

SUPPORT VECTOR MACHINE (SVM) :

It is a machine learning algorithm

which can be used for both classification
outlier detection
and regression.

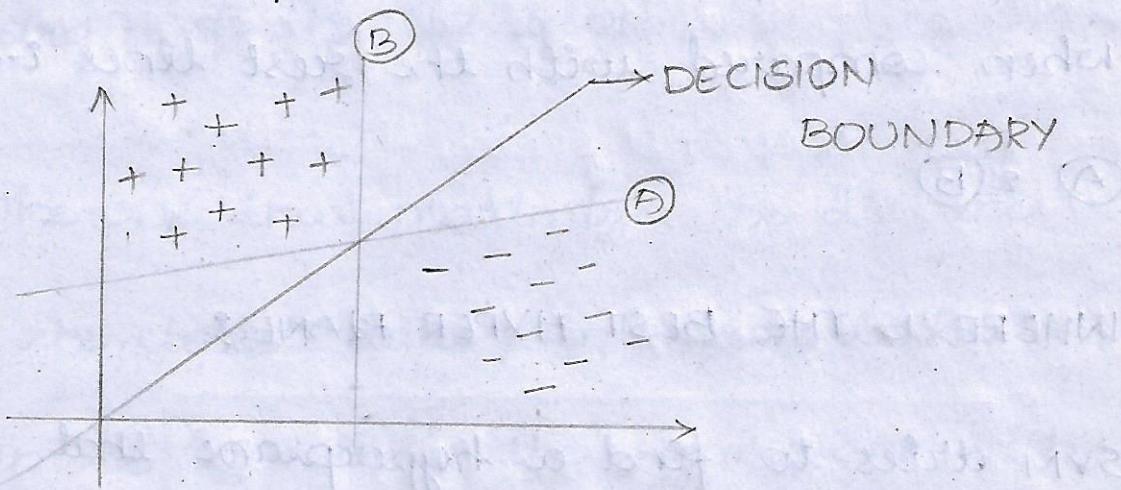
- Mostly used in classification problems.
- Support vectors are simply the co-ordinates of an individual observation.
- The SVM are particular linear classifiers which are based on the margin maximization principle.
- They perform structural risk minimization, which improves the complexity of the classifier with the aim of achieving excellent generalization performance.

(137)

→ A linear SVM classifier works by drawing a straight line between two classes.

PURPOSE OF SVM :

→ It uses a technique called the Kernel trick, to transform the data and then based on those transformations it finds an optimal boundary between the possible outputs.



DECISION BOUNDARY : Separation of classes.

→ A decision boundary is the region of a problem space in which the output label of a classifier is ambiguous.

→ If the decision surface is a hyperplane, then the classification problem is linear and the classes are linearly separable.

* → Decision boundaries are not always clear cut.

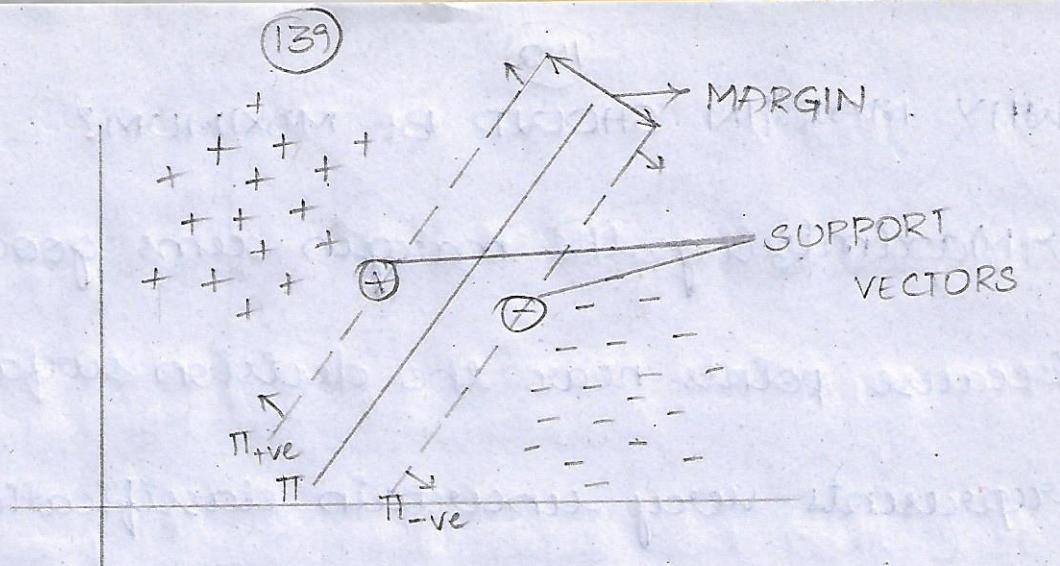
* → Decision boundary has the maximum distance for both +ve and -ve classes when compared with the rest lines i.e.,

(A) & (B)

WHERE IS THE BEST HYPER PLANE?

