

"Support Vector Machines"

Support Vector machine, abbreviated as SVM, can be used for both Regression & Classification tasks. But, it is widely used in Classification.

→ The Goal of SVM:-

The goal of the SVM algorithm is to Create the best line or decision boundary 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 a "hyperplane".

SVM chooses the extreme point/vectors that help in creating the hyperplane. These extreme cases are called as support vectors, & hence this algorithm is termed as "Support Vector Machine".

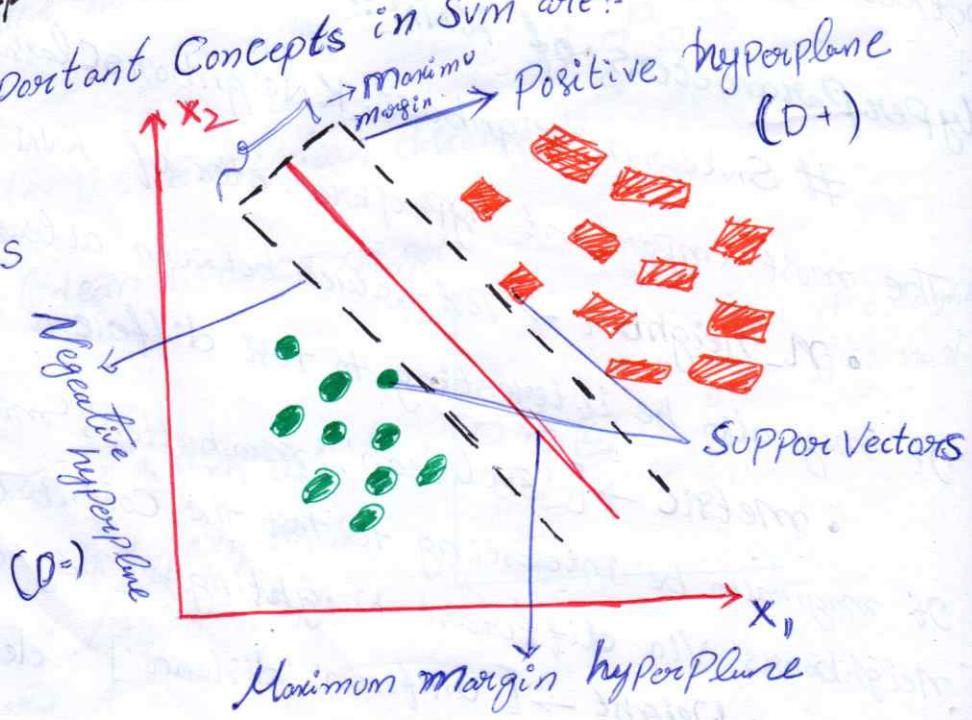
→ Important Concepts in SVM:-

The following are important concepts in SVM:-

↳ Support Vectors

↳ Hyper Plane.

↳ Margin.



→ HyperPlane :-

There can be many possible hyperplanes that could to be chosen, to separate the two classes of data points. Our objective is to find a plane that has "maximum margin", i.e., the maximum distance between data points of both classes. Maximizing the margin distance provides some reinforcement. So that future data points can be classified with more confidence.

SVM Selects the hyper-plane which classifies the classes accurately prior to maximizing margin.

Clases Segregates iteratively that SVM will generate hyperplanes correctly.

→ First, SVM will separate the classes in best way.

→ Then it will choose the hyperplane that separates the classes correctly.

The dimensions of the hyperplane depend on the features present in the dataset, which means if there are 2 features, then hyperplane will be a straight line, and if there are 3 features, then hyperplane will be 2-Dimensional plane. It becomes different to imagine when the number of features exceeds 3.

3.

→ This hyperplane is also called Decision boundary, mathematical equation of that decision boundary as a line

$$g_n = g(x) = w^T x + b = 0$$

$$y_n \in \{1, -1\}$$

↪ Support Vectors :-

The data points or vectors that are the closest to the hyperplane & which affect the position of the hyperplane are termed as Support Vector. Since these vectors support the hyperplane, hence called a Support Vector.

Using these Support Vectors, we maximize the margin of the classifier. If we do some geometry, we can figure out that the distance from any point to the decision boundary is the following.

$$r = \frac{|g(x)|}{\|w\|} = \frac{1}{\|w\|}$$

Our goal is to maximize r for the points closest to the optimal decision boundary.

↪ Margin :-

It may be defined as the gap between two lines on the closest data points of different classes.

It can be calculated as the perpendicular distance from the line to the support vectors, given by:-

$$\text{margin} = \frac{2}{\|w\|}$$

Large margin is considered as a "good margin" & a small margin is considered as a "bad margin".

There are two types of margins that are used in SVM.

Hard & Soft :

When the training dataset is linearly separable, SVM can

simply select two parallel lines that maximize the marginal distance; this is called a "Hard margin".

When the training dataset is not fully linearly separable,

then the SVM allows some margin violation. It allows some data points to stay on the wrong side of the hyperplane or between the margin & hyperplane so that the accuracy is not compromised, this is called a "Soft margin".

