

Assignment 5

Q.1 Discuss the need of SVM? Explain why they are called as optimal binary classifier?

→ ~~in next page : main points~~

- Support Vector Machine (SVM) is a type of supervised learning that can be used for classification or regression.

- Even if the data points are unseen, SVM can classify the data properly.

Q - The need of SVM :-

1) Effective in high dimensional spaces

→ SVMs perform well in case of where the number of features is large compared to the number of samples.

- This is particularly useful in areas like text classification, image recognition, and bioinformatics, where data often has many features.

Q 2) Robust to overfitting

→ SVMs aims to find the optimal separating boundary between classes with a maximum margin, making them less prone to overfitting, especially in high-dimensional spaces.

3) Handles Outliers

→ SVMs are relatively robust to outliers due to use of a margin and support vectors, focusing only on the most important data points near the boundary.

A financial plan

- 4) Well-suited for binary classification
→ SVMs are particularly effective for binary classification problems, as they aim to maximise (the margin between two classes, resulting in a clear and strong decision boundary.

- Why SVMs are called as Optimal Binary Classification
→ MVA in basic form
 - SVMs are termed as optimal binary classifier because they focus on finding the most effective way to separate two classes by constructing an optimal hyperplane.
 - This hyperplane is positioned in such a way that it maximizes the margin between the two classes with this margin being the distance between the hyperplane and the closest data points from each class, known as support vectors.
 - By maximizing this margin, SVM ensures that the model is not only accurate but also generalizes well to new data by minimizing the risk of overfitting.
 - Thus, SVM are known for their optimality in achieving both high accuracy and good generalization, especially in binary classification problem.

Q.2 Explain the following terminologies with the help of appropriate illustrations.

i) Optimal Decision Boundary

⇒ ~~Decision Rule~~

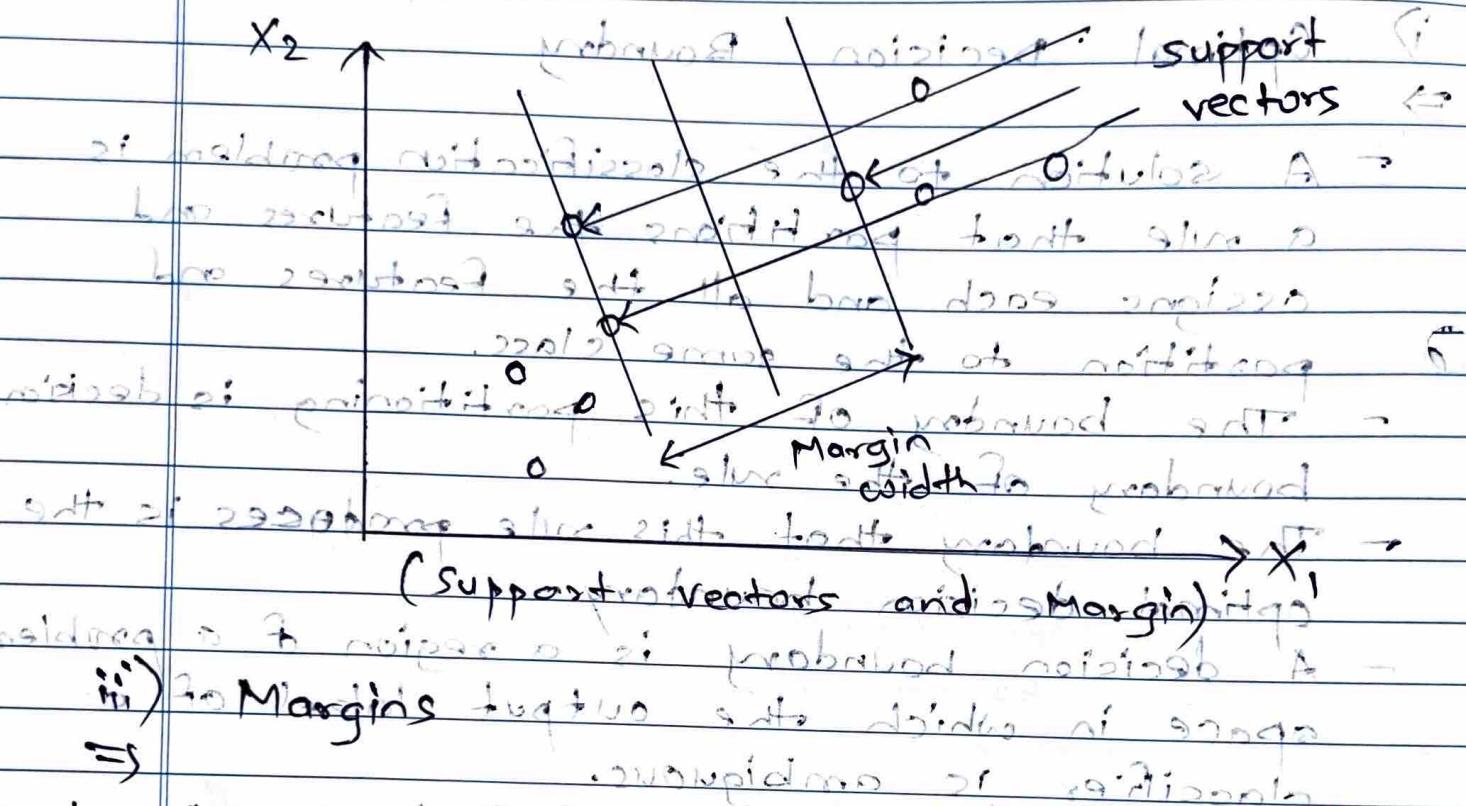
- A solution to the classification problem is a rule that partitions the features and assigns each and all the features and position to the same class.
- The boundary of this partitioning is decision boundary of the rule.
- The boundary that this rule produces is the optimal decision boundary.
- A decision boundary is a region of a problem space in which the output label of a classifier is ambiguous.
- If the decision surface is a hyperplane, then classification is linear and the classes are linearly separable.

