

1. Describe the Find-S algorithm and candidate elimination algorithm, explain by taking enjoy sport concept and training instance given below.

Example	sky	AirTemp	Humidity	wind	water	Fore cast	Enjoy sport
1.	sunny	warm	Normal	strong	warm	same	yes
2.	sunny	warm	high	strong	warm	same	yes
3.	Rainy	cold	high	strong	warm	change	no
4.	sunny	warm	high	strong	cool	change	yes

A. Find-S algorithm:

- This algorithm considers only positive examples
- Most specific hypothesis  $\Rightarrow \phi$
- Most general hypothesis  $\Rightarrow ?$

Algorithm:

step 1: Initialize  $h$  with most specific hypothesis ( $\phi$ ) in  $H$   
 $h_0 = \langle \phi, \phi, \phi, \phi, \phi \rangle \Rightarrow 5$  attributes (depends on attributes)

step 2: For each positive sample,

For each attribute,

if (value = hypothesis value)  $\Rightarrow$  ignore  
 else

Replace with the most general hypothesis (?)

step 3: output hypothesis  $h$

Example:

Applying Find-S algorithm on given Example enjoy sport

step 1:  $h_0 = \langle \phi, \phi, \phi, \phi, \phi, \phi \rangle$

step 2:  $h_1 = \langle \text{sunny, warm, normal, strong, warm, same} \rangle$

step 3:  $h_2 = \langle \text{sunny, warm, ?, strong, warm, same} \rangle$

step 4:  $h_3 = h_2$

step 5:  $h_4 = \langle \text{sunny, warm, ?, strong, ?, ?} \rangle$

Disadvantage:

- considers only positive (yes) values
- $h_4$  may not be sole hypotheses that fits the complete data.

## Candidate Elimination Algorithm:

- It uses the concept of version space
- It considers both positive and negative values (Yes & No)
- Both specific and general hypothesis.
- For positive samples, move from specific to general
- For negative samples, move from general to specific
- Specific values  $\langle \phi, \phi, \phi, \phi, \phi \rangle \rightarrow$  positive
- General uses  $\langle ?, ?, ?, ?, ? \rangle \rightarrow$  Negative

### Algorithm:

Step 1: Initialize both general and specific hypothesis ( $S$  and  $G$ )

$$S = \langle \phi, \phi, \phi, \phi, \phi \rangle$$

$$G = \langle ?, ?, ?, ?, ? \rangle \text{ depends on number of attributes}$$

Step 2: For each example,

if example is positive

Make specific to general

if example is negative

Make general to specific

### Example:

Applying candidate elimination algorithm on given example enjoy sport.

$$S_0 = \langle \phi, \phi, \phi, \phi, \phi, \phi \rangle \quad G_0 = \langle ?, ?, ?, ?, ?, ? \rangle$$

1. For positive:  $S_1 = \langle \text{'sunny'}, \text{'warm'}, \text{'normal'}, \text{'strong'}, \text{'warm'}, \text{'same'} \rangle$

$$G_1 = \langle ?, ?, ?, ?, ?, ? \rangle$$

2. For positive:  $S_2 = \langle \text{'sunny'}, \text{'warm'}, ?, \text{'strong'}, \text{'warm'}, \text{'same'} \rangle$

$$G_2 = \langle ?, ?, ?, ?, ?, ? \rangle$$

3. For negative:  $S_3 = \langle \text{'sunny'}, \text{'warm'}, ?, \text{'strong'}, \text{'warm'}, \text{'same'} \rangle$

$$G_3 = \langle \langle \text{'sunny'}, ?, ?, ?, ?, ? \rangle, \langle ?, \text{'warm'}, ?, ?, ?, ? \rangle, \langle ?, ?, ?, ?, ? \rangle, \text{'same'} \rangle$$

4. For positive:  $S_4 = \langle \text{'sunny'}, \text{'warm'}, ?, \text{'strong'}, ?, ? \rangle$

$$G_4 = \langle \langle \text{'sunny'}, ?, ?, ?, ?, ? \rangle, \langle ?, \text{'warm'}, ?, ?, ?, ? \rangle \rangle$$

∴ therefore  $S_4$  and  $G_4$  are final hypothesis.