Soft margin permits few of the above data points to get misclassified. Also, it tries to make the balance back & forth between finding a hyperplane that attempts to make less misclassifications & maximize the margin.

Types of SVM:

→ Linear SVM: Linear SVM is used for linearly separable data, which means if a dataset can be classified into two classes by using a single straight line, then such data is termed as linearly separable data, & classifier is used as Linear SVM.

→ Non-linear SVM: Non-linear SVM is used for non-linearly separated data, which means if a dataset cannot be classified by using a straight line, then such data is termed as non-linear data & classifier used is called Non-linear SVM classifier.

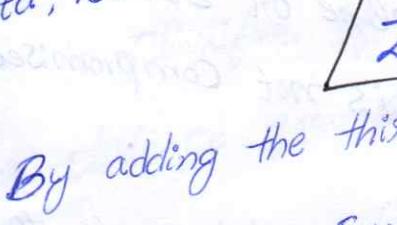
Non-Linear SVM:-

→ If data is linearly arranged, then we can separate it by using a straight line, but for non-linear data, we cannot draw a straight line.

→ To separate these data points, we need to add one more dimension.

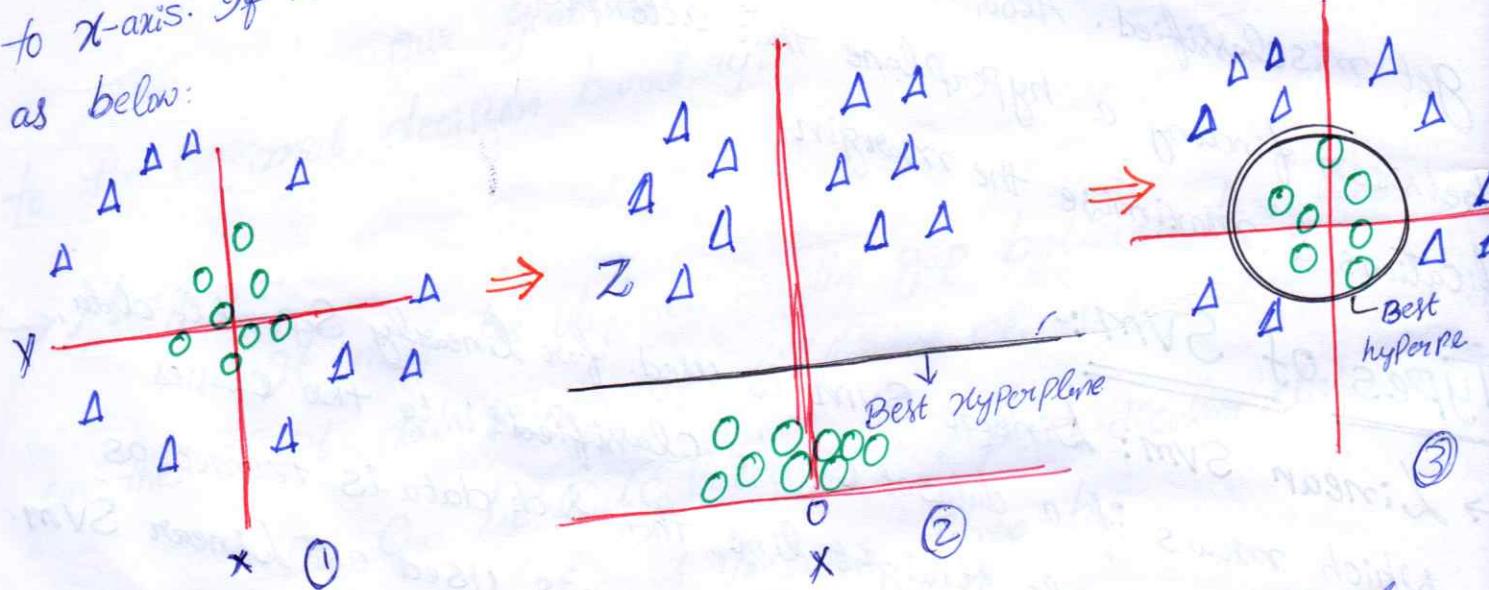
For linear data, we have used two dimensions x & y . So for non-linear data, we will add a third dimension. It can be calculated as:

$$Z = x^2 + y^2$$

→ By adding the third dimension, the sample space will become as

→ So, now SVM will divide the datasets into classes.

→ Since, we are in 3-d Space, hence it is looking like a plane parallel to x -axis. If we convert it in 2-d Space with $Z=1$, then it will become as below:



→ Hence we get a circumference of radius 1 in case of non-linear data.

* All values for Z , would be positive, always because Z is the squared sum of both x & y .

Ditheesh

Svm KERNELS - KERNEL TRICK :-

In Practice, Svm algorithm is implemented with Kernel that transforms an input data Space into the Required form.

Svm uses a technique called the "Kernel trick" in which takes a low dimensional input Space & transforms it into a higher dimension Space.

In Simple Words, Kernel Converts non-Separable Problems into Separable Problems by adding more dimensions to it. It is mostly useful in non-linear Separable Problem.