ii) Support Vectors

⇒ ~~A vector in a n-dimensional space~~

- Support vectors are the data points that lie closest to the decision boundary or hyperplane.
- These points are critical because they determine the position of the optimal hyperplane.
- The decision boundary is determined based only on these support vectors, while other points do not influence its position.

Support vectors are essential in SVM because they help maximize the margin between two classes.



(Support vectors and Margin)

iii) Margins divide the data into two classes correctly.

- Margin represents the distance between the decision boundary and the nearest data points of each class.
- SVM aims to maximize margin, which leads to a robust decision boundary.
- A larger margin implies a more confident and generalizable classifier, reducing the risk of overfitting.
- The margin is measured from the support vectors to the decision boundary, and maximizing it ensures the model's optimal performance.

Q.3 What do you understand by linear classifiers and non-linear classifiers? Can you use SVM as non-linear classifier? If yes explain how SVM can be used as non-linear classifier.

⇒ This is a weak form of the question.

- Linear classifier

- ⇒ A linear classifier is a classification algorithm that makes its predictions based on linear combination of input features.
- ⇒ The decision boundary it creates to separate different classes is a straight line or a plane, or a hyperplane.

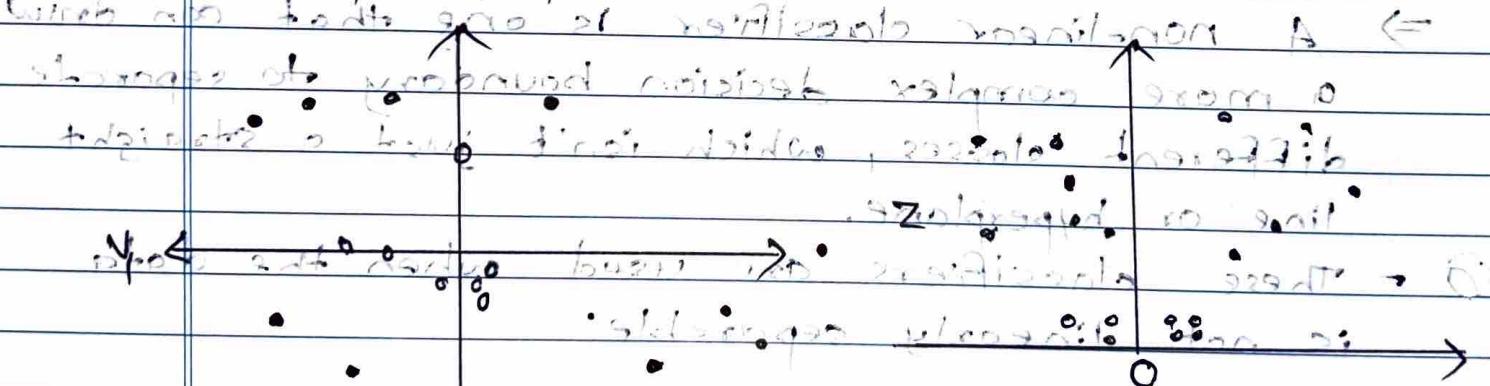
- Non-linear classifier

- ⇒ A non-linear classifier is one that can draw a more complex decision boundary to separate different classes, which isn't just a straight line or hyperplane.
- ⇒ These classifiers are used when the data is not linearly separable.

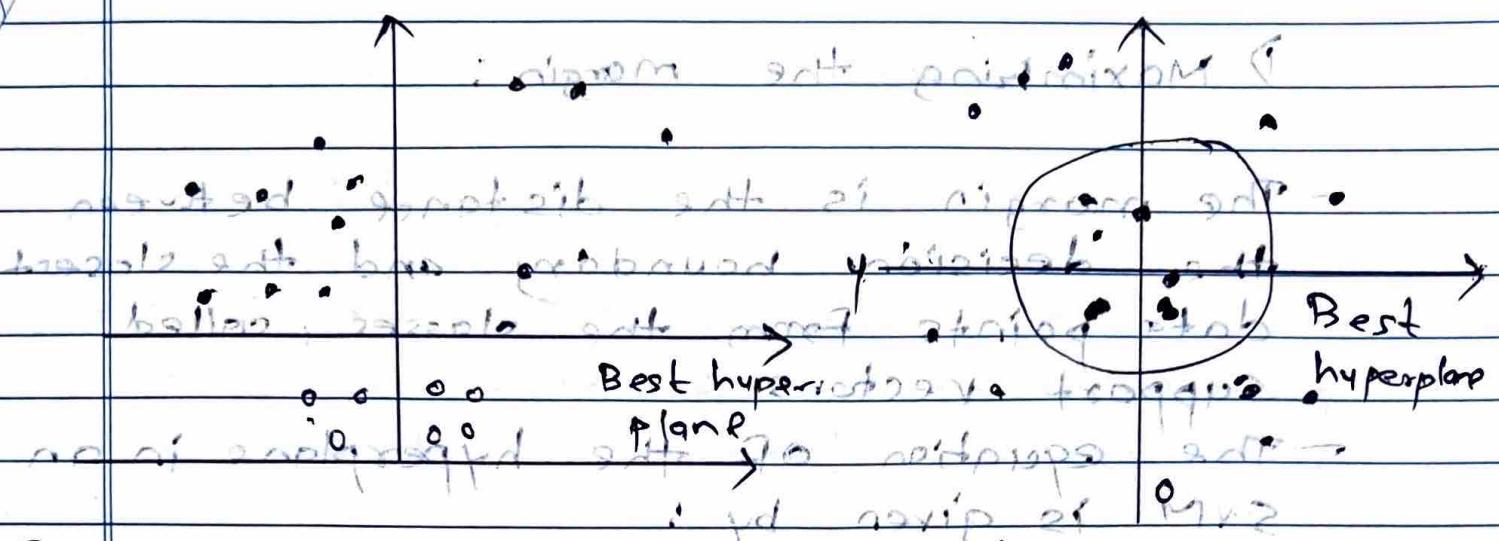
- Can SVM be used as Non-linear classifier?

- ⇒ Yes, SVM can be used as non-linear classifier.
- ⇒ Although the basic SVM algorithm is linear classifier, it can be extended to handle non-linear classification tasks through a method called the kernel trick.

- Non-linear SVM is based on the fact that if the data is non-linearly separable, we can't separate it by using a straight line, but if data is non-linear, then we cannot draw a single straight line.
- So to separate these data points we need one more dimension, i.e. add one more dimension to it.
- For linear data, we have used two dimensions, but for non-linear data, we have added a third dimension.
- It can be calculated using $z = xy^2$.
- After adding the third dimension, the sample space becomes:



Positional constraints on how we can position our data points on the axes and in the 3D space. Now, SVM will divide the datasets into classes in the following way:



Q - If we put $z=1$, then it will become as above, $0 = d + \text{margin}$
 $x^2 + y^2 = 1$.