2. Give Decision tree for the following set of training examples. calculate entropy, information gain and gain index.

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Steps: Create a root Node

→ The attribute that best classifies the training data, use this attribute at the root of the tree

→ calculate entropy

$$\text{entropy} = -\frac{p}{p+n} \log_2\left(\frac{p}{p+n}\right) - \frac{n}{p+n} \log_2\left(\frac{n}{p+n}\right)$$

→ calculate average information

$$\Delta(\text{attribute}) = \sum \frac{p_i+n_i}{p+n} \text{entropy}(A)$$

→ calculate Information gain

$$\text{Gain} = \text{entropy}(S) - \Delta(\text{attribute})$$

Building decision tree for a given example

$$P=9 \quad N=5 \quad \text{Total}=14$$

$$\text{Entropy} = -\frac{P}{P+N} \log_2 \left( \frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left( \frac{N}{P+N} \right)$$

$$\text{entropy}(S) = -\frac{9}{9+5} \log_2 \left( \frac{9}{9+5} \right) - \frac{5}{9+5} \log_2 \left( \frac{5}{9+5} \right)$$

$$\text{Entropy}(S) = 0.940$$

→ For each attribute: (let say outlook)  
calculate entropy for each values i.e for 'Sunny', 'Rainy', 'overcast'

outlook	p	n	entropy
sunny	2	3	0.971
Rainy	3	2	0.971
overcast	4	0	0

calculate entropy (outlook = 'value'):

$$E(\text{outlook} = \text{sunny}) = -\frac{2}{5} \log_2 \left( \frac{2}{5} \right) - \frac{3}{5} \log_2 \left( \frac{3}{5} \right) = 0.971$$

$$E(\text{outlook} = \text{overcast}) = -1 \log_2 (1) - 0 \log_2 (0) = 0$$

$$E(\text{outlook} = \text{rainy}) = -\frac{3}{5} \log_2 \left( \frac{3}{5} \right) - \frac{2}{5} \log_2 \left( \frac{2}{5} \right) = 0.971$$

→ calculate Average Information Entropy

$$I(\text{outlook}) = \frac{P_{\text{sunny}} + n_{\text{sunny}}}{P+N} \text{entropy}(\text{outlook} = \text{sunny}) + \frac{P_{\text{rainy}} + n_{\text{rainy}}}{P+N} \text{entropy}(\text{outlook} = \text{rainy}) + \frac{P_{\text{overcast}} + n_{\text{overcast}}}{P+N} \text{entropy}(\text{outlook} = \text{overcast})$$

$$I(\text{outlook}) = \frac{2+2}{9+5} * 0.971 + \frac{2+3}{9+5} * 0.971 + \frac{4+0}{9+5} * 0 = 0.693$$

→ calculate gain: attribute is outlook

$$\text{Gain}(\text{outlook}) = 0.940 - 0.693 = 0.247$$

→ For each attribute (let say temperature)

temperature	p	n	entropy
Hot	2	2	1
Mild	4	2	0.918
cool	3	1	0.811

→ calculate average information entropy

$$I(\text{Temperature}) = \frac{3}{15} \log_2 \frac{3}{15} + \frac{11}{15} \log_2 \frac{11}{15} = 0.918$$

→ calculate Gain: attribute is Temperature

$$\text{Gain}(\text{Temperature}) = 0.940 - 0.918 = 0.022$$

→ for each attribute: let say humidity

Humidity	P	n	entropy
high	2	4	0.722
normal	6	1	0.585

→ calculate average information entropy

$$I(\text{Humidity}) = \frac{3}{15} \log_2 \frac{3}{15} + \frac{6}{15} \log_2 \frac{6}{15} + \frac{6}{15} \log_2 \frac{6}{15} = 0.988$$

→ calculate gain: attribute is humidity

$$\text{Gain}(\text{Humidity}) = 0.940 - 0.988 = 0.052$$

→ for each attribute: let say windy

windy	P	n	entropy
strong	3	3	1
weak	6	2	0.918

→ calculate Average Information entropy

$$I(\text{windy}) = \frac{3}{15} \log_2 \frac{3}{15} + \frac{6}{15} \log_2 \frac{6}{15} + \frac{6}{15} \log_2 \frac{6}{15} = 0.918$$