SVM tries to find a hyperplane that separates +ves and -ves as widely as possible.

→ Also called as "MARGIN MAXIMIZING HYPERPLANE".



MARGIN :

The distance from the decision surface to the closest data point determines the margin of the classifier.

→ The one that maximizes the distance to the closest data points from both classes then we say it is the hyperplane with maximum margin.

WHY MARGIN SHOULD BE MAXIMUM?

→ Maximizing the margin seems good because points near the decision surface represents very uncertain classification decisions i.e.,

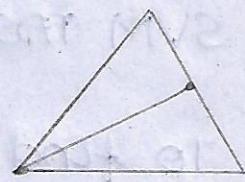
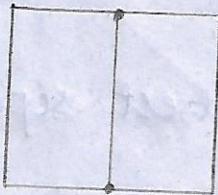
there is a 50% chance of the classifier deciding either way.

* → A classifier with a large margin makes no low certainty classification decisions.

MATHEMATICALLY,

$$\text{Margin} = \text{distance}(\pi_{+ve}, \pi_{-ve})$$

So, as the margin increases, generalization accuracy increases.



HOW TO ACTUALLY FIND THE MARGIN MAXIMIZING HYPERPLANE?

→ In order to understand it, we need have the idea behind the CONVEX HULL.

HULL: closed figure

From the above figures,

Draw any two points anywhere and join them by a line. This is called as a Convex Hull.

→ If the line goes out of the boundary then it is not a convex hull.

SVM TASK :

To find a line that best separates the +ves and -ves as widely as possible.

line that best separates - $y_{act} * y_{pred}$.

As widely as possible - maximizing the margin.

STEPS FOR SVM ALGORITHM:

STEP-1: Draw the convex hull for +ves and -ves

STEP-2: Find the shortest line connecting those hulls.

STEP-3: Bisect the line i.e., break them into two equal parts.

So, If $y_{act} * y_{pred} \geq 0 \rightarrow$ SIGNED DISTANCE IS POSITIVE & CORRECT

CLASSIFICATION!

(143)

\Rightarrow max. margin subjected to (s.t.)

$$y_i * \{m x_i + c\} \geq 0 \quad \forall i \text{ (for all)}$$

$$\Rightarrow \boxed{y_i * \{m x_i + c\} \geq 0 \quad \forall i}$$

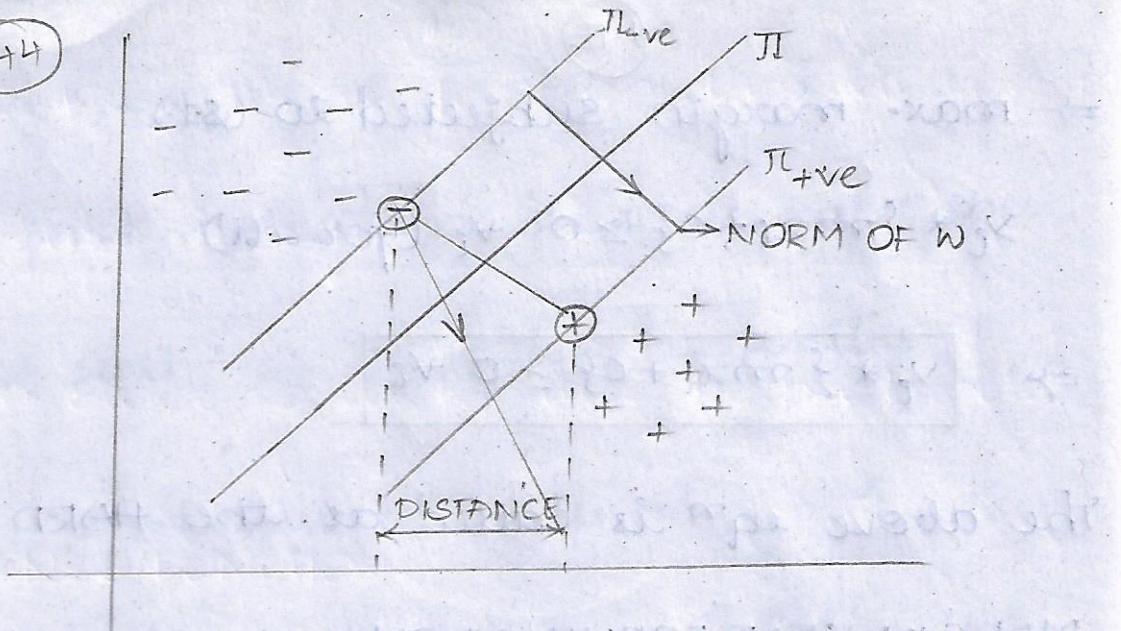
The above eqⁿ is called as the HARD MARGIN SVM FORMULATION.

