

Support Vector Machine: Classification

S. Sumitra

Department of Mathematics

Indian Institute of Space Science and Technology

Advanced Machine Learning

Support Vector Machines

- Vapnik and A. Chervonenkis in 1962 at AT& T Bell Laboratories
- Bernhard E. Boser, Isabelle M. Guyon and Vladimir N. Vapnik in 1992
- Corinna Cortes and Vapnik in 1993

Hypothesis Function

- $\tilde{f}(x_i) = f(x_i) + b$
- f lies in Reproducing Kernel Hilbert Space \mathcal{F} , $b \in \mathbb{R}$

Optimisation

$$\begin{aligned} & \min_{f \in \mathcal{F}} \frac{1}{2} \|f\|^2 \\ \text{subject to } & \frac{1}{N} \sum_{i=1}^N V(y_i, \tilde{f}(x_i)) \leq k \end{aligned}$$

Classification

Data $\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$, $x_i \in \mathbb{R}^n$, $y_i \in \{1, -1\}$ be the given data.

The SVM uses regularized cost function

$$J(\tilde{f}, b) = C \sum_{i=1}^N V(y_i, \tilde{f}(x_i)) + \frac{1}{2} \|f\|^2 \quad (1)$$

where $C > 0$ is the regularization parameter.

Prediction with Confidence

- $\tilde{f}(x) \geq 1, y = 1$, Correct Prediction, high confidence
- $\tilde{f}(x) < -1, y = -1$, Correct Prediction, high confidence
- $0 < \tilde{f}(x) < 1, y = 1$, Correct Prediction, low confidence
- $-1 < \tilde{f}(x) < 0, y = -1$, Correct Prediction, low confidence

SVM Loss Function: Hinge Loss Function

$\{x_i, y_i\}_{i=1}^N, x_i \in \mathbb{R}^n, y_i = \{-1, 1\}$ be the data points

- For correct prediction with high confidence
 - $\tilde{f}(x_i) \geq 1, y_i = 1$
 - $\tilde{f}(x_i) \leq -1, y_i = -1$
 - $y_i \tilde{f}(x_i) \geq 1$
- Hinge Loss Function $V(y, \tilde{f}(x)) = \max(0, 1 - y\tilde{f}(x))$.
- In this case, $V = 0$ for correct prediction with high confidence. Else, $V > 0$.

SVM: Signum Function

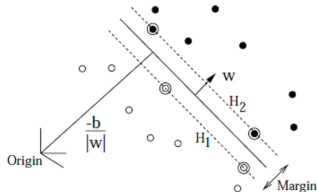
- Signum function is used to determine the class

$\tilde{f}(x) \geq 0$, x is assigned to positive class $\tilde{f}(x) < 0$, x is assigned to negative class

SVM Development

- Linear and loss equal to 0
- Linear and loss not equal to 0
- Nonlinear

Case 1: Hyperplane and $V \equiv 0$



- $\tilde{f}(x) = f(x) + b = \langle w, x \rangle + b$
- $V(y, \tilde{f}(x)) = \max(0, 1 - y\tilde{f}(x))$.
- $V \equiv 0$

Case 1: Problem Formulation

$$\begin{aligned} & \underset{f \in \mathcal{F}}{\text{minimize}} && \|f\|^2 \\ & \text{subject to} && 1 - y_i \tilde{f}(x_i) \leq 0, \quad i = 1, \dots, N. \end{aligned}$$

As $f(x) = \langle w, x \rangle$, $\|f\| = \|w\|$, the above optimization problem can be formulated as

$$\begin{aligned} & \underset{w \in \mathbb{R}^n}{\text{minimize}} && \|w\|^2 \\ & \text{subject to} && y_i(\langle w, x_i \rangle + b - 1) \geq 1, \quad i = 1, \dots, N. \end{aligned}$$

It can be proved that in this case finding a function which is smooth (minimum $\|w\|$) is equivalent to finding a hyperplane with largest margin.

