1.

a) False,

Naive bayes classifier assumes that each feature $x_i$ is conditionally independent of each other feature $x_j$ for $j \neq i$

b) ~~false~~ True

the success of pattern classification scheme using decission function depends on two factors :-

     i) The form of the decision function $d(x)$

     ii) ~~the~~ One's ability to determine the coefficients.

The first factor is governed by the geometrical properties of the classes under consideration.

c) ~~False, True,~~ False

     ~~The~~ In syntactic pattern recognition the set of ~~input~~ pattern primitives & ~~gos~~ the grammar is required for the classification.

Training phase is not an essential step.

d) False, The points that are ~~also~~ within the margin of the hyperplane can only be the support vectors.

e) True,

     Density - based clustering algorithms groups points together under one cluster when they occur in a high density region.

     Outliers generally occur in low ~~d~~ density regions, hence are not classified.

4). True,

Naive Bayes classifier assumes that the classes are conditionally independent.

For example consider a naive-bayes classifier to classify spam & important emails. The classifier will only consider the occurrence of certain key words in the mails in order to classify the mail. The ordering of the keywords is not words ~~pass~~ in not ~~in~~ paid importance, which is not true in real of life.

g) False,

~~sta~~ Syntactic Pattern Recognition attempts to classify ~~problems~~ patterns based on a set of extracted features called pattern primitives and their geometrical model ~~repres~~ represented through the grammar. ~~io~~

h) False,

Hierarchial clustering methods help in exploring data at different levels of ~~goa~~ granularity.

i) False,

A Hopfield net is mainly used for optimization.

j) False,

For a supervised pattern classification problem having M-classes ~~with~~ where the classes are pairwise seperable, the classifier needs to compute [M(M-1)]/2 number of decission surfaces.

k). False,

Data processing is required to ensure
- accuracy, consistency, completeness. timeliness, believability &
  &  is  interpretability.

2.

M - patterns.

prototypes : $Z_1, Z_2, \ldots Z_M$.     classes : $W_1, W_2, W_3 \ldots, W_M$

~~we don determine the~~

we use euclidean distance of a ~~g~~ point from any given prototype as
a measure of similarity of that point to ~~a particular~~ the particular
class belonging to the prototype.

$$D_i = \| x - z_i \| = \sqrt{(x - z_i)'(x - z_i)}$$

Minimum - distance classification is used here, i.e. a given point $x$
is assigned to the ~~to~~ class that has the minimum euclidean distance
from it.

   $x$ is assigned to class $w_i$ if $D_i < D_j \ \forall \ j \neq i$.

   (Ties are resolved arbitrarily)

$D_i^2 = \| x - z_i \|^2 = (x - z_i)'(x - z_i)$

   $= x'x - x'z_i - z_i'x + z_i'z$

   $= x'x - 2(x'z_i - \frac{1}{2} z_i'z)$

we need to choose $i$ such that $D_i$ is minimum
$\Rightarrow D_i^2$ is minimum     (as $D_i$ is always a +ve value)
$\Rightarrow$ choosing $i$ such that $(x'z_i - \frac{1}{2} z_i'z)$ is maximum.
         as $x'x$ is independent of $i$.

∴

∴ we can define decision function as follows:-

$$d_i(x) = x'z_i - \frac{1}{2} z_i' z \quad, i = 1, 2, \dots M$$

when $X$ is assigned to the class $W_i$

when $d_i(x) > d_j(x) \quad \forall i \neq j$

we can see that $d_i(x)$ is a linear decission function.

3.

Baye's theorem describes the probaility of an event, based on prior knowledge of conditions that might be related to the event.

Mathematically :-

Let $X$ be a data sample

Let ✶ H be a hypothesis that $X$ belongs to a class $C$.

$P(H) \rightarrow$ initial provability

$P(x) \rightarrow$ (evidence) probability that sample data is observed

$P(x|H) \rightarrow$ (likelyhood) the probability of observing $X$ given that the hypothesis holds.

$$P(H|x) = \frac{P(x|H) \, P(H)}{P(x)}$$

Let $D$ be a set of N tuples $\{x_1, x_2, \dots x_N\}$

when $x_i \in \mathbb{R}^n$.

$\rightarrow$ class K-classes $C_1, C_2, \dots C_K$.

∴ From bayes theorem

$$P(C_k|x) = \frac{P(x|C_k) \, P(C_k)}{P(x)} \qquad P(C_i|x) = \frac{P(x|C_i) \, P(C_i)}{P(x)}$$

Since the denominator is not dependent on $C_i$ and all values of the feature features are known, ∴ The denominator is effectively constant.