Simply Put, it does some extremely complex data transformations, then find out the process to separate the data based on the labels or inputs you have defined. It makes SVM more powerful, flexible & accurate. The following are some of the types of Kernels used by SVM:

$$\text{Kernel} = \{ \text{"linear"}, \text{"Poly"}, \text{"RBF"}, \text{"Sigmoid"} \}$$

Linear Kernel:-

It can be used as a dot product between any two observations.

The formula of linear Kernel is as below:

$$K(x, x_i) = \text{Sum}(x * x_i)$$

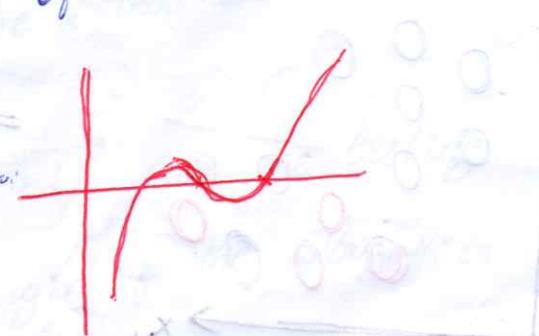
Polynomial Kernel:-

It is more generalized form of linear Kernel &

distinguish curved or non-linear input space. Following is the formula for Polynomial.

$$K(x, x_i) = 1 + \text{Sum}(x * x_i)^d$$

d - degree of Polynomial \rightarrow Specify manually.



Radial Basis Function (RBF) Kernel :-

RBF Kernel, mostly used in SVM Classification, maps input Space in indefinite dimensional Space. Following formula explains it mathematically-

$$K(x, x_i) = \text{Exp}(-\gamma \cdot \sum (x - x_i)^2)$$



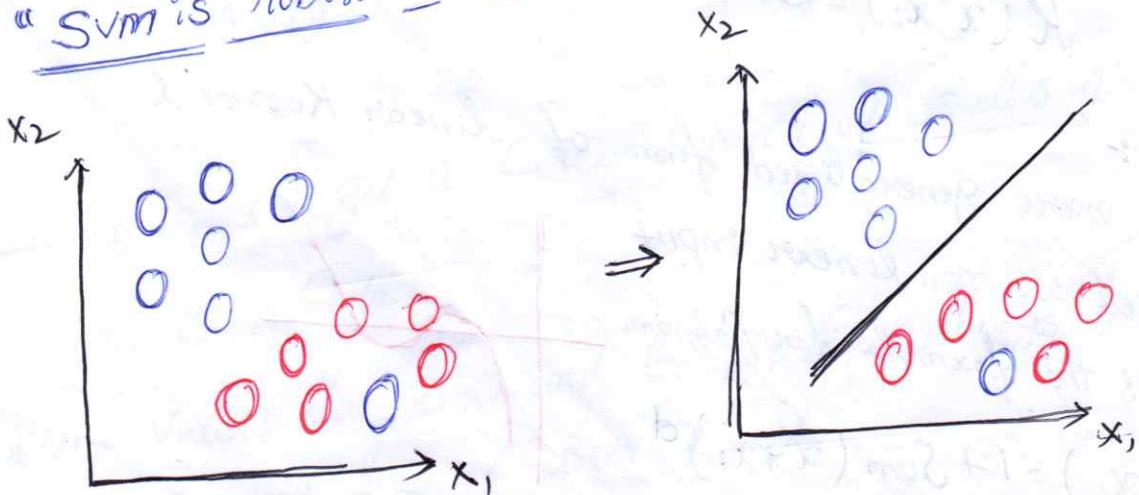
- Here, 'gamma' ranges from 0 to 1. We need to manually specify it in the learning Algorithms.
→ A good default value of gamma is 0.1.

SVM Robust To Outliers :-

→ SVM is Robust To Outliers :-

- Let Consider a Scenario like shown in below image:
→ Here we have one ball in the boundary of the red ball.
→ So, how does SVM Classify the data? It's simple: the blue ball in the boundary of red ones is an outlier of blue balls.
→ The SVM algorithm has the characteristic to ignore the outlier & finds the best hyperplane that maximizes the margin.

→ SVM is Robust To Outliers.



- So, in this type of Data Points, What SVM does, is it finds maximum margin as done with Previous data sets along with that it adds a Penalty Each time a Point crosses the margin. Margins in these type of Cases are called Soft margin.
- ~~With the margins~~ When there is a Soft margin to the dataset, the SVM tries to minimize: $(\frac{1}{margin} + \lambda(\epsilon_{Penalty}))$.
- Hinge loss is a commonly used Penalty. If no violations no hinge loss. If violations hinge loss proportional to the distance of Violation.

How SVM Works :-

In logistic regression, we take the output of the linear function & squash the value within the range [0, 1] using the sigmoid function. If the squashed value is greater than a threshold value (0.5) we assign it a label 1, else we assign it a label 0.

In SVM, we take the output of the linear function & if that output is greater than 1, we identify it with one class & if the output is -1, we identify it with another class. Since if $g(x) = 0$, it means that you testing sample belongs to the separating hyperplane & thus there cannot be some decision about its truth label. Note that, if $-1 < g(x) < 1$, then the testing datum lies inside the so-called margin of the classifier.