Geometrical Interpretation of Cost Function

Let the shortest distance from the separating hyperplane $w^T x + b = 0$ to the positive class be d^+ , while to the negative class be d^- .

The *margin* of a separating hyperplane is defined to be $\gamma = d^+ + d^-$.

[The distance of a point (x', y') to the plane $ax + by = 0$ is given by $\frac{|ax' + by'|}{\sqrt{a^2 + b^2}}$]

The distance of a point $x_i = (x_{i1}, x_{i2}, \dots, x_{in})$ to the hyperplane $w^T x + b = 0$ is given by

$$\begin{aligned} \text{dist}(x_i, \tilde{f}(x)) &= \frac{|w_1 x_{i1} + w_2 x_{i2} + \dots + w_n x_{in} + b|}{\sqrt{w_1^2 + w_2^2 + \dots + w_n^2}} \\ &= \frac{|w^T x_i + b|}{\|w\|} \end{aligned}$$

$w^T x + b \geq 1$ for all x in the positive class. Therefore $dist(x, \tilde{f}(x)) \geq \frac{1}{\|w\|}$, for all x in the positive class. The points in the positive class, that lie closest to the separating hyperplane $w^T x + b = 0$ lies in the hyperplane

$$H_1 : w^T x + b = 1$$

Therefore $d^+ = \frac{1}{\|w\|}$.

In the negative class, for all points $w^T x + b \leq -1$. Therefore $dist(x, \tilde{f}(x)) \geq \frac{1}{\|w\|}$ for all x in the negative class. The points in the negative class, that lie closest to the separating hyperplane $w^T x + b = 0$ lies in the hyperplane

$$H_2 : w^T x + b = -1$$

Therefore $d^- = \frac{1}{\|w\|}$.

Hence margin $\gamma = d^+ + d^- = \frac{2}{\|w\|}$

There will be more than one hyperplane that separates the given data points. Among all such hyperplanes, the support vector machine chooses that one, which has the largest margin. That is *for the linearly separable case, the support vector algorithm chooses for the separating hyperplane with largest margin*. That is to $\max \frac{2}{\|w\|}$, which is equivalent to $\min \frac{\|w\|}{2}$, which is equivalent to $\min \frac{\|w\|^2}{2}$.
Now

$$w^T x_i + b \geq +1 \text{ for } y_i = +1 \quad (2)$$

$$w^T x_i + b \leq -1 \text{ for } y_i = -1 \quad (3)$$

The above equations can be combined into one set of inequalities:

$$y_i(w^T x_i + b) \geq 1 \quad \forall i \quad (4)$$

Here the hyperplanes H_1 and H_2 are parallel and that no training points fall between them. Thus we can find the pair of hyperplanes which gives the maximum margin by minimizing $\|w\|^2$, subject to constraints (4).

Optimisation Problem

The general optimisation problem can be stated as follows:

Definition

(Primal optimisation problem) Given functions f , g_i , $i = 1, 2, \dots, k$, and h_i , $i = 1, 2, \dots, m$, defined on a domain $\Omega \subseteq \mathbb{R}^n$,

$$\begin{aligned} & \text{minimise} && f(w), w \in \Omega, \\ & \text{subject to} && g_i(w) \leq 0, i = 1, \dots, k, \\ & && h_i(w) = 0, i = 1, \dots, m \end{aligned}$$

where $f(w)$ is called the objective function, and the remaining relations are called, respectively, the inequality and equality constraints. The optimal value of the objective function is called the value of the optimisation problem.

Maximisation problems can be converted to minimisation ones by changing the the sign of $f(w)$.

Feasible Region

[Refer: An Introduction to Support Vector Machines and Other Kernel-based Learning Methods by Nello Cristianini and John Shawe-Taylor]

Definition

(feasible region) The region of the domain where the objective function is defined and where all the constraints are satisfied is called the feasible region, which can be denoted by

$$S = \{w \in \Omega : g(w) \leq 0, h(w) = 0\}$$

A solution of the optimisation problem is a point $w^* \in S$ such that there exists no other point $w \in S$ for which $f(w) < f(w^*)$. Such a point is also known as a global minimum.