∴ Only numerator is significant, which is equivalent to the joint probability model $P(C_k, x_1, \ldots x_n)$

$$P(C_i) P(x|C_i) = P(C_k, x)$$

or $P(C_i) P(x_1, x_2, \ldots x_n | C_k) = P(C_k, x_1, v_2 \ldots x_n)$

with $[x = (x_1, x_2, x_3 \ldots x_n)]$

Using chain rule for repeated application of definition of conditional probability:

$$P(C_i, x_1, u_2 \ldots x_n) = P(x_1, x_2 \ldots, x_n, C_i)$$

$$= P(x_1 | x_2, \ldots x_n, C_i) P(x_2, x_3 \ldots x_n, C_i)$$

$$\vdots$$

$$= P(x_1 | x_2 \ldots x_n, C_i) P(x_2 | x_3, \ldots x_n, C_i) \ldots P(x_{j-1} | x_j, \ldots x_n, C_i)$$

$$\ldots P(x_{n-1} | x_n, C_k) P(x_n | C_k) P(C_k)$$

Naive Bayes assumes that each feature $x_i$ is conditionally independent of every other feature $x_j$ for $i \neq j$.

$$\Rightarrow P(x_j | x_{j+1}, \ldots x_n, C_i) = P(x_j | C_i)$$

Then the joint model can be expressed as :-

$$P(C_i, x_1, \ldots x_n) \propto P(C_k, x_1, \ldots x_n)$$

$$= P(C_k) P(x_1 | C_k) P(x_2 | C_k) P(x_3 | C_i)$$

$$P(C_i, x_1, x_2 \ldots x_n) = P(C_i) \cdot P(x_1 | C_i) P(x_2 | C_i) \ldots P(x_n | C_i)$$

$$= P(C_i) \prod_{j=1}^{n} P(x_j | C_i)$$

In order to classify $X$, we pick the hypothesis that is most probable, . ee.

∴ The corresponding classifier, a Bayes classifie, is the function that assigns a class label $\hat{y} = c_i$ for some $i$ as follows :

$$\hat{y} = \underset{i \in \{1, \dots k\}}{\text{argmax}} \; P(c_i) \prod_{j=1}^{n} P(x_j | c_i) .$$

Let $A_j$ denote the $j$th feature of a given data sample $X$.

Now $A_j$ can be either categorical or continous valued.

$P(x_j | c_i)$ has to be computed differential differently for the avowe mentiond two cases.

If $A_j$ is categorical, $P(x_j | c_i)$ is the ∞ number of tuples in $c_i$ having value $x_j$ of $A_j$ divided by $6$ $|c_i|$ (s $|c_{i,D}|$) |ets, number of tuples of $c_j$ $c_i$ in $D$.

If $A_j$ is continous-valued, $P(x_j | c_i)$ is usually computed based on a Gaussian distribution with a mean $\mu$ and $\emptyset$ SD $\sigma$.

$$g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi \sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$
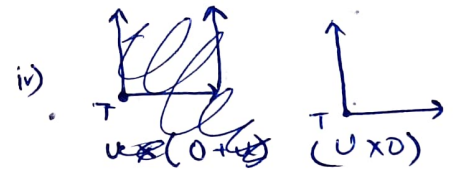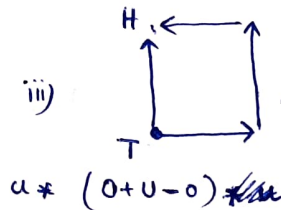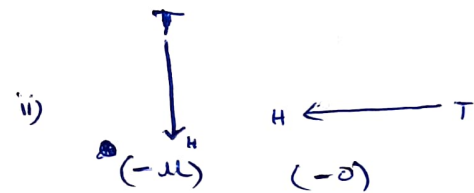
so, $P(x_j | c_i) = g(x_j, \mu_{c_i}, \sigma_{c_i})$

5.


A C P F
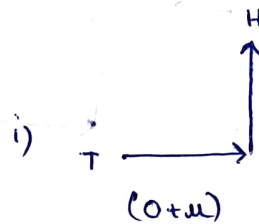
**Primitives required :-**

$\{0, u\}$

$T \xrightarrow{0} H$


H
↑
| u
T

**Operation allowed on primitives :-**

$\{ *, -, +, X \}$.

i)

$T \longrightarrow$ H
$(0+u)$

ii)

$(-u)$   $H \longleftarrow T$   $(-0)$

iii)

$u * (0 + U - 0)$

iv)

$u * (0 + u)$   $(U \times 0)$

* $\longrightarrow$ represents head - head & tail-tail attachment.