We get a circumference of radius 1 as hyperplane surface? (in case of non-linear data). margin = 10

Q.4 Express SVM as a constrained optimization problem. Discuss how predictions can be done by using SVM.

⇒ $\max_{\alpha} \frac{1}{2} \alpha^T \alpha$ s.t. $\alpha_i \leq 1$ for all i

- SVMs aims to find the optimal separating hyperplane that maximizes the margin between two classes in a dataset.

- To achieve this, SVMs are formulated as a constrained optimization problem.

- constrained optimization is a set of method to identify efficiently and systematically the best solution to a problem.

- The problem is characterized by a number of potential solution in the absence presence of identified constrained.

D) Maximizing the margin:

- The margin is the distance between the decision boundary and the closest data points from the classes, called support vectors.
- The equation of the hyperplane in an SVM is given by:

$$w \cdot x + b = 0$$

w = weight vector

b = bias (offset)

x = input vector

The decision boundary separates the classes such that it is far away from both classes.

$$w \cdot x_i + b \geq 1, \text{ for } y_i = +1$$

2) Constrained Optimization problem:

Thus, SVM optimization problem can be written as:

Minimize $\frac{1}{2} \|w\|^2$

subject to $y_i(w \cdot x_i + b) \geq 1$

for all $i = 1, 2, \dots, n$

subject to the constraints is

$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$

and at least b is a margin of separation.

This is a convex optimization problem because both the objective function and linear constraints are convex, i.e. "linear".

Thus, the minimum value is being sought.

Q 3) Soft Margin SVM

→ allows margins with some error.

If the data is not perfectly separable, we introduce slack variables ξ_i to allow for some misclassifications.

The optimization problem becomes:

$$\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_i \xi_i$$

subject to $y_i(w \cdot x_i + b) \geq 1 - \xi_i \quad \forall i$

Q Subject to:

what soft margin does it mean?

$y_i(w \cdot x_i + b) \geq 1 - \xi_i$

$\xi_i \geq 0 \quad \forall i$ shows if $\xi_i > 0$ →

$$w \cdot x_i + b = (w; b) \cdot x_i$$

means ξ_i margin

whereas $w \cdot x_i + b$ is the linear boundary

margin given by two parallel lines

Q.5 What are kernel functions? List some kernel functions. What is their use in SVM?

⇒ $y = \sum (w_i x_i + b)$ if

- Kernel functions is a method used to take data as input and transform it into the required form of processing data.
- "Kernel" is used due to set of mathematical functions used in SVM providing the window to manipulate the data.
- So, kernel function generally transforms the training set of data so that a non-linear decision surface is able to transform to a linear equation in a higher number of dimension spaces.
- Basically, it returns the inner product between two points in a standard feature dimension.
- Some common kernel functions:

1) Linear Kernel

- ⇒ This kernel is used when the data is linearly separable in the original space.
- It's formula is:

$$K(x_i; x_j) = x_i \cdot x_j$$

2) Polynomial Kernel

- ⇒ The polynomial kernel allows SVM to create non-linear decision boundaries by using polynomial functions of the input data.

$K(x_i, x_j) = \gamma (x_i \cdot x_j + b)^d$ with analogy
between dot product and sigmoid function

3) Gaussian kernel

⇒ The gaussian kernel maps the data into an infinite-dimensional space, allowing for highly complex decision boundaries.

$$K(x_i, x_j) = \gamma \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$$

4) Sigmoid kernel

⇒ This kernel function resembles the activation function of a neural network.

$$K(x_i, x_j) = \tanh(\alpha(x_i \cdot x_j) + \beta)$$

When the relationship between features is similar to the neural network structure or in cases where sigmoid-like nonlinearities are expected.

11w11 11w11

11w11

11w11

11w11

Q. 6

Explain the concept of Kernel trick? Discuss with example, how kernel tricks can speed up the computations.

