(1)(0) = \cancel{-0.6} 0.6$$

$$w_2 = 0.7 - (0.4)(1)(1) = 0.3$$

$$w_0 = -0.5 - (0.4)(1)(1) = -0.9$$

$$w_1 = 0.6 - (0.4)(1)(1) = -0.2$$

$$w_2 = 0.3 - (0.4)(1)(0) = 0.3$$

$$w_0 = -0.9 - (0.4)(-1)(1) = -0.5$$

$$w_1 = 0.2 - (0.4)(-1)(1) = 0.6$$

$$w_2 = 0.3 - (0.4)(-1)(0) = 0.7$$

\Rightarrow weights are ITERATION #2 are $w_0 = -0.5$

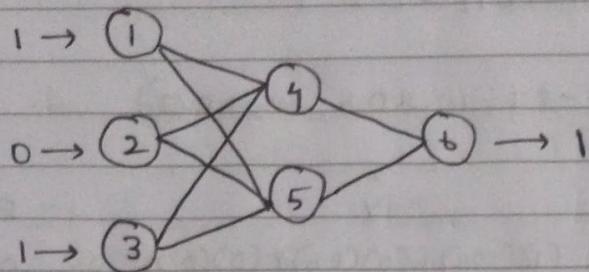
$$w_1 = 0.6$$

$$w_2 = 0.7$$

— x —

(4)

MULTILAYER PERCEPTRON:



$$w_{14} = 0.2 \quad w_{15} = -0.3 \quad w_{46} = -0.5$$

$$w_{24} = 0.4 \quad w_{25} = 0.1 \quad w_{56} = -0.2$$

$$w_{34} = -0.5 \quad w_{35} = 0.2$$

$$b_4 = 0.218 \quad b_5 = 0.194 \quad b_6 = 0.408$$

ITERATION #1

$$\begin{aligned} I_4 &= (1)(0.2) + (0)(0.4) + (1)(-0.5) + (1)(0.218) \\ &= -0.082 \end{aligned}$$

$$O_4 = \frac{1}{1 + e^{-0.082}} = 0.479$$

$$\begin{aligned} I_5 &= (1)(-0.3) + (0)(0.1) + (1)(0.2) + (1)(0.194) \\ &= 0.094 \end{aligned}$$

$$O_5 = \frac{1}{1 + e^{-0.094}} = 0.523$$

$$\begin{aligned} I_6 &= (0.479)(-0.5) + (0.523)(-0.2) + (1)(0.408) \\ &= 0.064 \end{aligned}$$

$$O_6 = \frac{1}{1 + e^{-0.064}} = 0.516$$

$$E_6 = (0.516)(1-0.516)(1-0.516) = 0.121$$

$$E_5 = (0.523)(1-0.523)(-0.2 \times 0.121) = -0.006$$

$$E_4 = (0.479)(1-0.479)(-0.5 \times 0.121) = -0.015$$

~~Initial values~~

$$w_{14} = 0.2 + (0.2)(1)(-0.015) =$$

$$w_{24} = 0.4 + (0.2)(0)(-0.015) =$$

$$w_{34} = -0.5 + (0.2)(1)(-0.015) =$$

$$w_{15} = -0.3 + (0.2)(1)(-0.006) =$$

$$w_{25} = 0.1 + (0.2)(0)(-0.006) =$$

$$w_{35} = 0.2 + (0.2)(1)(-0.006) =$$

$$w_{46} = -0.5 + (0.2)(0.479)(0.121) =$$

$$w_{56} = -0.2 + (0.2)(0.523)(0.121) =$$

$$b_4 = 0.218 + (0.2)(1)(-0.015) =$$

$$b_5 = 0.194 + (0.2)(1)(-0.006) =$$

$$b_6 = 0.408 + (0.2)(1)(0.121) =$$

— — — — — X — — — — —

NAME: BAYES

$$\textcircled{5} \quad P(\text{YES}) = \frac{1}{10} \quad P(\text{NO}) = \frac{3}{10}$$

* CGPA	P(YES)	P(NO)
≥ 9	$3/7$	$1/3$
≥ 8	$4/7$	$0/3$
< 8	$0/7$	$2/3$

* INTERACTIVENESS	P(YES)	P(NO)
YES	$5/7$	$1/3$
NO	$2/7$	$2/3$

* PRACTICAL-KNOW	P(YES)	P(NO)
EXCELLENT	$2/7$	$0/3$
GOOD	$4/7$	$1/3$
AVERAGE	$1/7$	$2/3$

* COMM-SKILL	P(YES)	P(NO)
GOOD	$4/7$	$1/3$
MODERATE	$3/7$	$0/3$
POOR	$0/7$	$2/3$

$x = \{ \geq 9, \text{ YES}, \text{ AVERAGE}, \text{ GOOD} \}$.