* Since the threshold values are changed to 1 & -1 in SVM, we obtain this reinforcement range of value [-1, 1] which acts as margin.

Maths Behind SVM :-

Basic Linear Algebra:-

Support Vector machine basically helps in sorting the data into two or more categories with the help of a boundary to differentiate similar categories.

So, first let's revise the formulae for how each data is represented in space, & what is an equation of a line which will help in segregating the similar categories, & lastly the distance formula between a data point & the line (hyperplane).

Vectors: (A point in the space) :

Vectors are mathematical quantity which has both magnitude & direction. A point in the 2D plane can be represented as a vector between origin & the point.

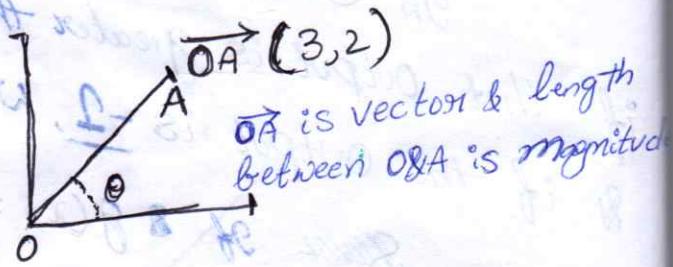
Length of vectors :-

Length of vectors are also,

Called as norms. It tells how far vectors are from the origin.

Length of vector $\vec{x}(x_1, x_2, x_3)$ is calculated as:

$$\|\vec{x}\| = \sqrt{x_1^2 + x_2^2 + x_3^2}$$



Director of Vectors :-

Direction of Vector $x(x_1, x_2, x_3)$ is calculated as:

$$\left\{ \frac{x_1}{\|x\|}, \frac{x_2}{\|x\|}, \frac{x_3}{\|x\|} \right\}$$

+ It is worth noting that the norm of a direction vector is always equal to 1. Because of this, the direction vector w is called the Unit Vector.

Dot Product:

The dot product of two vectors returns a scalar. It gives us some insights into how the two vectors are related.

The geometric formula of a dot product is defined as:

$$x \cdot y = \|x\| \cdot \|y\| \cos(\theta) = \frac{x_1 y_1 + x_2 y_2}{\|x\| \|y\|} \times \|x\| \|y\|$$

$\hookrightarrow \cos(\theta)$ is simplified.

$$= x_1 y_1 + x_2 y_2$$

This is the algebraic formula of dot product. In general, dot product can be computed as the following for

$$x \cdot y = \sum_{i=1}^n x_i y_i$$

Hyperplane :-

The two dimensional linearly separable data can be separated by a line. The function of line is

$$y = ax + b$$

We rename x with x_1 & y with x_2 & we get

$$ax_1 - x_2 + b = 0$$

If we define $x = (x_1, x_2)$ & $w = (a, -1)$, we get

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

$$\Rightarrow h(x_i) = w^T x_i + b = 0$$

This is Equations of hyperplane.

w^T → Unit Vector.

x_i → i^{th} training Example.

y_i → i^{th} Label Corresponding to x_i

Classifier :-

Once, we have the hyperplane, we can then use the hyperplane to make prediction, we define the hypothesis

function h as :-

$$h(x_i) = \begin{cases} +1 & w \cdot x_i + b \geq 1 \\ -1 & w \cdot x_i + b \leq -1 \end{cases}$$

Combining above two equation, it can be written as.

$$\hat{y}_i(w \cdot x_i + b) - 1 \geq 0 \quad y_i \in \{-1, 1\}$$

The point above hyperplane will be classified as

Class +1 (Positive Class) & the point below the hyperplane

will be classified as Class -1 (negative class)

$$\pi_{-ve} = w^T x_i + b = -1 \rightarrow ①$$

$$\pi = w^T x_i + b = 0 \rightarrow ②$$

$$\pi_{+ve} = w^T x_i + b = +1 \rightarrow ③$$

But here the signs were about training data, which is how we are training our model. That for positive class, give a positive sign & for negative give a negative sign.

But while testing this model on test data, if we predict a Positive class (Positive sign) correctly as Positive, then two positive makes positive & hence greater than zero result. The same applies if we correctly predict the negative group since two negatives will again make a positive.

But, if the model miss classifies the positive group as a negative group then one plus & one minus makes a minus, hence overall less than zero.