→ calculate Gain: attribute is windy

$$\text{Gain}(\text{windy}) = 0.940 - 0.918 = 0.022$$

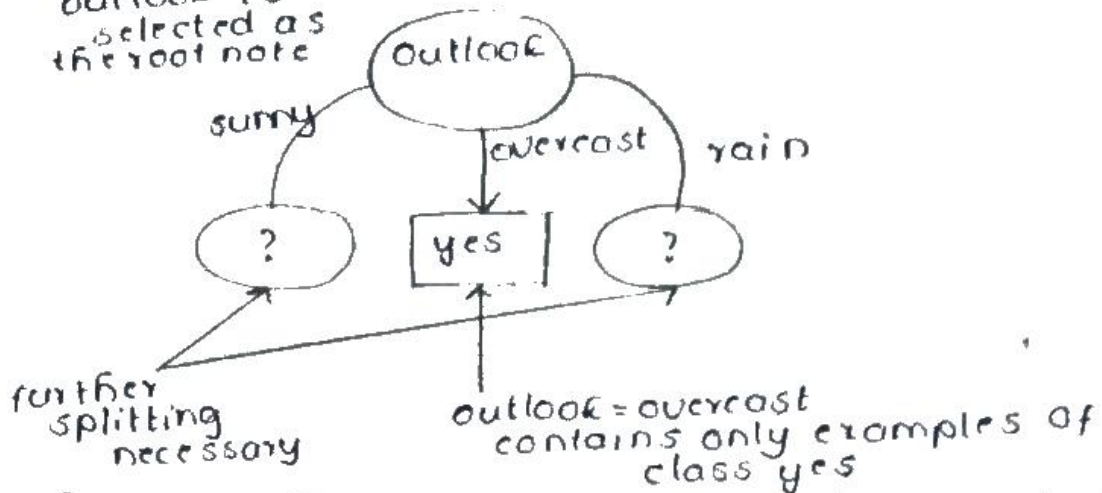
→ Pick the highest gain attribute

Attributes	Gain
cutlook	0.247
temp probab	0.029
humidity	0.152
windy	0.048

ROOT NODE: cutlook

outlook	Temperature	Humidity	windy	play tennis
overcast	hot	high	weak	yes
overcast	cool	normal	strong	yes
overcast	mild	high	strong	yes
overcast	hot	normal	weak	yes

outlook is selected as the root node



→ Repeat the same thing for subtrees till we get the tree  
 outlook = sunny

$$P = 2, N = 3, \text{Total} = 5$$

$$\text{entropy}(s_{\text{sunny}}) = \frac{-2}{2+3} \log_2 \left( \frac{2}{2+3} \right) - \frac{3}{2+3} \log_2 \left( \frac{3}{2+3} \right) = 0.971$$

→ for each attribute: (let say humidity)

→ calculate entropy for each humidity, i.e. for 'high', 'normal'

