

II. Dual SVM and Kernel SVM

More general formulation
of an SVM.

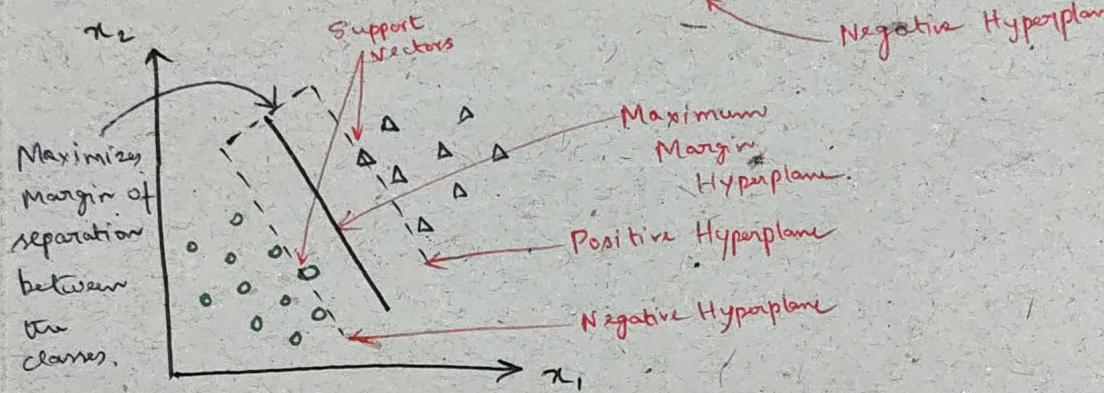
Recall, the problem to determine the SVM classifier with maximum margin is

$$\min \|\bar{a}\|$$

subject to the constraints

$$c_0 : \bar{a}^T \bar{x}_i + b \geq 1 ; i = 1, 2, \dots, M$$

$$c_1 : \bar{a}^T \bar{x}_i + b \leq -1 ; i = M+1, M+2, \dots, 2M$$



This is a Convex Optimization problem that can be solved very easily.

Let the response be defined as

$$c_0 : y_i = 1 ; i = 1, 2, \dots, M$$

$$c_1 : y_i = -1 ; i = M+1, M+2, \dots, 2M$$

The constraints can be expressed as

$$c_0 : 1 (\bar{a}^T \bar{x}_i + b) \geq 1$$

$$\Rightarrow \boxed{y_i (\bar{a}^T \bar{x}_i + b) \geq 1}$$

$$c_1 : -1 (\bar{a}^T \bar{x}_i + b) \geq 1$$

$$\Rightarrow \boxed{y_i (\bar{a}^T \bar{x}_i + b) \geq 1}$$

Constraints become same for both the classes!

Therefore, the constraints can be combined as

$$\begin{aligned} y_i (\bar{\alpha}^T \bar{x}_i + b) &\geq 1 \\ \Rightarrow -y_i (\bar{\alpha}^T \bar{x}_i + b) &\leq -1 \\ \Rightarrow \boxed{- (y_i (\bar{\alpha}^T \bar{x}_i + b) - 1)} &\leq 0 ; i = 1, 2, \dots, 2M \end{aligned}$$

This is the final version of the constraint.

Therefore, finally, the SVM classifier problem can be recast as

$$\min \|\bar{\alpha}\| \equiv \min \frac{1}{2} \|\bar{\alpha}\|^2$$

subject to the constraints

$$\boxed{- (y_i (\bar{\alpha}^T \bar{x}_i + b) - 1) \leq 0, \text{ for all } i}$$

Equivalent SVM optimization problem

Needless to say, this is also convex optimization problem.

Now, we approach the Lagrangian for this constrained optimization problem.

The Lagrangian function is

$$\frac{1}{2} \|\bar{\alpha}\|^2 + \sum_{i=1}^{2M} \lambda_i (- (y_i (\bar{\alpha}^T \bar{x}_i + b) - 1))$$

Remember, to minimize this Lagrangian; we have to take gradient w.r.t the unknown quantity $\bar{\alpha}$ and set equal to 0

$$\Rightarrow \nabla_{\bar{\alpha}} \left(\frac{1}{2} \bar{\alpha}^T \bar{\alpha} - \sum_{i=1}^{2M} \lambda_i (y_i (\bar{\alpha}^T \bar{x}_i + b) - 1) \right) = 0$$

$$\Rightarrow \nabla_{\bar{\alpha}} \left(\frac{1}{2} \bar{\alpha}^T \bar{\alpha} \right) - \sum_{i=1}^{2M} \left(\nabla_{\bar{\alpha}} (y_i \lambda_i \bar{\alpha}^T \bar{x}_i + \lambda_i y_i b - 1) \right) = 0$$

$$\Rightarrow \nabla_{\bar{\alpha}} \left(\frac{1}{2} \bar{\alpha}^T \bar{\alpha} \right) - \sum_{i=1}^{2M} \left(\nabla_{\bar{\alpha}} (y_i \lambda_i \bar{\alpha}^T \bar{x}_i) + \nabla_{\bar{\alpha}} (\lambda_i y_i b) - \nabla_{\bar{\alpha}} (1) \right) = 0$$

