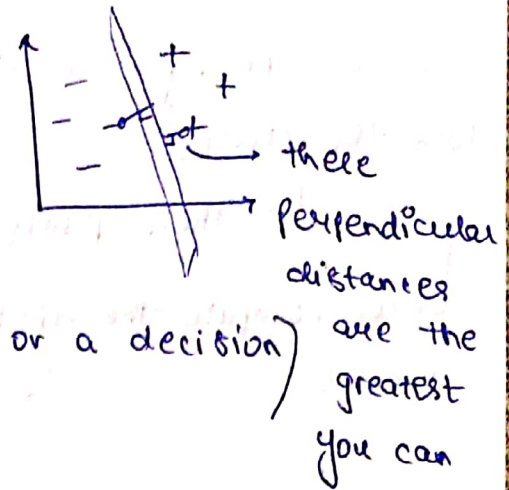


Support vector machine

used for binary classifiers

separates into two groups.



aim: to find best separating hyperplane or a decision boundary.

hyperplane: separates two sets of points
↳ a surface that acts like a boundary in an n -dim space.

for final classification, we just need to check which side of the boundary does the new image / input falls on.

use cases:

- | | |
|---------------------------------------|-----------------------|
| ① classification |] <u>supervised</u> |
| ② regression (time series prediction) | |
| ③ outlier detection |] <u>unsupervised</u> |
| ④ clustering | |

Hyperparameters:

- ① kernel: kernel helps us in finding a hyperplane in higher dimensional space without increasing the computational cost.
- ② Hyperplane: decision boundaries that classify data points into their respective classes in multi-dimensional space.

SVM is way to find the best line:

Acc. to SVM algo, we find the points closest to the line from both the classes.

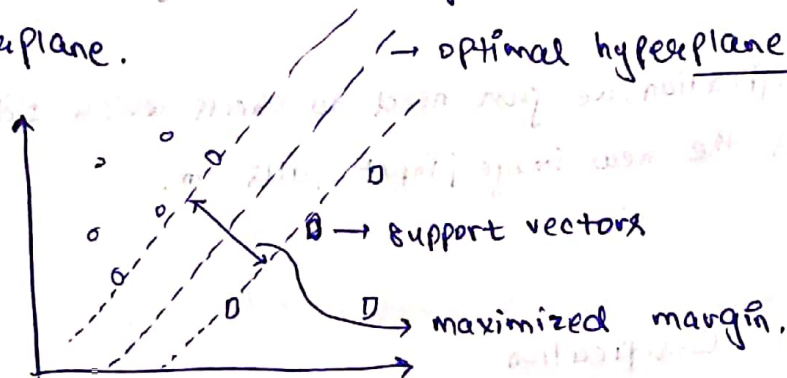
↳ These points are called "support vectors".

↳ we compute the distance between the line and the support vectors.

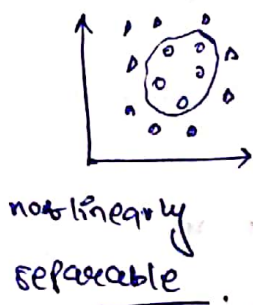
↳ This distance is called "margin".

So, our aim / goal is to maximize the margin.

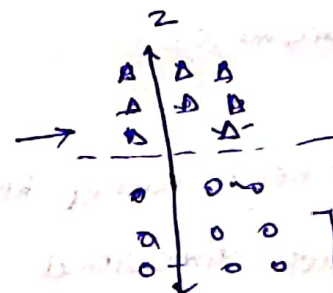
The hyperplane for which the margin is maximum is called the optimal hyperplane.



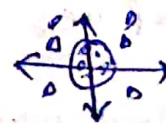
Thus SVM tries to make a decision boundary in such a way that the separation between the two classes is as wide as possible.



we cannot draw
→ a straight line.
→ but this data can
be converted to
linearly separable data
in higher dimension.



Project into 2D



$z = x^2 + y^2$ (square of distance from origin)
↙
new dimension

Thus we can classify data by adding an extra dimension to it so that it becomes linearly separable and then projecting the decision boundary back to original dimensions using mathematical transformation.

Note: Finding the correct transformation for any given dataset isn't easy.

we use kernels in sklearn to do that.