Summing up the above concept:

"The product of a predicted & actual label would be greater than zero on correct prediction, otherwise less than zero"

* $\mathbf{w}^T \rightarrow$ Transpose of (\mathbf{w}).

$$\hat{y}_i [\mathbf{w}^T \cdot \mathbf{x}_i + b] = \begin{cases} \geq 0 & \text{if correct} \\ < 0 & \text{if incorrect} \end{cases}$$

Margin :-

Margin is nothing but distance between two classes,

therefore subtracting Eq ② & ①, that gives:

$$\mathbf{w}^T \mathbf{x}_{i+ve} - \mathbf{w}^T \mathbf{x}_{i-ve} = 2$$

$$\Rightarrow \mathbf{w}^T (\mathbf{x}_{i+ve} - \mathbf{x}_{i-ve}) = 2.$$

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

$$\Rightarrow \text{margin} = \boxed{\frac{2}{\|w\|} = M}$$

Here $\|w\|$ is the Euclidean norm for the length of w given
(magnitude)

by:

$$\|w\| = \sqrt{w_1^2 + w_2^2 + w_3^2 + \dots + w_n^2}$$

Taking distance to any one of the planes/side from the hyperplane is given by:-

$$\boxed{d = \frac{1}{\|w\|}}$$

Optimization: (Primal & Dual) :-

The goal of SVM is to get a margin with maximum distance. Therefore, in order to find the optimal margin, what we want to do is to maximize d . With these considerations, we can finally state the optimization

Consider Considerations, we can finally state the optimization problem for the model.

$$d = \max \left(\frac{1}{\|w\|} \right) \rightarrow \min \|w\|$$

Well, from Calculus we know that in order to find a maximum value, we need to use derivatives, to make things for us. When computing derivatives, we notice it tells:

and as ℓ_2 optimization are often more stable than ℓ_1 optimization so again above equation can be written as :-

$$\min \frac{\|\mathbf{w}\|^2}{2} \text{ such that } y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 \geq 0$$

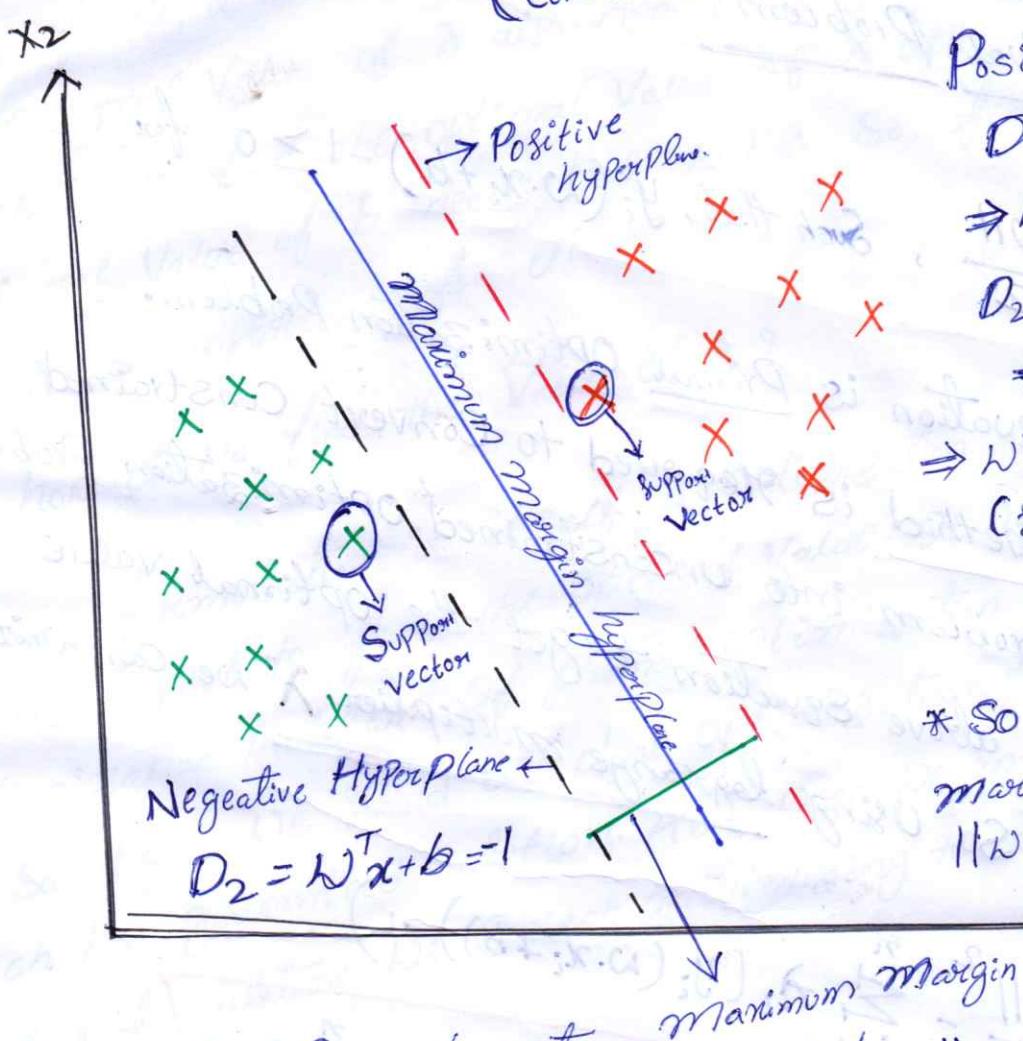
maximize margin

separating hyperplane

To Sum-UP :-

This is called as hard margin support vector formulation.