The above eqⁿ states that for all the data points, it should be true such that every point in the data is correctly classified.

$$\Rightarrow y_i * \{w^T x_i + w_0\} \geq 0 \quad \forall i \rightarrow \text{other form of representation.}$$

Now, let's see the mathematical representation of margin (least distance).

144



$$\Rightarrow \Pi_{-ve} = w^T x_c + w_0 = -1 \rightarrow \textcircled{1}$$

$$\Rightarrow \Pi = w^T x_c + w_0 = 0$$

$$\Rightarrow \Pi_{+ve} = w^T x_c + w_0 = 1 \rightarrow \textcircled{2}$$

Subtracting
Adding $\textcircled{1} + \textcircled{2}$

$$\Rightarrow w^T x_c + w_0 = +1 \rightarrow \textcircled{2}$$

$$\Rightarrow \underline{\underline{w^T x_c + w_0 = -1}} \rightarrow \textcircled{1}$$

$$\Rightarrow w^T x_{c+ve} - w^T x_{c-ve} = 2$$

$$\Rightarrow w^T (x_{c+ve} - x_{c-ve}) = 2$$

$$\Rightarrow w^T (x_{+ve} - x_{-ve}) = 2$$

$$\Rightarrow \boxed{\frac{w^T}{\|w\|} (x_{+ve} - x_{-ve})} = \boxed{\frac{2}{\|w\|}} \rightarrow \text{MARGIN}$$

↳ PROJECTION DISTANCE ON \vec{w}

(14b)

∴ The mathematical representation of margin is

$$\Rightarrow \max \left\{ \frac{2}{\|w\|} \right\} \text{ s.t. } y_i * \{ w^T x_i + w_0 \} \geq 1 \forall i$$

This is called as Hard margin support vector formulation.

→ There is a dual form of hard margin SVM formulation i.e.,

$$\Rightarrow \max \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j \text{ s.t.}$$

$$0 \leq \alpha_i \leq C, \quad \sum_{i=1}^n \alpha_i y_i = 0$$

To solve the dual form of hard margin we use "LAGRANGE'S MULTIPLIER" which is a constrained optimisation.

Now, we have a problem i.e., Dual form of hard margin.

→ In order to it, we use Kernel Function.

KERNEL FUNCTION:

SVM algorithms use a set of mathematical functions that are defined as the Kernel.

→ The function of Kernel is to take data as input and transform it into required form.

→ For example,

linear, non-linear, polynomial, radial basis function (RBF) and sigmoid.

WHAT IS THE PURPOSE OF KERNEL TRICK?

→ It is used to offer a more efficient and less expensive way to transform

data into higher dimensions.

→ With that saying, the application of the Kernel trick is not limited to the SVM algorithm.

→ Any computations involving the dot products (x, y) can utilize the Kernel trick.

WHAT KERNEL IS USED IN SVM?

→ Use linear SVM (LOGISTIC REGRESSION) for
↳ linear separation in the data
linear problems

→ Use Non-linear kernels such as radial basis function kernel for non-linear problems.

so, now we use the kernel in the Dual form of soft margin

\Rightarrow KERNEL TRICK

$$= \max \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

For Quadratic, $K(x_i, x_j) = K(x_i \cdot x_j)^2$

For Polynomial, $K(x_i, x_j) = K(x_i \cdot x_j)^d$

where $d = 1, 2, 3, \dots, d_{n-1}, d_n$

$K(x_i \cdot x_j)^2$ transforms to 3D form.

