

S.V.M. - Mathematics

Objective: Optimal hyperplane which linearly separates the data points in 2 components by maximizing the margin.

Hyperplane: linearly divides n -dimensional data points in 2 components

$$W = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \quad x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

~~hyperplane~~

Let $x = (x, y)$ & $W = (a, -1)$ then hyperplane is

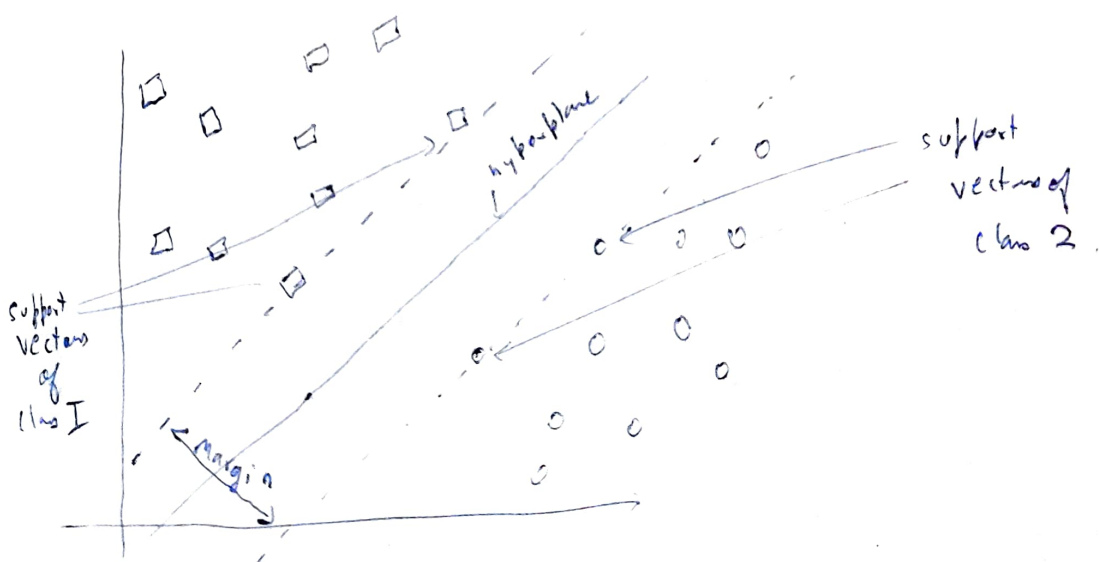
$$W \cdot X + b = 0$$

$$\boxed{W^T X = 0} \leftarrow n\text{-dimensional hyper-plane.}$$

If data is not linearly separable. We add extra dimension like $z = x^2 + y^2$ to make data separable.

Kernel Trick.

• Many Hyperplanes. Need to find Optimal Hyperplane



Margin: If solid line is optimal hyperplane of 2 dotted line are some hyperplane passing through nearest data points of optimal hyperplane. Then distance between hyperplane & optimal hyperplane is margin.

Support Vectors: Closest data-points to optimal hyperplane which help for formulation of other 2 hyperplane

When optimal hyperplane^{is} selected we choose hyperplane which has highest distance from closest data points.

Goal is to maximize margin while selecting optimal hyperplane.

• Let there be L training eg's

• Each eg x is D dimensional & has label $y \in \{+1, -1\}$;
(2 class).

& it is linearly separable data.

Training data: $\{x_i, y_i\}$ where $i = 1, \dots, L$
 $y_i \in \{-1, 1\}$ &
 $x \in \mathbb{R}^D$

Goal of SVM is to orient the hyperplane as far as possible from closest members of both classes.

$$w \cdot x_i + b \geq 1 \quad \forall i : y_i = +1$$

$$w \cdot x_i + b \leq -1 \quad \forall i : y_i = -1$$

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

$$H1: w \cdot x + b = -1$$

$$H2: w \cdot x + b = 1$$

} Two hyperplanes passing through support vectors of class 1 & -1

Distance of optimal hyperplane from origin is

$$\frac{-b}{|w|}$$

Then distance of H_1 is $\frac{-1-b}{|w|}$

& H_2 is $\frac{1-b}{|w|}$

~~Then~~ margin

$$M = \frac{1-b}{|w|} - \frac{-(-1-b)}{|w|} = \frac{2}{|w|}$$

Margin is $\frac{M}{2}$ as M is distance between H_1 & H_2 .

$$\text{Hence margin} = \frac{1}{|w|}$$

To find optimal hyperplane $\text{objective function} = \max \frac{1}{||w||}$

$$\text{or } \min ||w||$$

$$\text{or } \min \frac{||w||^2}{2} \text{ s.t. } y_i(x_i \cdot w) \geq 1$$

$$y_i(w \cdot x_i + b) - 1 \geq 0 \quad \forall i=1, \dots, l$$

$$\therefore \text{Objective function} : \min \frac{||w||^2}{2}$$

$$\text{Constraint} : y_i(w \cdot x_i + b) - 1 \geq 0 \quad \forall i$$

For constrained Optimization we use Lagrange Multiplier

Do Not use gradient descent as for constrained optimization better dual optimization algorithm preferred.