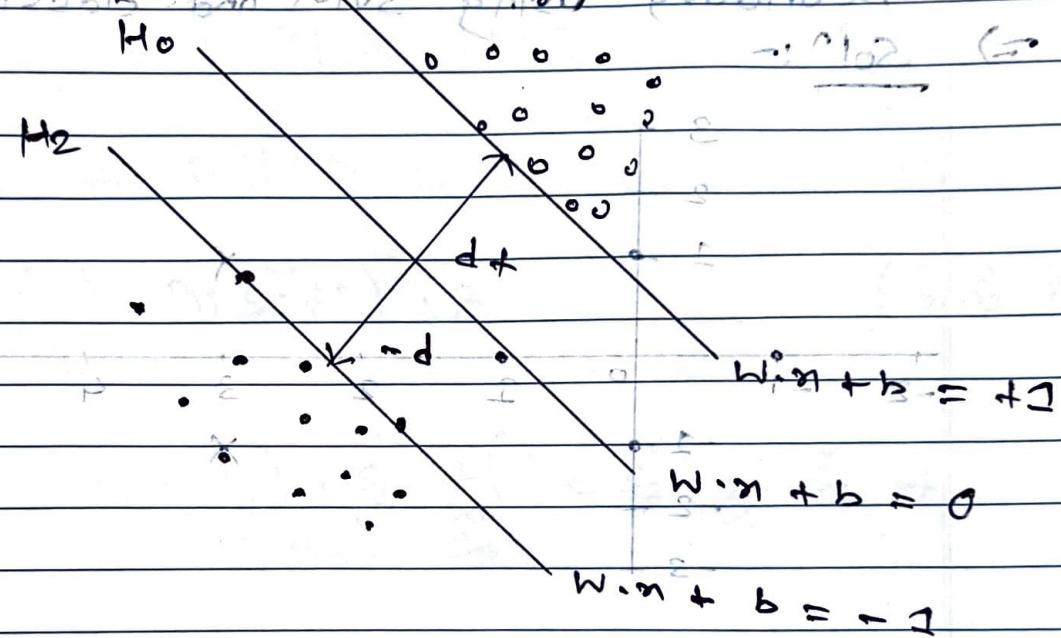
→ ~~to what are support vectors and margins?~~
 - A kernel trick is a simple method where a Non-linear data is projected onto a higher-dimensional space so as to make it easier to classify the data. Then it is linearly divided by a plane.

- This is mathematically achieved by Lagrangian formula using Lagrangian multipliers.
 Our basic idea of kernel trick is to find the planes which can separate, classify or split the data with maximum margin is also called street width.

The distance from the point (x_1, y) to a line $Ax + By + c = 0$ is $\frac{|Ax + By + c|}{\sqrt{A^2 + B^2}}$.
 In order to maximize the margin, the distance between H_1 and H_2 is then $\frac{2}{\|w\|}$, so the total distance between H_1 and H_2 is $\frac{2}{\|w\|}$.

In order to maximize the margin, we need to minimize $\|w\|$.
 Thus we are trying to optimize the margin or street width by maximizing the distance by maximizing distance

(18) Between two support vectors w is said to be between two points $(1,2)$, $(1,0)$, $(-1,2)$, $(0,1)$, $(-1,0)$, $(1,0)$, $(0,1)$ and $(-1,-1)$ so that H_1 is margin with best separating plane can be written as



→ How Kernel Trick speeds up computation in which kernel function $\phi(x)$ maps data into high-dimensional space.