Humidity	P	N	Entropy
high	0	3	0
normal	2	0	0

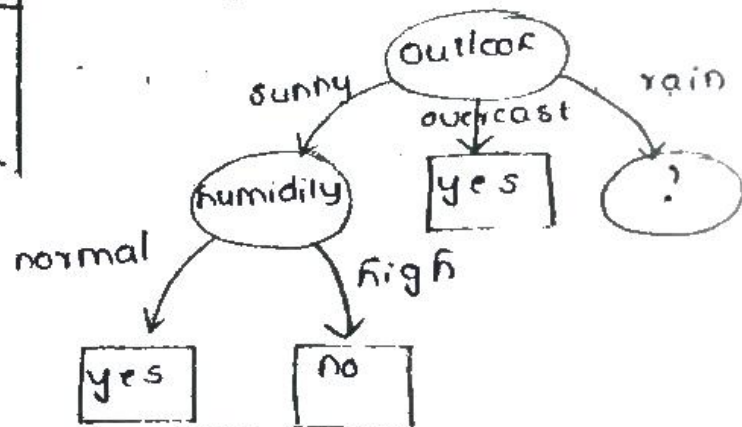
$$I(\text{humidity}) = 0$$

$$\text{Gain} = 0.971$$

→ pick the highest gain attribute

Attributes	Gain
Temperature	0.571
humidity	0.971
windy	0.02

Next Node: insunny-humidity



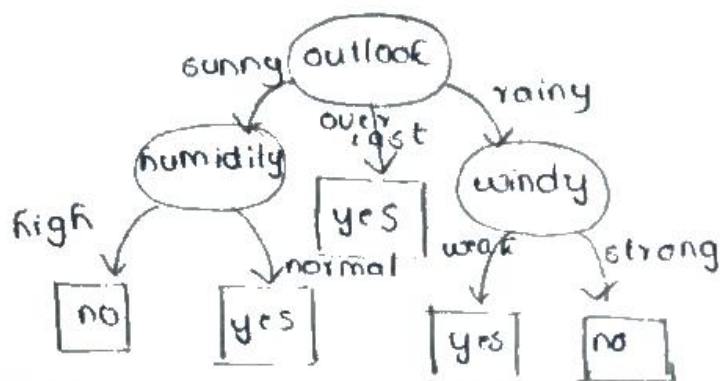
No further expansion necessary

→ Repeat the same process and the gain we get is

Attributes	Gain
Humidity	0.02
windy	0.971
Temperature	0.02

Next Node in  
 Rainy:  
 windy

Final decision tree



With a neat sketch explain how the two hyper planes in SUPPORT VECTOR MACHINE works.

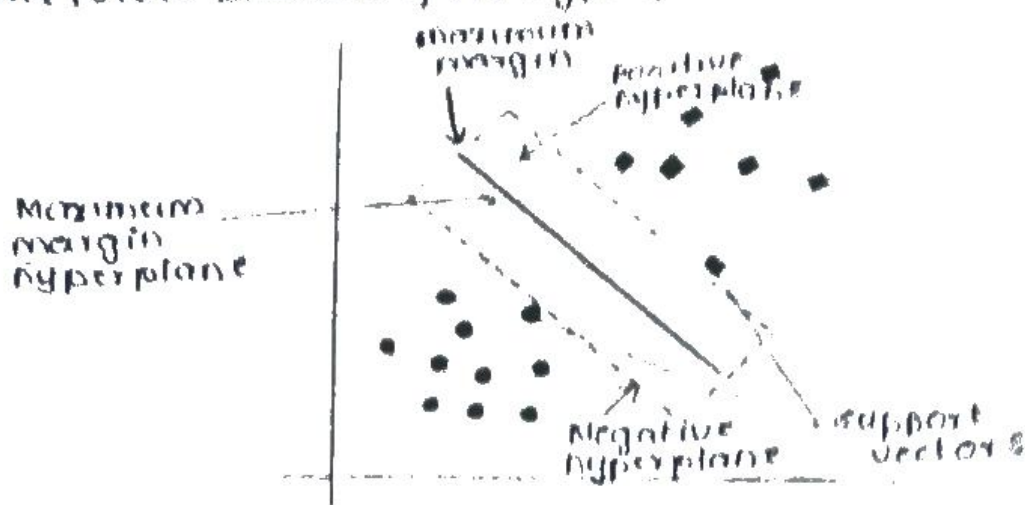
Support vector machine:

Support vector machine is one of the most popular supervised learning algorithms, which is used for classification as well as regression problem. However, it is used for classification problems in machine learning.

The goal of SVM algorithm is to create the best line that can segregate n-dimensional space into classes so that we can easily put the new data point in the correct category in the future. This best decision boundary is called hyper plane.

SVM chooses the extreme points/vectors that help in creating the hyperplane. These extreme cases are called as support vectors and hence algorithm is termed as support vector machine.

consider below diagram in which there are two different categories that are classified using a decision boundary or hyperplane



Hyperplane and support vectors in the SVM algorithm:  
Hyperplane:

There can be multiple lines/decision boundaries to segregate the classes in  $n$ -dimensional space, but we need to find out the best decision boundary that helps to classify the data points. This best boundary is known as the hyperplane of SVM.

→ The dimensions of hyperplane depend on the features present in the dataset, which means if there are 2 features, then hyperplane will be a 2-dimension plane.

→ we always create a hyperplane that has a maximum margin, which means the maximum distance between the data points.