+ $\longrightarrow$ represents head to tail concatenation

- $\longrightarrow$ represents head-tail reversal

X $\longrightarrow$ represents tail-tail attachment.

The grammar consists of :

1. Start symbol S

2. Set of terminals $\{ u, 0, *, -, +, X, 1, (, ) \}$

    ( where "1" is or-operator, &, "(" & ")" are opening & closing parentheses).

3. Set of non-terminals $\rightarrow \{ S, A, C, P, F \}$.

4. Set of production values $P = \{$   $S \rightarrow A | C | P | F$,

            $A \rightarrow u + ((0u + 0 + -u) * 0) + - u$,
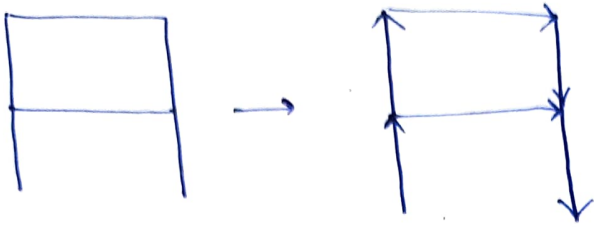
            $C \rightarrow -0 + u + u + 0$,

            $P \rightarrow u + ((u + 0 + -u) * 0)$,

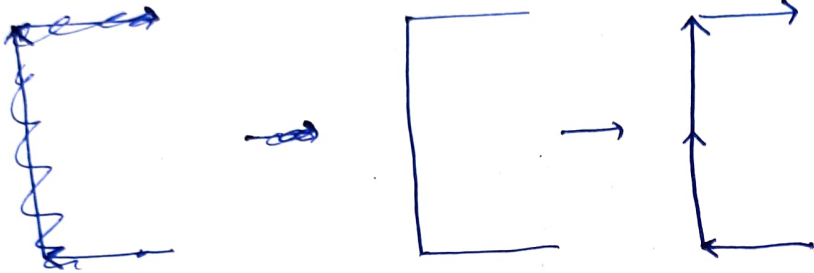            $F \rightarrow u + (0 \times u) + 0$

     $\}$

Premitive representation & interconnections :-



$A = u + ((u + 0 + -u) * 0) + -u$



$C = -0 + u + u + 0$



$P = u + ((u + 0 + -u) * 0)$



$F = u + (0 \times u) + 0$

6. Pattern recognition is the automated recognition of patterns and regularities in data.

It can be divided into two use cases:

i) Recognizing concrete items.

   Ex - Pictures.

ii) Recognizing abstract items.

   Ex - voical audio.

Most conventional approaches of pattern recognition are based on direct computation through machines which are math-related techniques. These conventional approaches include :- feature extraction, classification, clustering etc.

We can also use application of biological concets of neurons inside in humans for computing. This lead to the development of neural networks.

ANN (Artifical Neural Networks) :-

An Artificial Neural Network is a paralleled distributed information processing structures in the form of a directed graph.

It consists of a larger number of simple processing units (perceptors) with a high degree of dist arranged into layers one or more layers with a high degree of interconnection between each layer of units.