1) Avoid Explicit Transformation

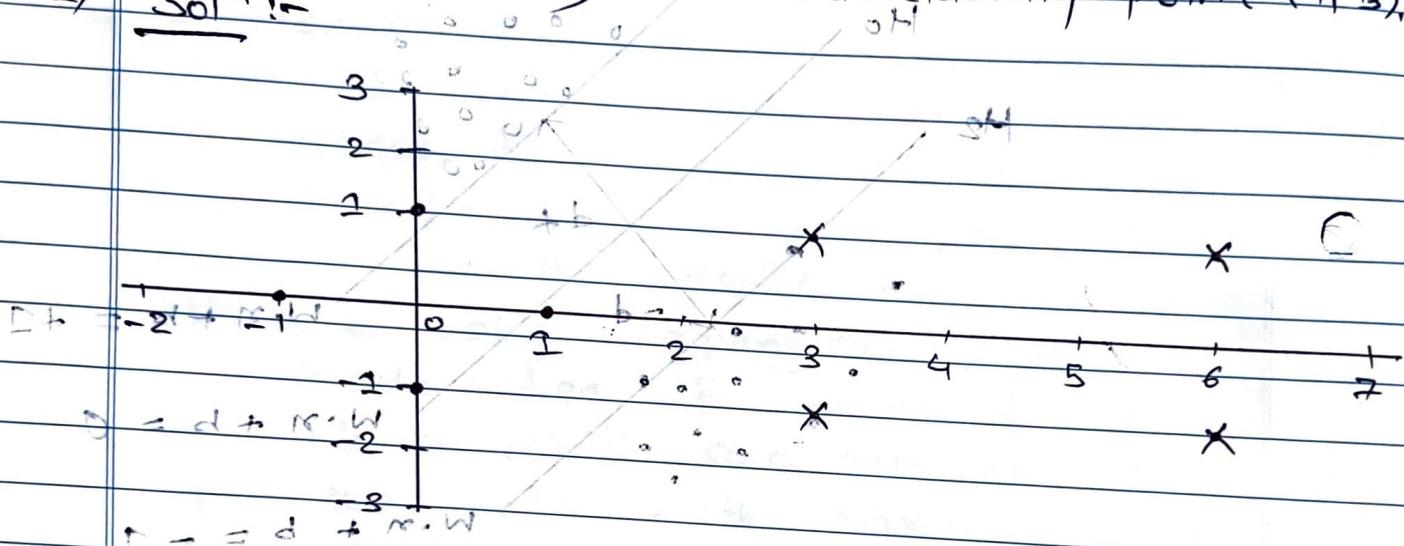
→ In cases where the transformation $\phi(x)$ maps the data into the high-dimensional or even infinite-dimensional space, explicitly calculating the transformed data is impractical.

2) Efficient Computation

→ The kernel trick allows SVM to compute the dot product between data points in the transformed space without the need to explicitly calculate the coordinates of the data in that space.

Q.7 Give 4+1ly labelled data points as $\{(3,1), (3,-1), (6,1), (6,-1)\}$ and 1+1ly labelled data points as $\{(1,0), (0,1), (0,-1), (-1,0)\}$. Find the parameters of the decision boundary using SVM and classify point $(1,3)$.

\Rightarrow Soln:-



The circles are negatively labelled data points and cross are positively labelled data points.

Find three support vectors.

$$\text{Let } \mathbf{S}_1 = \begin{pmatrix} 3 \\ 1 \\ 0 \end{pmatrix}, \mathbf{S}_2 = \begin{pmatrix} 1 \\ 3 \\ 1 \end{pmatrix}, \mathbf{S}_3 = \begin{pmatrix} 3 \\ -1 \\ 2 \end{pmatrix}$$

After adding bias input into S_1, S_2, S_3

$$\text{Let } \mathbf{S}_1 = \begin{pmatrix} 2 \\ 0 \\ 1 \end{pmatrix}, \mathbf{S}_2 = \begin{pmatrix} 3 \\ 1 \\ 0 \end{pmatrix}, \mathbf{S}_3 = \begin{pmatrix} 3 \\ -1 \\ 2 \end{pmatrix}$$

Graph of bias with function $y = \frac{3}{2}x - 1$

- Now, seek form of general solution -

- Eq 2.9

$$D(z) = w^T z + b$$

$$(z) \rightarrow r = (\sum x_i \tilde{s}_i t_i) \cdot z + b \quad (1)$$

$$(p) \rightarrow l = \alpha z_1 + \beta z_2 + \gamma z_3 \quad (\text{ignore } b)$$

$$\therefore (D(z)) = (\sum x_i \tilde{s}_i t_i) \cdot z \quad \leftarrow (1)$$

Eq in (n), (o), (e) go parallel work

- We need to find $\lambda_1, \lambda_2, \lambda_3$ using eq (1).

$$y_1 = \left(\sum_{i=1}^3 \lambda_i \tilde{s}_i t_i \right) \cdot z \quad \leftarrow \text{Eq } (1)$$

$$\therefore y_1 = (\lambda_1 \tilde{s}_1 t_1) \cdot \tilde{s}_1 + (\lambda_2 \tilde{s}_2 t_2) \cdot \tilde{s}_2 + (\lambda_3 \tilde{s}_3 t_3) \cdot \tilde{s}_3 = w$$