$$P(\text{YES} | x) = \frac{3}{7} \times \frac{5}{7} \times \frac{1}{7} \times \frac{4}{7} \times \frac{7}{10} = 0.0174$$

$$P(\text{NO} | x) = \frac{1}{3} \times \frac{1}{3} \times \frac{2}{3} \times \frac{1}{3} \times \frac{3}{10} = 0.0074.$$

=> RESULT: YES
classmate

LINEAR REGRESSION

DATE

x	y	xy	x^2
1	4	4	1
2	5	10	4
3	8	24	9
4	2	8	16

$$\bar{x} = 2.5 \quad \bar{y} = 4.75 \quad \bar{x^2} = 7.5$$

$$\bar{xy} = 11.5 \quad \bar{x} \cdot \bar{y} = 11.875$$

$$\theta_0 = \frac{11.875 - 11.5}{6.25 - 7.5} = \frac{0.375}{-1.25} = -0.3$$

$$\boxed{\theta_0 = -0.3}$$

$$\theta_1 = 4.75 + 0.3 \times 2.5 = 5.5$$

$$\boxed{\theta_1 = 5.5}$$

$$\boxed{y = -0.3x + 5.5}$$

 $\rightarrow x$

7 LOSS FUNCTION SUM

 $\rightarrow x$

* CONVOLUTION NEURAL NETWORK:

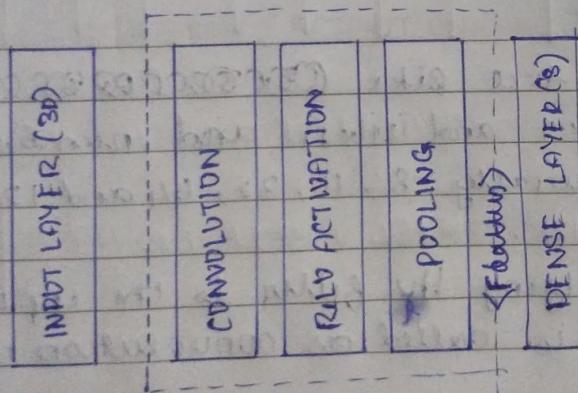
→ A fully connected neural network creates massive number of connection while modelling image data. (i.e large number of neurons and large number of connections).

→ disadvantage of Normal ANN to work with image

→ In CNN, we can arrange neurons in 3-D structure based on

- (i) width \Rightarrow image width in Px
- (ii) Height \Rightarrow image height in Px
- (iii) Depth \Rightarrow Number of Channels (RGB=3, GRAY=1)

→ common pattern in CNN architectures are :



the general pattern of a CNN model is convolution layer followed by pooling layer

$\xrightarrow{\text{(will be repeating)}}$

→ CONVOLUTION LAYER:

Some of the important concepts in a convolution layer are:

- (i) Filters
- (ii) Strides
- (iii) Zero Padding
- (iv) Convolution operation

→ A convolution layer can contain a Set of filters of fixed size.

→ These filters are also called as Kernels. Two Hyperparameters are used in a convolution layer to control the filters / kernels they are.

- (i) number of filters
- (ii) filter size

→ Filters have common size like ~~different kernel sizes~~, 3×3 , 5×5 , 7×7 , 9×9 and 11×11 and number of filters used are generally 8, 16, 32, 64 and 128.

→ The process of applying the filter to the input of a convolution layer is called as convolution operation.

→ This operation is done in a sliding window fashion. The amount by which the filter should move is controlled by Stride value.

→ The main operation performed in convolution is DOT PRODUCT.

1	2	3				
4	5	6				
7	8	9				

3	1	5
2	4	7
4	1	1

FILTER OF SIZE 3x3

IMAGE OF SIZE 6x6

$$\Rightarrow 1 \times 3 + 2 \times 1 + 5 \times 3 + 4 \times 2 + 3 \times 4 + \\ 6 \times 1 + 7 \times 4 + 8 \times 1 + 9 \times 1$$

$$= 135$$

This will be repeated until whole image is convolved based on the filter values.

- The output obtained as a result of the convolution operation from a convolution layer is called as FEATUREMAP / ACTIVATIONMAP.
- Each of the filters present in a convolution layer will identify a specific feature from an image \rightarrow Filter extract feature from the image.