Active Constraint

An inequality constraint $g_i(w) \leq 0$ is said to be active (or tight) if the solution w^* satisfies $g_i(w^*) = 0$, otherwise it is said to be inactive. Hence, equality constraints are always active.

Convex Function

Definition

A real-valued function $f(w)$ is called convex for $w \in \mathbb{R}^n$ if, $\forall w, u \in \mathbb{R}^n$, and for any $\theta \in (0, 1)$,

$$f(\theta w + (1 - \theta)u) \leq \theta f(w) + (1 - \theta)f(u)$$

If a strict inequality holds, the function is said to be strictly convex. A function that is twice differentiable will be convex provided its Hessian matrix is positive semi-definite. An affine function is one that can be expressed in the form

$$f(w) = Aw + b$$

for some matrix A and vector b . Note that affine functions are convex as they have zero Hessian.

Convex Set

Definition

A set $\Omega \subseteq \mathbb{R}^n$ is called convex if, $\forall w, u \in \Omega$, and for any $\theta \in (0, 1)$, the point $(\theta w + (1 - \theta)u) \in \Omega$.

If a function f is convex, any local minimum w^* of the unconstrained optimisation problem with objective function f is also a global minimum. To prove this take $u \neq w^*$. Then by the definition of a local minimum there exists θ sufficiently close to 1 such that

$$\begin{aligned} f(w^*) &\leq f(\theta w^* + (1 - \theta)u) \\ &\leq \theta f(w^*) + (1 - \theta)f(u) \\ f(w^*)(1 - \theta) &\leq f(u)(1 - \theta) \end{aligned}$$

Therefore $f(w^*) \leq f(u)$.

Convex Quadratic Programmes

Definition

An optimisation problem in which the set Ω , the objective function and all of the constraints are convex is said to be convex.

For SVM, here we are considering the case where the constraints are linear, the objective function is convex and quadratic and $\Omega = \mathbb{R}^n$, that is, convex quadratic programmes.

Lagrangian Theory

The Lagrangian theory was developed by Lagrange in 1797 for determining the solution of an optimisation problem which there are no inequality constraints. In 1951 Kuhn and Tucker extended the theory by incorporating inequality constraints which is known as Kuhn-Tucker theory.

Definition

(Fermat) A necessary condition for w^ to be a minimum of $f(w)$, $f \in C^1$, is $\frac{\partial f(w)}{\partial w^*} = 0$. This condition, together with convexity of f , is also a sufficient condition.*

The Lagrangian is defined as the objective function plus a linear combination of the constraints, where the coefficients of the combination are called the Lagrange multipliers.

Lagrangian Theory

Definition

Given an optimisation problem with objective function $f(\mathbf{w})$, and equality constraints $h_i(\mathbf{w}) = 0, i = 1, \dots, m$, we define the Lagrangian function as

$$L(\mathbf{w}, \beta) = f(\mathbf{w}) + \sum_{i=1}^m \beta_i h_i(\mathbf{w})$$

where the coefficients β_i are called the Lagrangian multipliers.

Theorem

(Lagrange) A necessary condition for a normal point w^ to be a minimum of $f(w)$ subject to $h_i(w) = 0, i = 1, \dots, m$, with $f, h_i \in C^1$, is*

$$\frac{\partial L(w^*, \beta^*)}{\partial w} = 0$$

$$\frac{\partial L(w^*, \beta^*)}{\partial \beta} = 0$$

for some values β^* . The above conditions are also sufficient provided that $L(w, \beta^*)$ is a convex function of w .