(Convex Quadratic Optimization)



Positive Hyperplane

$$D_1 = \mathbf{w}^T \mathbf{x} + b = 1$$

$$\Rightarrow \mathbf{w}^T \mathbf{x} + b - 1 = 0$$

$$D_2 = \mathbf{w}^T \mathbf{x} + b = -1$$

$$\Rightarrow \mathbf{w}^T \mathbf{x} + b + 1 = 0$$

$$\Rightarrow \mathbf{w}^T \mathbf{x} + b - 1 - (\mathbf{w}^T \mathbf{x} + b + 1) \\ \text{(Solve algebraically)}$$

$$\Downarrow \\ \frac{1}{2} \|\mathbf{w}\|^2$$

* So to increase the margin minimize $\|\mathbf{w}\| \Rightarrow \frac{1}{2} \|\mathbf{w}\|^2$

In layman term Sum optimization :-

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

if $y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 = 0$
then \mathbf{x}_i, y_i are support vectors

Elif, $y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 > 0$:
then save the parameters \mathbf{w}, b

Elif, $y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 < 0$:
then update the parameters \mathbf{w}, b .

- ⇒ We have to find the values of Parameters (w, b) , which minimizes $\|w\|$ for optimal hyperplane.
- ⇒ The problem of finding the values of w & b is called an optimization problem.
- ⇒ Any optimization problem can be formulated in two ways
- Primal & Dual Problem
- ⇒ $l = \min \frac{\|w\|^2}{2}$, such that, $y_i(w \cdot x_i + b) - 1 \geq 0$ for $i=1 \dots l$
- ⇒ The above equation is Primal Optimization Problem.
- ⇒ The above equation is required to convert constrained optimization problem into unconstrained optimization problem.
- ⇒ The goal of above equation to get the optimal value for w & b , so using Lagrange multiplier λ we can write.
- $$l = \frac{\|w\|^2}{2} - \sum_{i=1}^n \lambda_i (y_i(w \cdot x_i + b) - 1)$$
- $$= \frac{\|w\|^2}{2} - \sum_{i=1}^n \lambda_i (y_i(w \cdot x_i + b))$$
- $$\boxed{= \frac{\|w\|^2}{2} - \sum_{i=1}^n \lambda_i y_i w \cdot x_i - \sum_{i=1}^n \lambda_i y_i}$$
- $$\Rightarrow \frac{\partial l}{\partial w} = \nabla_w l(w, b, \lambda) = w - \sum_{i=1}^n \lambda_i y_i x_i = 0 \rightarrow ①$$
- $$\text{from } ① \Rightarrow w = \sum_{i=1}^n \lambda_i y_i x_i \rightarrow *$$

$$\Rightarrow \frac{\partial l}{\partial \lambda} = \nabla_{\lambda} l(w, b, \lambda) = \sum_{i=1}^n y_i (w \cdot x_i + b) - 1 = 0 \rightarrow \textcircled{2}$$

$$\Rightarrow \frac{\partial l}{\partial b} = \nabla_b l(w, b, \lambda) = - \sum_{i=1}^n x_i \cdot y_i = 0 \rightarrow \textcircled{3}$$

~~From equations ② & ③ we can find the values of w & b.~~

As, from the above formulation, we only able to find the optimal value of w , & that is to dependent on λ . So we need to find the optimal value of λ also. And finding optimal value of b needs to find the optimal value of λ also, & finding optimal value of w needs both w & λ . So, finding the value of λ is important for us.

So how do we find the value of λ ?

Above, formulation is itself a optimization algorithm, but it not helpful to find the optimal value. It is Primal Optimization Problem. If Primal Optimization doesn't result in Solution, We should use dual optimization formulation, which has guaranteed Solution. Also when we move from Primal to dual formulation we switch minimizing to maximizing the loss function.

We will divide the complete formulation into three parts to easier to understand:

- Problem formulation & Substitution value from Primal

- Simplify the loss function equation after Substitution

- final optimization formulation to get the value of λ .

$$f_{\text{dual}} = \frac{\|w\|^2}{2} - \sum_{i=1}^n \lambda_i (y_i (w \cdot x_i + b)) + \sum_{i=1}^n \lambda_i$$

Substituting the value of $w = \sum_{i=1}^n \lambda_i y_i x_i$ into above equation,

$$f_{\text{dual}} = \frac{\sum_{i=1}^n \|\lambda_i y_i x_i\|^2}{2} - \sum_{j=1}^n \sum_{i=1}^n \lambda_j y_j (\|\lambda_i y_i x_i\|) \cdot x_j + b + \sum_{i=1}^n \lambda_i$$

* Above equation is called Dual optimization Problem.

$$\Rightarrow f_d = \frac{\sum_{i=1}^n (\lambda_i y_i x_i)^T \cdot (\lambda_i y_i x_i)}{2} - \sum_{j=1}^n \sum_{i=1}^n \lambda_j \lambda_i y_j \cdot y_i \cdot x_i \cdot x_j + \sum_{j=1}^n \sum_{i=1}^n \lambda_j y_j - \sum_{i=1}^n \lambda_i$$

Since constraint is $\lambda_j \geq 0$, $\sum_{i=1}^n \lambda_j y_i = 0 \Rightarrow 3rd term become zero$

$$\Rightarrow f_d = \sum_{i=1}^n \lambda_i - \frac{\sum_{i,j=1}^n \lambda_j \lambda_i y_j \cdot y_i \cdot x_i \cdot x_j}{2}$$

* Simplified equation of above dual optimization Problem

$$\boxed{\max f_d = \sum_{i=1}^n \lambda_i - \frac{\lambda^T \cdot \lambda K}{2}}$$

Where, $K(i, j) = y_i \cdot y_j \cdot x_i \cdot x_j \rightarrow$ In vector, $K = y^T y \cdot x^T x$

↓ nothing but dot product of input variable x

Once we get λ from above, we can get w & b.