Support vectors:

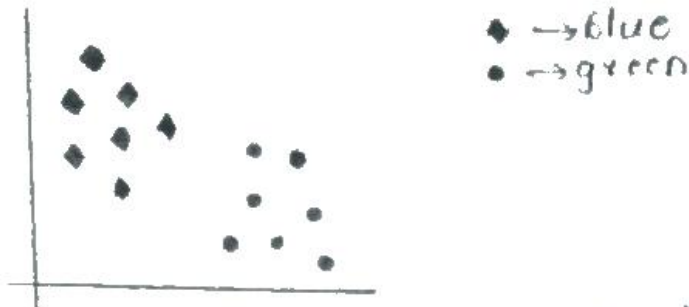
The data points or vectors that are the closest to the hyperplane and which affect the position of the hyperplane are termed as support vector since these vectors support the hyperplane, hence called a support vector.

How does SVM works:

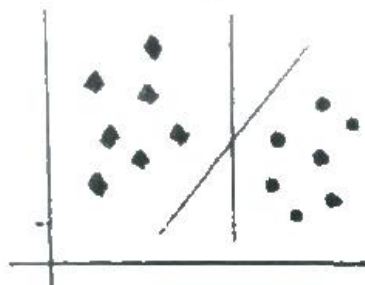
The working of the SVM algorithm can be understood by using an example. Suppose we have a dataset that has two tags (green and blue), and the dataset has two features  $x_1$  and  $x_2$ . we want



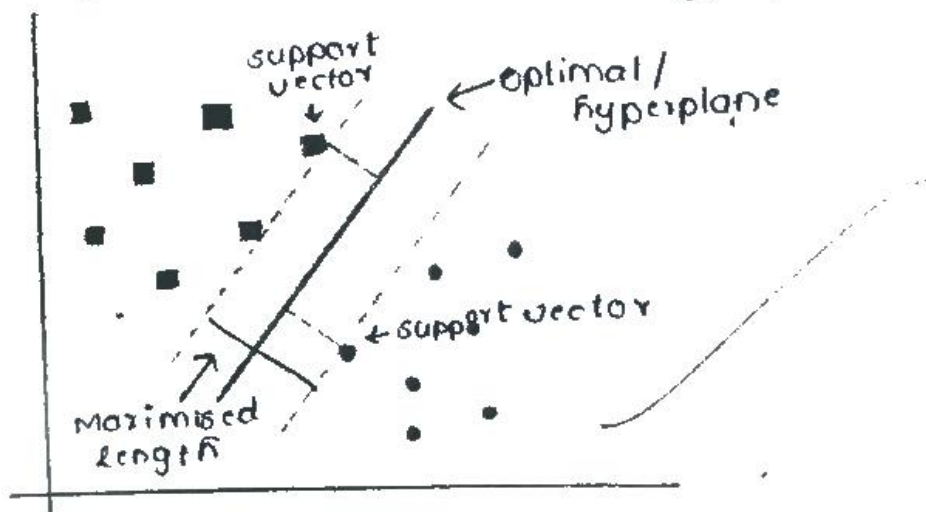
classifier that can classify the pair  $(x_1, x_2)$  of coordinates in either green or blue. consider the below image.



So as it is 2d-space so by just a straight line. we can easily separate these two classes. But there can be multiple lines that can separate these classes. consider the below image



Hence, the SVM algorithm helps to find the best line or decision boundary, this best boundary or region is called as a hyperplane. SVM algorithm finds the closest point of the lines from both the classes. These points are called support vectors. The distance between the vectors and the hyperplane is called as margin. And the goal of SVM is to maximize this margin. The hyperplane with maximum margin is called the optimal hyperplane.