⑧ ZERO PADDING:

- Applying the convolution operation reduce the size of our image if image size is $N \times N$ and filter size is $F \times F$ the output image will be $(N-F+1) \times (N-F+1)$.

$$T_w - F_w + 1 = A_w$$

$$\Leftrightarrow (T_w + P) - F_w + 1 = T_w$$

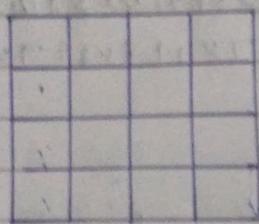
$$P = \frac{F_w - 1}{2}$$

DATE

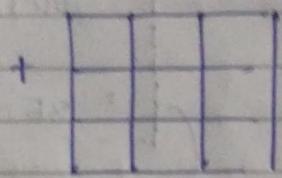
→ Two main problem we face is:

- (i) Reduction in the Size
- (ii) Losing valuable information on the edge.

→ So two overcome the above mentioned issue we use:
use zero padding.

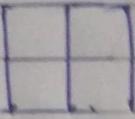


4x4 IMAGE



3x3 FILTER

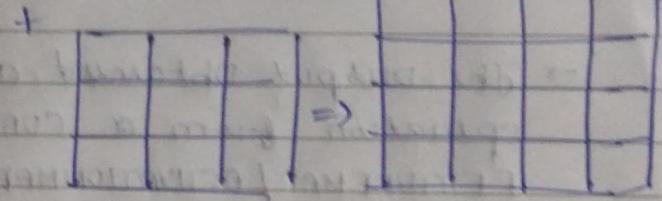
=>



2x2 ACTIVATION MAP

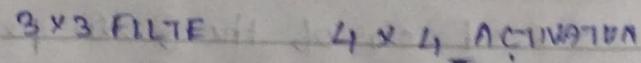
0	0	0	0	0	0
0					0
0	.	.			0
0					0
0					0
0	0	0	0	0	0

4x4 IMAGE + 0-PADS



3x3 FILTER

4x4 ACTIVATION



WHY RELU ACTIVATION IN CONVOLUTION LAYER ???

GIVEN:

Image Size: $N \times N$
Input Size: $N \times N$ ~~Input Size: $N \times N$~~

Filter Size: $F \times F$

Padding : P

Stride : S

→ Size of activation

$$\left(\frac{N - F + 2P + 1}{S} \right) \times \left(\frac{N - F + 2P + 1}{S} \right)$$

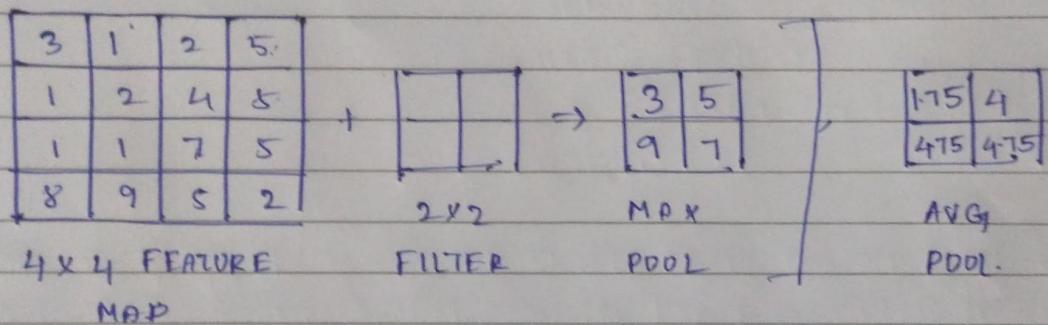
\Rightarrow POOLING LAYERS:

\rightarrow Pooling layers are inserted between successive convolution layers to reduce the size of activation map. gradually the process is called down Sampling

\rightarrow Two main types of pooling are.

- ① Max Pooling ② Average Pooling.

\rightarrow Zero-padding will not be used in pooling layer.



GIVEN:

Feature Map of Size: $N \times N$.

Filter Size: $F \times F$

Stride : S .

$$\Rightarrow \text{Output Size} = \left(\frac{N - F + 1}{S} \right) \times \left(\frac{N - F + 1}{S} \right)$$

* APPLICATIONS OF CNN:

- \rightarrow Image Classification
- \rightarrow NLP Applications
- \rightarrow Image Segmentation
- \rightarrow Image Generation
- \rightarrow Image Feature extraction

* SOME CNN ARCHITECTURE

- \rightarrow VGG16
- \rightarrow VGG19
- \rightarrow ResNet 50
- \rightarrow AlexNet
- \rightarrow LeNet