Problem $\min \frac{\|w\|^2}{2} \quad \text{s.t. } y_i(w \cdot x_i + b) - 1 \geq 0$

Solⁿ:

Using Lagrange Multipliers

$\forall i \in \{1, \dots, l\}$

$$\begin{aligned} \mathcal{L} &= \frac{\|w\|^2}{2} - \sum_{i=1}^l \lambda_i (y_i (w \cdot x_i + b) - 1) \\ &= \frac{\|w\|^2}{2} - \sum_{i=1}^l \lambda_i (y_i (w \cdot x_i + b)) + \sum_{i=1}^l \lambda_i \end{aligned}$$

$$\frac{\partial \mathcal{L}}{\partial w} = w - \sum_{i=1}^l \lambda_i (y_i x_i) = 0 \quad \text{--- (1)}$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_i} = y_i (w \cdot x_i + b) - 1 = 0$$

$$\frac{\partial \mathcal{L}}{\partial b} = \sum_{i=1}^l \lambda_i \cdot y_i = 0 \quad \text{--- (2)}$$

Primal Optimization Solⁿ. We cannot get Solⁿ. So we.

Dual optimization formulation

Computing Problem to Dual Problem

$$ld = \frac{\|w\|^2}{2} - \sum_{i=1}^l \lambda_i (y_i (w \cdot x_i + b)) + \sum_{i=1}^l \lambda_i$$

Using (1)

$$\begin{aligned} ld &= \frac{\sum_{i=1}^l \|\lambda_i \cdot y_i \cdot x_i\|^2}{2} - \sum_{j=1}^l \sum_{i=1}^l \lambda_j y_j (\|\lambda_i y_i x_i\| \cdot x_j + b) + \sum_{i=1}^l \lambda_i \\ &= \frac{\sum_{i=1}^l (\lambda_i y_i x_i)^T (\lambda_i y_i x_i)}{2} - \sum_{j=1}^l \sum_{i=1}^l \lambda_j \lambda_i y_j y_i x_i x_j + \sum_{j=1}^l \sum_{i=1}^l \lambda_j y_j - \sum_{i=1}^l \lambda_i \end{aligned}$$

Using (2) $\sum_{j=1}^l \sum_{i=1}^l \lambda_j y_j = 0$

$$ld = \sum_{i=1}^l \lambda_i - \frac{\sum_{i,j=1}^l \lambda_j \lambda_i y_j y_i x_i x_j}{2} \quad (9)$$

\therefore Simplified eqn of dual opti problem...

$$K(i,j) = y_i y_j x_i x_j \quad K = y^T y \cdot x^T x$$

$$\therefore \max ld = \sum_{i=1}^l \lambda_i - \frac{\lambda^T \lambda K}{2}$$

So Problem is to find λ to $\max^m ld$

After optimization using SMO (sequential minimization optimization)

we get value of λ .

Using that $w = \sum_{i=1}^l \lambda_i y_i x_i$

Now with w & b we calculate b .

Any support vector x_s will have

$$y_s(w \cdot x_s + b) - 1 = 0$$

Substituting w .

$$y_s \left(\sum_{m \in S} \lambda_m y_m x_m \cdot x_s + b \right) - 1 = 0$$

S denotes indices of support vectors
(as dot product)

$$y_s y_s \left(\sum_{m \in S} \lambda_m y_m x_m \cdot x_s + b \right) = y_s$$

$$y_s y_s = 1 \text{ as } (\pm 1)^2 = 1$$

$$\sum_{m \in S} \lambda_m y_m x_m \cdot x_s + b = y_s$$

$$b = y_s - \sum_{m \in S} \lambda_m y_m x_m \cdot x_s$$

Use avg of all support vectors.

$$b = \sum_{s \in S} \frac{\left(y_s - \sum_{m \in S} \lambda_m y_m x_m \cdot x_s \right)}{N_S}$$

where S is set of all support vectors

Now we have value of both w & b

Optimal hyperplane is

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

$$\text{Predict}(x) = \underset{\substack{\uparrow \\ \text{(sign function)}}}{\text{sign}(w \cdot x + b)}$$

* Previous solⁿ is based on fact that no datapoint is allowed inside data-points (support-vectors).

Hard Margin solⁿ

Now we will try soft margin solⁿ
i.e. \swarrow few data points are allowed inside margin.