- Now, $y_1 = 1$

$$\therefore (1) \cdot z + (1) \cdot z = 1$$

$$\lambda_1 (\tilde{s}_1 \cdot \tilde{s}_1) t_1 + \lambda_2 (\tilde{s}_2 \cdot \tilde{s}_1) t_2 + \lambda_3 (\tilde{s}_3 \cdot \tilde{s}_1) t_3 = 1$$

$$\lambda_1 (\tilde{s}_1 \cdot \tilde{s}_2) t_1 + \lambda_2 (\tilde{s}_2 \cdot \tilde{s}_2) t_2 + \lambda_3 (\tilde{s}_3 \cdot \tilde{s}_2) t_3 = 1$$

$$\lambda_1 (\tilde{s}_1 \cdot \tilde{s}_3) t_1 + \lambda_2 (\tilde{s}_2 \cdot \tilde{s}_3) t_2 + \lambda_3 (\tilde{s}_3 \cdot \tilde{s}_3) t_3 = 1$$

- Putting the values in above equations,

$$\lambda_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} t_1 + \lambda_2 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} t_2 + \lambda_3 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} t_3$$

to simplify our work, taking $t_1 = 1, t_2 = 1, t_3 = 1$

$$= -2 \cdot 1 + 2 \cdot 1 + 2 \cdot 1 = 2$$

Let's do step by step of simplification

$$2\lambda_1 + 4\lambda_2 + 4\lambda_3 = -1 \quad (2)$$

- Similarly solving for rest two equations, we get -

$$\begin{aligned} d + \vec{s}^T w &= (s) \alpha \\ 4x_1 + 11x_2 + 11x_3 &= 1 \quad (3) \\ (d + \vec{s}^T w) + 9x_1 + 9x_2 + 11x_3 &= 1 \quad (4) \\ (4) \rightarrow \vec{s}^T (\vec{s} + \vec{x}) &= ((s)\alpha) \end{aligned}$$

Now, solving eq (2), (3), (4) we get,

(i) $\vec{s} + \vec{x}$ is not at bias, so \vec{x} must be basis vector

$$x_1 = -3.5$$

$$(2, 2, 2) \rightarrow x_2 = 0.75 \cdot (1 + 2 + 2) = 3.75$$

$$x_3 = 0.75$$

$$\begin{aligned} &+ 2 \cdot (x_1 + x_2 + x_3) + 2 \cdot (1 + 2 + 2) = 14 \\ \therefore \tilde{w} &= \sum_{i=1}^3 \vec{s}_i x_i \cdot (x_1 + x_2 + x_3) \end{aligned}$$

$$\begin{aligned} \therefore \tilde{w} &= -3.5(1) + 0.75(3) + \\ &+ 8 \cdot (1 + 2 + 2) \cdot x_1 + 8 \cdot (1 + 2 + 2) \cdot x_2 + 8 \cdot (1 + 2 + 2) \cdot x_3 \\ &+ 8 \cdot (1 + 2 + 2) \cdot 0.75(1 + 2 + 2) \cdot x_1 + 8 \cdot (1 + 2 + 2) \cdot 0.75(1 + 2 + 2) \cdot x_2 + 8 \cdot (1 + 2 + 2) \cdot 0.75(1 + 2 + 2) \cdot x_3 \\ &= 8 \cdot (1 + 2 + 2) \cdot 0.75(-3) + 8 \cdot (1 + 2 + 2) \cdot 0.75(1) \end{aligned}$$

so \tilde{w} is sum of positive and negative -

$$\begin{aligned} &= \cancel{8 \cdot (1 + 2 + 2) \cdot (-3)} + \cancel{8 \cdot (1 + 2 + 2) \cdot 0.75} \tilde{w} + 8 \cdot (1 + 2 + 2) \cdot 0.75 \\ &= b + \tilde{w} \end{aligned}$$

- Since the support vectors were aligned with bias, entries, entries in \tilde{w} correspond to hyperplane with offset b ,

$$\therefore \rightarrow b = \tilde{w}^T \vec{s} + \tilde{w}^T \vec{x}$$

where in, $y = w^T x + b$,

$$w = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \text{ and } b = -2.$$

- The point $(1, 3)$ is less than 2, since the decision boundary is $w_1 x_1 + b = 2$, this point lies on the negative side of the decision boundary.