* Dual forms - Rewrites the same problem using a different set of variables, so the alternate formulation will help in eliminating the dependence & reducing the effect will be done by Kernelization

* Lagrange Multiplier Method :- It is a strategy to find the local minima or maxima of a function subject to the condition that one or more equations have to be satisfied exactly by the chosen values of variables.

Till now we have discussed, linear classification - Hard margin, no data points are allowed in the margin areas.

The problem with hard margin SVM is that it only works for linearly separable data also leads to underfitting.
* we can simply use optimization libraries CVXOPT to solve optimal parameters

Soft Margin SVM - Optimization :-

The problem with hard margin SVM is that it only works for linearly separable data, however, this would not be the case in the real world. It is most likely that in practical cases the data will contain some noise & might not be linearly separable. We'll look at how Soft margin SVM handles this problem.

Basically, the trick Soft margin SVM is using very simple, it adds Slack Variables ξ_i to the constraints of the optimization problem. The constraints now become:

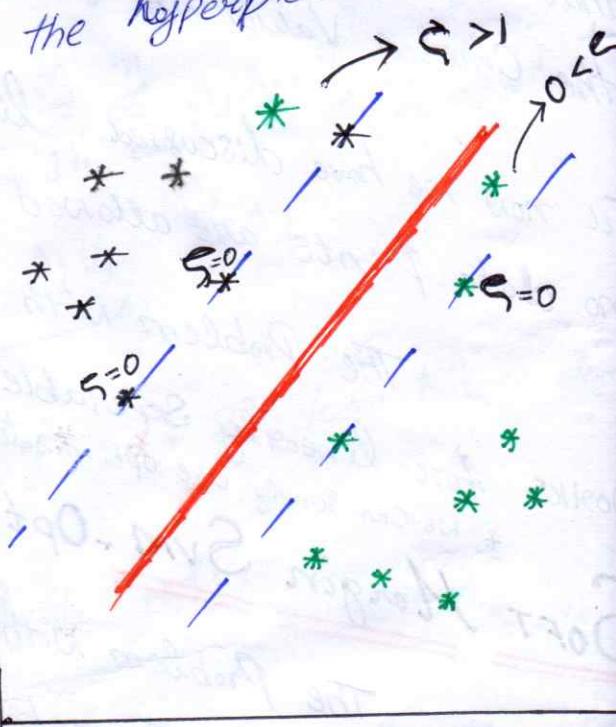
$$y_i(\omega \cdot x_i + b) \geq 1 \rightarrow \text{Hard Margin Optimization}$$

Now, the Constraints become :-

$$y_i(\omega \cdot x_i + b) \geq 1 - \xi_i ; \quad i=1 \dots n$$

- * Here, the slack variable is defined as a tolerance variable. Points that lie on the margin will have $\xi = 0$. On the other hand, the further we are from the correct boundary, the larger the value of ξ will be. Ultimately, for points that are on the wrong side of the hyperplane, we expect ξ to be greater than 1.

* By adding the slack variables, when minimizing the objective function, it is possible to choose a large enough value of ξ . So that all the examples will satisfy the constraints.



* One technique to handle this is to use Regularization. For example, we could use L2 Regularization to penalize large value of ξ . The Regularized optimization problem becomes:

$$\min_{\omega, b, \xi} \frac{1}{2} \|\omega\|^2 + \sum_{i=1}^n \xi_i$$

$$\text{Subject to } y_i(\omega \cdot x_i + b) \geq 1 - \xi_i ; \quad i=1 \dots n$$

Also, we want to make sure that we do not minimize the objective function by choosing negative values of ξ . We add the constraints $\xi_i \geq 0$.

We also add a regularization parameter C to determine how important ξ should be, which means how much we want to avoid misclassifying each training example. The regularized optimization problem becomes:

$$\min_{w, b, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

Subject to $y_i(w \cdot x_i + b) \geq 1 - \xi_i ; \xi_i \geq 0, i=1 \dots n.$

By considering these new variables in the original primal problem we get:

$$L = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i [y_i(w \cdot x_i + b) - 1 + \xi_i] - \sum_{i=1}^n \mu \xi_i$$

Here, μ is a lagrange multiplier we have defined for the tolerance values of ξ . We except each μ greater than or equal to zero.