- Soft margin - underfitting.
- Hard margin - Overfitting

Mathematically, we relax margin by introducing a positive slack variable ξ_i

$$w \cdot x_i + b \geq 1 - \xi_i \quad \text{for } (y_i = +1)$$

$$w \cdot x_i + b \leq -1 + \xi_i \quad \text{for } (y_i = -1)$$

$$y_i(w \cdot x_i + b) - 1 + \xi_i \geq 0 \quad \text{for } y_i = \pm 1$$

$$\xi_i \geq 0 \quad \text{for } i = 1, 2, \dots, l$$

$$\sum_{i=1}^l \xi_i = C$$

- γ_i tells where i^{th} observation is located relative to hyperplane & margin. ~~also $\gamma_i \leq 0$~~
- if $0 \leq \gamma_i \leq 1$ - observation is between inner side of margin & correct side of hyperplane. (Margin violation)
- If $\gamma_i > 1$, observation is on incorrect side of both hyperplane & margin.
- Points on wrong side of margin have penalty which increases with distance.
- C is parameter that controls trade-off between length of margin & num of misclassification.
- $C = 0 \Rightarrow$ hard margin classification
- $C > 0$, ~~more~~ means no more than C observation can violate margin. $C \uparrow$ margin \uparrow

$$J_0 = \min \frac{\|w\|^2}{2} + C \sum_{i=1}^l \gamma_i \quad \text{s.t.}$$

$$y_i (w \cdot x_i + b) - 1 + \gamma_i \geq 0; \gamma_i \geq 0$$

• Primal Gradient based optimization method

- Make optimization problem into unconstrained form - without Lagrange Multipliers.
- Then solve using ~~the~~ gradient descent etc.

$$J = \min \frac{\|w\|^2}{2} + C \sum_{i=1}^I \eta_i \quad \text{s.t.}$$

$$y_i (w \cdot x_i + b) - 1 + \eta_i \geq 0, \eta_i \geq 0$$

$$\text{let } f(x_i) = w \cdot x_i + b$$

$$\therefore y_i f(x_i) \geq 1 - \eta_i$$

$$\Rightarrow \eta_i \geq 1 - y_i f(x_i)$$

$$\text{also } \eta_i \geq 0$$

$$\Rightarrow \eta_i = \max(1 - y_i f(x_i), 0)$$

$$\text{Then } J = \min \left(\underbrace{\frac{\|w\|^2}{2}}_{\text{regularization term}} + C \underbrace{\sum_{i=1}^I \max(1 - y_i f(x_i), 0)}_{\text{loss function term}} \right)$$

$$\therefore c(w, b)_{\min(w, b)} = \frac{\lambda}{2} \|w\|^2 + \frac{1}{N} \sum_{i=1}^n \max(0, 1 - y_i f(x_i))$$

$$\text{where } \lambda = \frac{2}{nC} \quad \& \quad f(x) = w \cdot x + b$$

• Can use gradient descent on loss function

Dual Quadratic programming approach
(same as Hard margin soln).

$$\max d = \sum_{i=1}^l \lambda_i - \frac{\lambda^T \lambda K}{2}$$

$$0 \leq \lambda \leq C, \quad \sum_{i=1}^l \lambda_i y_i = 0$$

Dual Optimization

- Not scalable for large dataset.
- Kernel trick can be applied.

Gradient Descent

- Converges easily.
- Online learning

For Non Linear Data Points

Mercer's Theorem: If function $k(a,b)$ satisfies all constraints called mercer's constraints, then there exists a function that maps a & b into higher dimension.

$$k(a,b) = \phi(a)^T \cdot \phi(b)$$

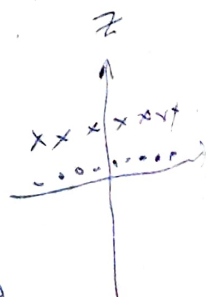
$\phi()$ is a kernel

Kernels: 'linear', 'polynomial', 'radial basis function'



$$z = x^2 + y^2$$

2D data $\left(\because (\log x) \right)$ is more distant z will always be higher for $(\log x)$. But z



Dual optimization problem of ^(soft) margin SVM.

$$\max_{\lambda} L_d = \sum_{i=1}^l \lambda_i - \frac{\lambda^T \lambda K}{2} \quad \text{subject to}$$

$$0 \leq \lambda \leq C, \quad \sum_{i=1}^l \lambda_i y_i = 0.$$

$$K(i, j) = y_i y_j x_i x_j$$

$$\& \because y_i y_j = 1 \text{ so}$$

$$K(i, j) = x_i x_j$$

↑
Linear Kernel.

$$K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$$

(ϕ is kernel)

Notes: Instead of this we can transfer input data also straight but it will increase computational cost & space by a lot. Kernel trick helps optimize code.

Gaussian Radial Basis Function: Best when no prior knowledge of data.

$$K(x_i, x_j) = e^{-\gamma \|x_i - x_j\|^2}$$

Gaussian function:

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

Polynomial Kernel function

$$K(x_i, x_j) = (x_i \cdot x_j + a)^b$$

Given kernel

$$k(x_i, x_j) = x_i \cdot x_j$$

Choice of SVM. kernel: linear \prec poly \prec radt.

ability to fit any data linear \prec poly \prec radt.

risk of overfitting: linear \prec poly \prec radt.

" " underfitting: radt \prec poly \prec linear.

no. of hyperparameters: linear(0) \prec radt(2) \prec

poly(3)

Case I (a). Support vectors x_p & x_a .