$$\Rightarrow \bar{\alpha} - \sum_{i=1}^{2M} \lambda_i y_i \bar{x}_i = 0$$

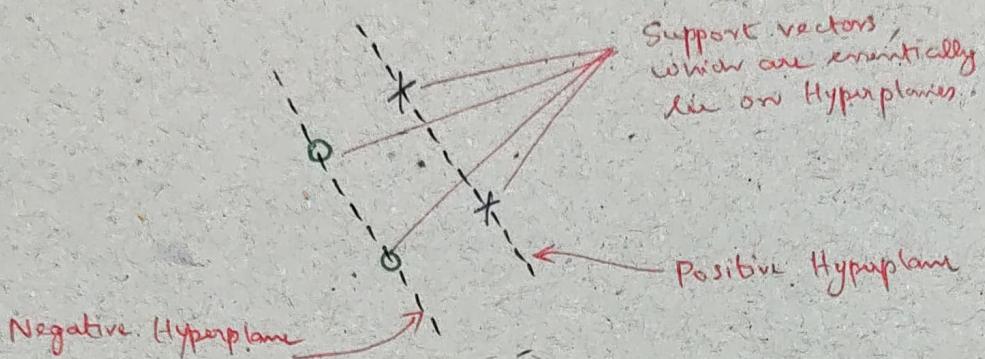
$$\Rightarrow \boxed{\bar{\alpha} = \sum_{i=1}^{2M} \lambda_i y_i \bar{x}_i}$$

$$\bar{a} = \sum_{i=1}^{2M} \lambda_i y_i \bar{x}_i$$

Linear combination of support vectors, since $\lambda_i \neq 0$.

$$\Rightarrow \bar{a} = \sum_{i=1}^{2M} \alpha_i \bar{x}_i, \text{ where } \alpha_i = \lambda_i y_i$$

Thus, \bar{a} can be expressed as linear combination of \bar{x}_i .
The points for which $\lambda_i \neq 0$ are termed as "Support Vectors".



Now, we take Gradient w.r.t. b , and set to 0.

$$\begin{aligned}
 & (\textcircled{i}) \nabla_b \left(\frac{1}{2} \bar{a}^T \bar{a} - \sum_{i=1}^{2M} \lambda_i (y_i (\bar{a}^T \bar{x}_i + b) - 1) \right) = 0 \\
 \Rightarrow & \nabla_b \left(\frac{1}{2} \bar{a}^T \bar{a} \right) - \sum_{i=1}^{2M} \left(\nabla_b (\lambda_i y_i \bar{a}^T \bar{x}_i) + \nabla_b (\lambda_i y_i b) - \nabla_b (1) \right) = 0 \\
 \Rightarrow & \boxed{\sum_{i=1}^{2M} \lambda_i y_i = 0}
 \end{aligned}$$

Now, substitute the expression for \bar{a} in the Lagrangian.

We have, the Lagrangian function,

$$\begin{aligned}
 & \frac{1}{2} \bar{a}^T \bar{a} - \sum_{i=1}^{2M} \lambda_i (y_i (\bar{a}^T \bar{x}_i + b) - 1) \\
 \Rightarrow & \frac{1}{2} \left(\sum_i \lambda_i y_i \bar{x}_i \right)^T \left(\sum_i \lambda_i y_i \bar{x}_i \right)
 \end{aligned}$$

$$- \sum_i \lambda_i \left(y_i \left(\left(\sum_i \lambda_i y_i \bar{x}_i \right)^T \bar{x}_i + b \right) - 1 \right)$$

MESSY!

$$\Rightarrow \sum_i \lambda_i - \frac{1}{2} \sum_{i,j=1}^{2M} \lambda_i \lambda_j y_i y_j \bar{x}_i^T \bar{x}_j - b \sum_i \lambda_i y_i = 0$$

$$\Rightarrow \boxed{\sum_i \lambda_i - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \bar{x}_i^T \bar{x}_j}$$

This is the Dual Objective function, which have to be maximized to formulate Dual Problem.

Remember:

For the Primal problem, we did minimization. Now,

For the Dual problem, we have to do maximization.

Therefore, the Dual problem can be formulated as

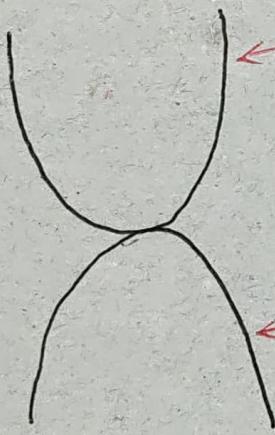
$$\max \sum_{i=1}^{2M} \lambda_i - \frac{1}{2} \sum_{i,j=1}^{2M} \lambda_i \lambda_j y_i y_j \bar{x}_i^T \bar{x}_j$$

subject to constraints

$$\textcircled{1} \quad \lambda_i \geq 0 \quad (\text{a) Lagrange Multipliers Cannot be Negative})$$

$$\textcircled{2} \quad \sum_{i=1}^{2M} \lambda_i y_i = 0.$$

This is the DUAL PROBLEM for SVM, which is equivalent to the Primal problem.