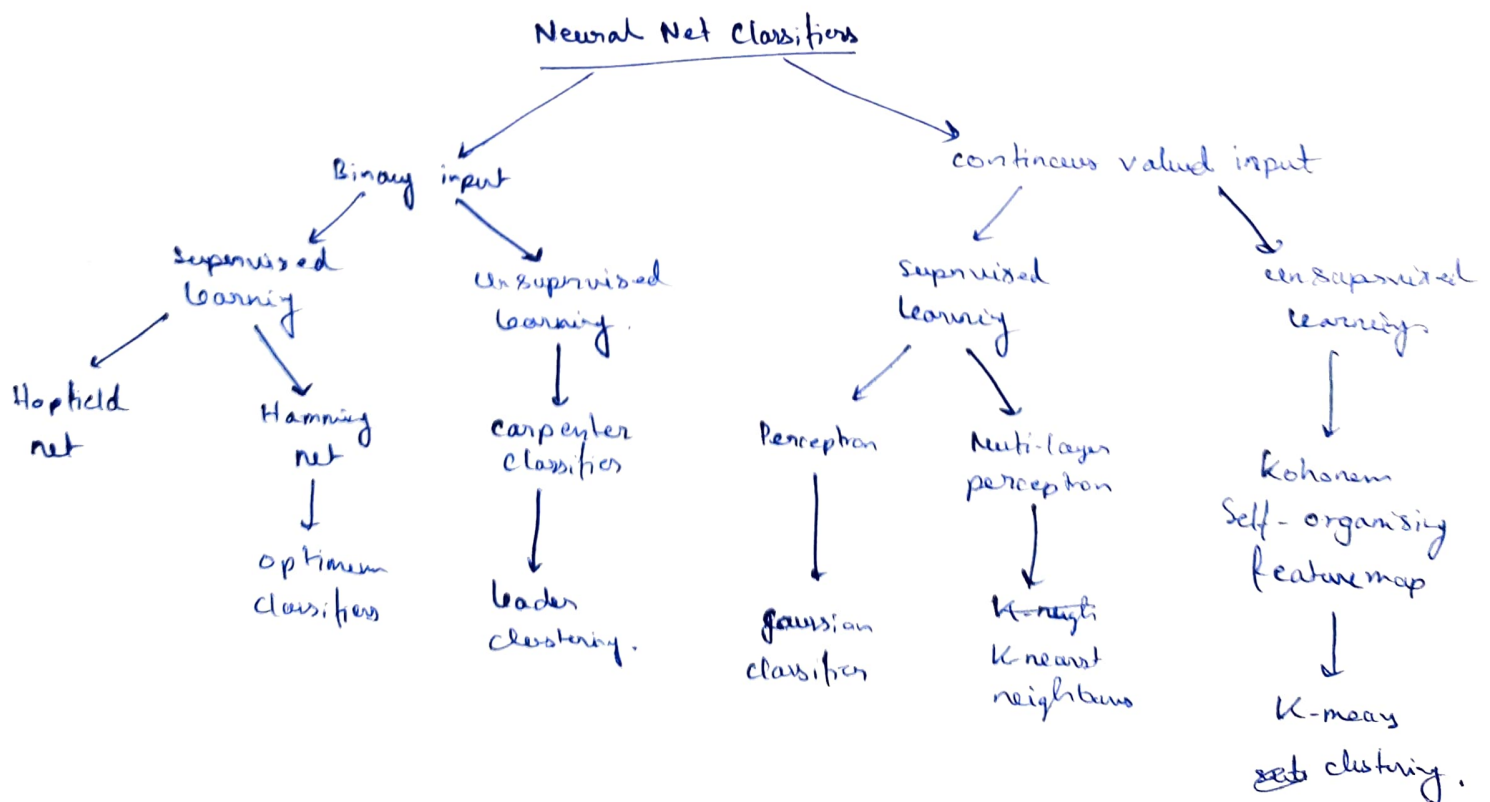The processing units works parallelly and in coordination with each other.

The design & function of neural networks simulate some functionality of biological brains and neural system.

Pattern recognition can be done using both conventional computers and neural networks.

→ Neural network require simple processors as opposed to general computers which use few complex processes.

→ Neural networks are fewer processing steps.

→ Neural networks use the concept of distributed processing making them faster.

→ Neural nets are trained by example helping them to achieve far better results even for unknown data.

→ Neural nets are tolerable to noisy patterns.

Due to the adaptive-learning, self organizing and fault tolerance capabilities of neural nets, ANNs are used for various pattern recognition applications.

Neural networks as classifiers:-

Neural Net Classifiers

Binary input

continuous valued input

Supervised learning

Unsupervised learning

Supervised learning

unsupervised learning

Hopfield net

Hamming net

carpenter classifier

Perceptron

Multi-layer perceptron

Kohonen Self-organising feature map

optimum classifier

leader clustering

gaussian classifier

K-nearst neighbours

K-means clustering

We can see the importance of ANN in Pattern Recognition by looking at the diversity of applications that ANN's have In Pattern Recognition problems.

| Algorithms | Type | usage |
|---|---|---|
| Hopfield | recursive | optimization |
| Multi-layer perceptron | feed forward | classification |
| Kohonan | self - organising | data - coding |
| Temporal differences | predictive | fore casting. |

∴ Neural networks provides the following advantages:-

- Can work with incomplete information once trained
- Fault tolerance (robust to outlies)
- Distributed
- Parallel
- Can learn non-linear and complex relationships
- trained by example
- Generalizability & i.e. can infer unseen relationships once trained.