RADIAL BASIS FUNCTION KERNED SVM:

A function whose values depends on the distance from the origin (or) from some point.

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

KERNELTRICK \rightarrow Take the data & transforms

to d' where $d' > d$.

\hookrightarrow In those d' - we figure out linear svm

(149)

→ so, now we have hard margin.

→ Apart from it we have soft margin too.

→ lets see what does soft margin does.

SOFT MARGIN SVM :

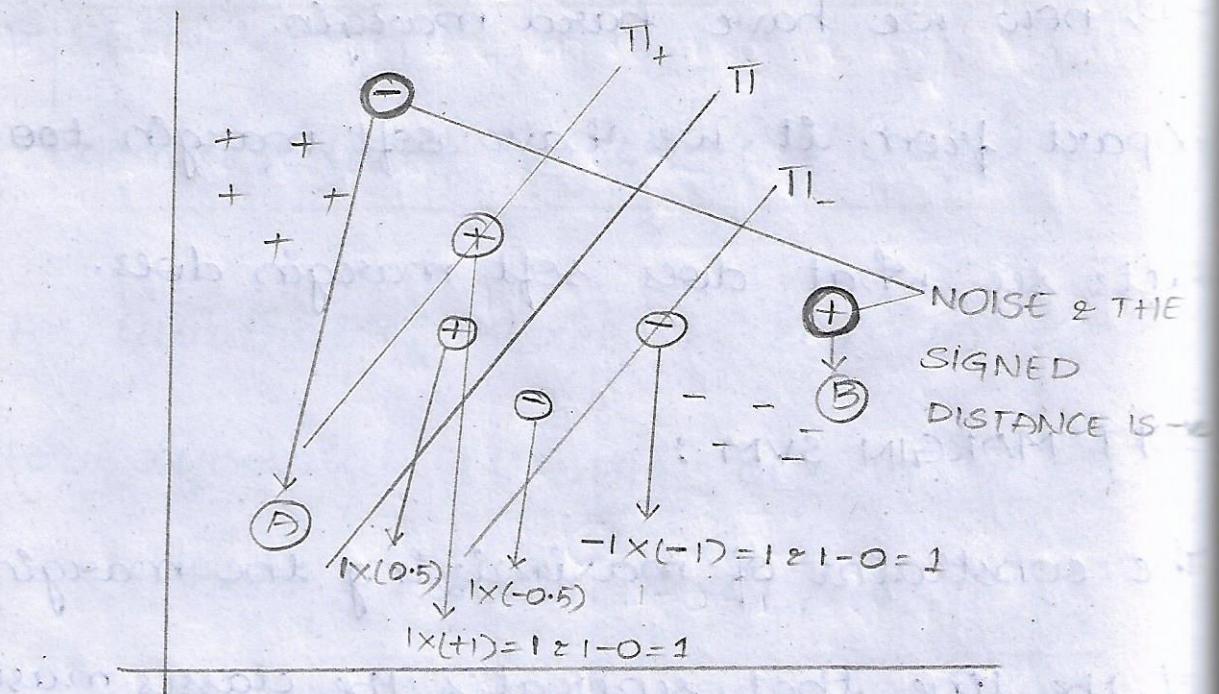
The constraint of maximizing the margin of the line that separates the classes must be relaxed. This is often called as soft margin classifier.

→ This change allows some points in the training data to violate the separating line.

Now lets see the formulation and

the rest things (Noise).

(150)



$$① \Rightarrow y_c * \{ w^T x_i + w_0 \} = 1 \times 2.5 = -2.5$$

$$② \Rightarrow y_c * \{ w^T x_i + w_0 \} = 1 \times -2.5 = -2.5$$

NOISE (OR) NOISY DATA :

Noisy Data is a data that has relatively signal-to-noise ratio.

→ The error is referred to as Noise.

→ Noise creates trouble for machine

learning algorithms because if not trained

properly, algorithms can think of noise

(151)

to be a pattern and can start generalizing from it, which of course is undesirable.