By doing all the hard math, the optimization problem

could be transformed to a dual problem.

$$\max_{\lambda} \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j x_i \cdot x_j$$

Subject to $0 \leq \lambda_i \leq C, i=1 \dots n, \sum_{i=1}^n \lambda_i y_i = 0$

Here the constraint $\lambda_i \geq 0$, has been changed to $0 \leq \lambda_i \leq C$.

Regularization Parameter C :

C is the Parameter which Controls the Trade-off between the length of margin & number of the misclassification on the training data, in other words, it gives you control of how the SVM will handle Errors.

For $C=0$, it is not letting any misclassification to occur which is nothing but hard margin classification & it result in narrow margin,

For $C>0$, it mean no more than C, observation can violate the margin, as C increases margins also widens

Cross-Validation & it is this Parameter Only, that can result in Bias-Variance Trade-off in SVM

If the Value of $C=0$, which results in high Variance but if the Value of $C=\infty$ it results in high-Bias.

Loss-Function :

Let us Consider, the Value of ϵ_i for the case of $C \neq 0$. Because the objective will always try to minimize ϵ_i as much as possible, the Equation must hold as an Equality & we have,

$$\xi_i = \begin{cases} 1 - y_i(\omega \cdot x_i + b) & \text{if } y_i(\omega \cdot x_i + b) < 1 \\ 0 & \text{if } y_i(\omega \cdot x_i + b) \geq 1 \end{cases}$$

This is equivalent to the following closed form:

$$\xi_i = \max(1 - y_i(\omega \cdot x_i + b), 0)$$

If we plug this closed form into the objective of our SVM optimization problem, we obtain the following unconstrained version as loss function & regularizer.

$$\min_{\omega, b} \frac{\|\omega\|^2}{2} + C \sum_{i=1}^n \max[1 - y_i(\omega \cdot x_i + b), 0]$$

Hinge-loss
L2 regularizer

This formulation allows us to optimize the SVM parameters (ω, b) just like logistic regression (e.g. through gradient descent). The only difference is that we have hinge loss instead of the logistic loss. * The model produced by support vector classification depends only on a subset of training data, because the cost function for building the model does not care about training points that lie beyond the margin.

Extensions of SVM:

Before moving on, it's worth pointing out that SVM's are among the most powerful machine learning algorithms for classification tasks and are used extensively for applications ranging from "Computer Vision to NLP".

* Support-Vector Clustering → UnSupervised learning

* Multi-Class SVM → One vs one or One vs Rest

* Transductive SVM → Semi-supervised learning.

* Regression → Support Vector Regressor

* Bayesian SVM → Data Augmentation.

Pros of SVM:

- High Stability due to dependency on support vectors & not the data points
- Does not get influenced by outliers
- No assumptions made of the datasets
- Numerical Predictions Problem can be dealt with SVM

Cons of SVM:

- BlackBox method
- Inclined to Overfitting method
- Very rigorous Computation
- SVM doesn't directly provide Probability Estimates, these are calculated using an expensive five-fold cross-validation.

HyperParameters of SVC:-

`Sklearn.Svm.SVC(C=1.0, Kernel='rbf', degree=3, gamma='scale')`

↓
Regularization Parameters
↓
Kernel type
↓
Kernel
Penalty.
↓
Degree of Polynomial Kernel function
↓
Kernel Coefficient

⇒ Kernels in `['linear', 'poly', 'rbf', 'Sigmoid']`
↓ default

⇒ C in `[100, 10, 1.0, 0.1, 0.001]`:

The C regularization Parameters accept Values from 100, 10, 1.0, 0.1, 0.001. You to allow a certain level of misclassification in each training data set. Higher C regularization values lead to small margins hyperplane & do not allow much misclassification. Lower values, on the other hand, leads to high-margin & greater misclassification.

⇒ Gamma in `[1, 0.1, 0.01, 0.001, 0.0001]`

Gamma is a hyperparameter used with non-linear SVM. One of the most commonly used non-linear kernels is the radial basis function (RBF). Gamma parameter of RBF controls "the distance of the influence" of a single training point.

Low value of gamma will indicate a large similarity radius which results in more points being grouped together. For high values of gamma, the points needs to be very close to each other in order to be considered in the group (or class).

Therefore, models with very large gamma value tend to overfit.

In other words, The gamma Parameter defines the range of support vectors that will impact the positioning of the hyperplane. High gamma Value Considers only nearby data points & low value considers far away points.

gamma: { 'Scale', 'auto' } or float
default: $\frac{1}{n_features} \times \text{Var}(x)$

Gamma vs C Parameter:

For a linear Kernel, we just need to optimize the 'C' parameter. However, if we want to use an RBF Kernel, both C & gamma parameters need to optimized simultaneously. If gamma is large, the effect of C becomes negligible.

If gamma is small, C affects the model just like how it affects a linear model. Typical values for C & gamma are as follows:

$$0.0001 < \text{gamma} < 10 ; 0.1 < C < 100$$

Hyper-Tunning:

{ 'C': [0, 1, 10, 100, 1000],

Param-grid:
'gamma': [1, 0.1, 0.01, 0.001, 0.0001],
'Kernel': ['rbf'] }