The original / Primal problem
is CONVEX

The DUAL problem is
CONCAVE.

Both have the same optimal solution.

$$\Rightarrow \text{Duality Gap} = 0.$$

How to calculate b ?

We have, the inequality

$$y_i (\bar{a}^\top \bar{x}_i + b) \geq 1$$

When $\lambda_i \neq 0$, (i) λ_i is strictly greater than 0,
this inequality becomes Equality.

$$\Rightarrow y_i (\bar{a}^\top \bar{x}_i + b) = 0$$

This is known as COMPLEMENTARY SLACKNESS property.

By solving this, one can determine b .

Note that, the quantity $\bar{x}_i^\top \bar{x}_j$ denotes the Inner Product,
which can be represented as $\langle \bar{x}_i, \bar{x}_j \rangle$.

Using this notation, the Dual SVM problem can be defined as

$$\max \sum_{i=1}^{2M} \lambda_i - \frac{1}{2} \sum_{i,j=1}^{2M} \lambda_i \lambda_j y_i y_j \underbrace{\langle \bar{x}_i, \bar{x}_j \rangle}_{\text{inner product}}$$

subject to constraints

$$\circ \lambda_i \geq 0$$

$$\circ \sum_{i=1}^{2M} \lambda_i y_i = 0$$

The Dual SVM problem
depends only on the
inner product.

For any new point \bar{x} , y can be calculated as $\bar{a}^\top \bar{x} + b$.

$$\Rightarrow \left(\sum_i \lambda_i y_i \bar{x}_i \right)^\top \bar{x} + b \quad (\because \bar{a} = \sum_i \lambda_i y_i \bar{x}_i)$$

$$\Rightarrow \sum_i \lambda_i y_i \underbrace{\langle \bar{x}_i, \bar{x} \rangle}_{\text{inner product}} + b$$

The output also depends
only on the inner product.

One can now replace $\langle \bar{x}_i, \bar{x}_j \rangle$ by a Kernel.

$$\begin{aligned} K(\bar{x}_i, \bar{x}_j) &= \phi^\top(\bar{x}_i) \phi(\bar{x}_j) \\ &= (\bar{x}_i^\top \bar{x}_j)^2 \end{aligned}$$

Previously, we were restricted to the linear operation.

Now, using a Kernel, we can go into a non-linear operation..

(ii) We can explore higher order powers of the vectors.

$$\bar{x}_i^T \cdot \bar{x}_j \rightarrow \text{Linear}$$

$$\phi \rightarrow \text{Non-linear}$$

The SVM that is obtained using a Kernel is called as Kernel SVM.

The quantity $\phi(\bar{x}_i)$ is termed as a Feature Mapping.

This can be used to model non-linear features.

Using this notation, the Kernel SVM problem can be defined as

$$\max \sum_{i=1}^{2M} \lambda_i - \frac{1}{2} \sum_{i,j=1}^{2M} \lambda_i \lambda_j y_i y_j K(\bar{x}_i, \bar{x}_j)$$

subject to constraints

$$\circ \quad \lambda_i \geq 0$$

$$\circ \quad \sum_{i=1}^{2M} \lambda_i y_i = 0.$$

Kernel SVM use non-linear feature mapping to realize non-linear features. So, it is basically separation in a non-linear space, which can be more efficient and can yield higher accuracy of class picture.

For example, for $N = 3$, (i) $\bar{x}_i = \begin{bmatrix} x_{i(1)} \\ x_{i(2)} \\ x_{i(3)} \end{bmatrix}$,

$$\text{The non-linear feature map, } \phi(\bar{x}_i) = \begin{bmatrix} x_{i(1)} & x_{i(1)} \\ x_{i(1)} & x_{i(2)} \\ x_{i(1)} & x_{i(3)} \\ x_{i(2)} & x_{i(1)} \\ x_{i(2)} & x_{i(2)} \\ x_{i(2)} & x_{i(3)} \\ x_{i(3)} & x_{i(1)} \\ x_{i(3)} & x_{i(2)} \\ x_{i(3)} & x_{i(3)} \end{bmatrix}$$

For this kernel, we need to know only Kernel!
Don't even need to know the feature mapping.

Another interesting kernel is the Gaussian Kernel defined as,

$$K(\bar{x}_i, \bar{x}_j) = \exp\left(-\frac{\|\bar{x}_i - \bar{x}_j\|^2}{2\sigma^2}\right)$$

Here, we can directly evaluate the kernel!

Advantage of Gaussian Kernel

Gaussian Kernel helps realize an infinite dimensional non-linear feature vector, which comprises of all non-linear (e.g.) all possible powers of feature vector. So, it maps it to a highly non-linear space, where it is much easier to do classification.

Example

Handwritten digit recognition, from 16×16 images.

$$2715 \Rightarrow 2715$$

Gaussian Kernel SVM's yield very good performance!

The Dual SVM is exploited to develop Kernel SVM, which replaces the inner product with a more general Kernel (inner product between the feature maps of the original feature).