→ so, from $A \geq B$, we need to convert into $1 - \text{something}$.

where something is "ZETA" and represented as ' ζ '

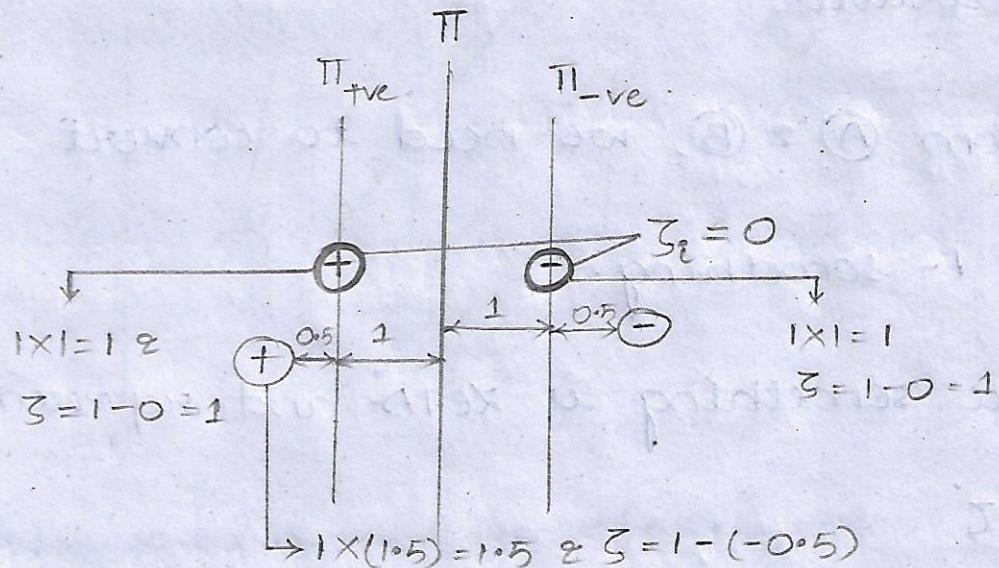
$$\Rightarrow \text{From } A \geq B \rightarrow 1 - \zeta = 1 - (3.5) = -2.5$$

For both +ve & -ve, (Noisy points) the Zeta value is 2.5

$\zeta \rightarrow$ Adding such a term is called as regularization.

→ Due to which we will make smallest possible error with hyperplane that maximizes the margin.

The eqn's can be minimized by adding negative values of ξ_i to it.



From the above figure,

$\xi_i < 0$ for all the correctly classified classes.

According to $\Pi_{+ve} \neq \Pi_{-ve}$ the classes are misclassified.

$\rightarrow \xi_i \geq 0$ - The point is misclassified according to the $\Pi_{+ve} \neq \Pi_{-ve}$

where ξ_i is called as slack variable.

(153)

So, if ξ_i increases, the point is far away
in the incorrect direction i.e., $\Pi_{\text{true}} \geq \Pi_{\text{true}}$

According to $\Pi_{\text{true}} \geq \Pi_{\text{true}}$:

$$\text{CORRECTLY CLASSIFIED} - y_i * \{ w^T x_i + w_0 \} \geq 1$$

$$\xi_i < 0$$

$$\text{MISCLASSIFIED} - y_i * \{ w^T x_i + w_0 \} \leq 1 \rightarrow \xi_i \geq 0$$

Now we are trying to max. the margin.

$$\Rightarrow \text{max margin} = \min \frac{1}{\text{margin}}$$

[From Optimization theory, w.k.t]

$$\max f = \min (-f) \geq \max (f) = \min (f')$$

$$\Rightarrow \text{max. margin} = \min \left\{ \frac{1}{2} \|w\|^2 + C \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \right\}$$

$$\text{s.t. } y_i * (w^T x_i + w_0) \geq 1 - \xi_i \quad \forall i$$

where C - HYPERPARAMETER OF SVM

If we have noises,

(154)

$$\Rightarrow \text{max. margin} = \min \left\{ \frac{1}{2} * \|w\| + C * \frac{1}{n} \sum_{i=1}^n \xi_i \right\}$$

s.t. $y_i * \{w^T x_i + w_0\} \geq 1 - \xi_i \geq 0 \quad \forall i$

The above eqⁿ is called as the soft margin of support vector machine formulation.

TIME COMPLEXITY OF ALGORITHM:

KNN > SVM > LOGISTIC REGRESSION



training takes a lot of time.

→ If we have points on line then how to do SVM on it?