A hyperplane in an n -dimensional Euclidean space is a flat ' $n-1$ ' dimensional subset of that space that divides the space into two disconnected parts.

ex: line \rightarrow 1 dim

point \rightarrow 0 dim \rightarrow hyperplane for the line.

from sklearn.svm import SVC

clf = SVC(kernel = 'linear')

clf.fit(X, y)

optimal hyperplane will be the one with the biggest margin, because a larger margin ensures that slight deviations in the data points should not affect the outcome of the model.

Support vectors are data points that are closer to the hyperplane and influence the position and orientation of the hyperplane.

Maths:

point: $x \in \mathbb{R}^D$

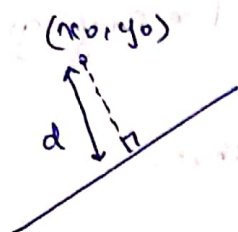
$\phi: \mathbb{R}^D \rightarrow \mathbb{R}^M$

$\phi(x) \in \mathbb{R}^M$

transform into higher dimension.

Decision boundary: $H: \omega^T \phi(x) + b = 0$
 (n-1 dim. separator)

distance measure:



line: $ax + by + c = 0$

distance of point from line

$$d = \frac{|ax_0 + by_0 + c|}{\sqrt{a^2 + b^2}}$$

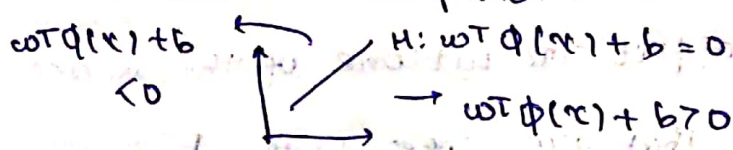
point
 \downarrow
 $d_H(\phi(x_0)) = \frac{|\omega^T \phi(x_0) + b|}{\|\omega\|_2}$

what are we optimizing:

case-1: perfect separation

$$\omega^* = \arg \max [\min d_H(\phi(x_n))]$$

maximize min. distance of support vector points to our plane.



$y_n [\omega^T \phi(x_n) + b] \rightarrow \geq 0 \rightarrow \text{correct}$

$< 0 \rightarrow \text{incorrect}$

$$\omega^* = \arg \max \frac{1}{\|\omega\|_2} [\min y_n [\omega^T \phi(x_n) + b]]$$

\rightarrow distance of closest point to H.

Primal form of SVM = $\min_{\omega} \frac{1}{2} \|\omega\|_2^2$

assumptions

data easily separable

100% accurate classification,

can lead to overfitting.

Case - 2 Non-perfect separation.

new
Primal form

$$\min \frac{1}{2} \|\omega\|_2^2 + c \cdot \sum_n \xi_n$$

Penalty term.

for correct prediction / classification $\rightarrow \xi = 0$ (zero penalty)

" incorrect " " $\rightarrow \xi \geq 1$

$$\xi_n \geq 0 \quad \forall n$$

c = hyperparameter.

[convex quadratic optimization]

$c = 0 \Rightarrow$ less complex boundary. \rightarrow may lead to underfit.

$c = \infty \Rightarrow$ even a small error is highly penalized.

more strict / complex boundary \rightarrow may lead to overfitting.

here we need to know ϕ :

use kernels to avoid use of ϕ .

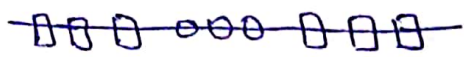
use Lagrange multipliers.

kernelization:

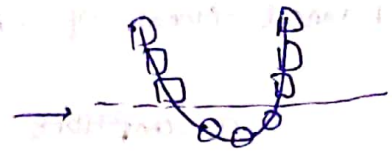
function with 2 properties:

① symmetry

② +ve semi-definite.



$$\rightarrow \phi(x) = x^2$$



linearly non-separable

Kernel Trick:

we have seen how higher dimensional transformations can allow us to separate data in order to make classification predictions.

Not all the function count as a kernel trick. The kernel function has a special property that makes it particularly useful in optimizing non-linear support vector classifiers is called a 'kernel trick'.

There might be many features in the data and applying transformations that involve many polynomial combinations of these features will lead to extremely high computation costs.

kernel methods represent the data only through a set of pairwise similarity comparisons between the original data (X) instead of explicitly applying the transformation $\phi(x)$.

here,

data set (X) is represented by an $n \times n$ kernel matrix of pairwise similarity comparisons.

The kernel function acts as a modified dot product