4. Enumerate in detail how Ada Boost algorithm works.

AdaBoost:

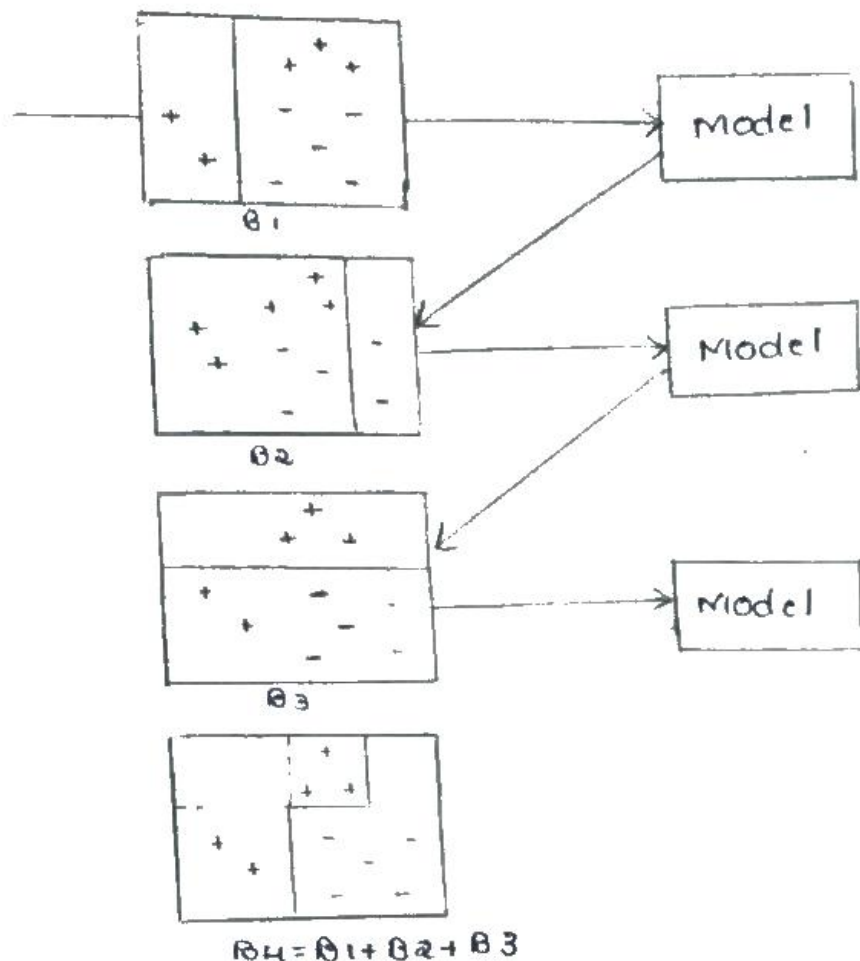
Boosting is an ensemble modeling technique that attempts to build a strong classifier from the number of weak classifiers. It is done by building a model by using weak models in series. Firstly, a model is built which tries to build from training data. The second model is built which tries to correct the errors present in the first model. This procedure is continued and models are added until either the complete training data set is predicted correctly or the maximum number of models are added.

Ada Boost was the first really successful boosting algorithm developed for the purpose of binary classification. AdaBoost is short for Adaptive Boosting and is a very popular boosting technique that combines multiple "weak classifiers" into a single "strong classifier". It was formulated by Yoav Freund and Robert Schapire. They also won the 2003 Godel Prize for their work.

Algorithm:

1. Initialise the dataset and assign equal weight to each of the data point.
2. provide this as input to the model and identify the wrongly classified data points
3. Increase the weight of the wrongly classified data points
4. if (got required results)  
    goto step 5  
else  
    goto step 2

Example:



Explanation:

The above diagram explains the AdaBoost algorithm in a very simple way. Let's try to understand it in a stepwise process:

- $B_1$  consists of 10 data points which consist of two types namely plus(+) and minus(-) and 5 of which are plus(+) and the other 5 are minus(-) and each one has been assigned equal weight initially. The first model tries to classify 3 plus(+) as minus(-)
- $B_2$  consists of the 10 data points from the previous model in which the 3 wrongly classified plus(+) are weighted more so that the current model tries more to classify these plus(+) correctly. This model generates a vertical separator line that correctly classifies the

previously wrongly classified pluses(+) but in this attempt, it wrongly classifies three minuses(-).

→ B3 consists of the 10 data points from the previous model in which the 3 wrongly classified minuses(-) are weighted more so that the current model tries more to classify these minuses(-) correctly. This model generates a horizontal separator line that correctly classifies the previously wrongly classified minuses(-).

→ B4 combines together B1, B2 and B3 in order to build a strong prediction model which is much better than any individual model used.

Verdict

MDU  
2/3/2023