Lagrangian Function

Now let us consider generalized Lagrangian. The definition of generalised Lagrangian function is:

Definition

Given an optimisation problem with domain $\Omega \subseteq \mathbb{R}^n$,

$$\begin{aligned} & \text{minimise} && f(\mathbf{w}) \quad , \mathbf{w} \in \Omega \\ & \text{subject to} && g_i(\mathbf{w}) \leq 0, i = 1, 2, \dots, k \\ & && h_i(\mathbf{w}) = 0, i = 1, 2, \dots, m \end{aligned}$$

The generalized lagrangian function is defined as

$$L(\mathbf{w}, \alpha, \beta) = f(\mathbf{w}) + \sum_{i=1}^k \alpha_i g_i(\mathbf{w}) + \sum_{i=1}^m \beta_i h_i(\mathbf{w}), \alpha_i \geq 0$$

[For constraints of the form $c_i \geq 0$, the constraint equations are multiplied by positive Lagrange multipliers and subtracted from the objective function, to form the Lagrangian. For equality constraints, the Lagrange multipliers are unconstrained].

Dual Problem

Definition

The Lagrangian dual problem of the above given primal problem is the following problem:

$$\begin{array}{ll} \text{maximize} & \theta(\alpha, \beta) \\ \text{subject to} & \alpha \geq 0 \end{array}$$

where $\theta(\alpha, \beta) = \inf_{w \in \Omega} L(w, \alpha, \beta)$. The value of the objective function at the optimal solution is called the value of the problem.

Theorem

Let $w \in \Omega$ be a feasible solution of the primal problem and (α, β) a feasible solution of the dual problem. Then $f(w) \geq \theta(\alpha, \beta)$.

$$\begin{aligned}\theta(\alpha, \beta) &= \inf_u L(u, \alpha, \beta) \\ &\leq L(w, \alpha, \beta) \\ &= f(w) + \alpha'g(w) + \beta'h(w) \\ &\leq f(w)\end{aligned}$$

- If $f(w^*) = \theta(\alpha^*, \beta^*)$, where $\alpha^* \geq 0$, and $g(w^*) \leq 0$, $h(w^*) = 0$, then w^* and (α^*, β^*) solve the primal and dual problems respectively. In this case $\alpha_i^* g_i(w^*) = 0$, for $i = 1, \dots, k$.
- The difference between the values of the primal and dual problems is known as the duality gap.
- The triple (w^*, α^*, β^*) is a saddle point of the Lagrangian function for the primal problem, if and only if its components are optimal solutions of the primal and dual problems and there is no duality gap, that is, $f(w^*) = \theta(\alpha^*, \beta^*)$.

Strong Duality Theorem

Theorem (strong duality theorem)

Given an optimisation problem with convex domain $\Omega \subseteq \mathbb{R}^n$,

$$\begin{aligned} & \text{minimise} && f(w) \quad , w \in \Omega \\ & \text{subject to} && g_i(w) \leq 0, i = 1, 2, \dots, k \\ & && h_i(w) = 0, i = 1, 2, \dots, m \end{aligned}$$

with $f \in C_1$ convex and g_i, h_i affine, then the duality gap is zero.

The Karush-Kuhn-Tucker Conditions (KKT)

(Kuhn-Tucker) Given an optimisation problem with convex domain $\Omega \subseteq \mathbb{R}^n$,

$$\begin{aligned} & \text{minimise} && f(w) \quad , w \in \Omega \\ & \text{subject to} && g_i(w) \leq 0, i = 1, 2, \dots, k \\ & && h_i(w) = 0, i = 1, 2, \dots, m \end{aligned}$$

with $f \in C_1$ convex and g_i, h_i affine, necessary and sufficient conditions for a normal point w^* to be an optimum are the existence of α^*, β^* such that

$$\begin{aligned} \frac{\partial(L(w^*, \alpha^*, \beta^*))}{\partial w} &= 0 \\ \frac{\partial(Lw^*, \alpha^*, \beta^*)}{\partial \beta} &= 0 \\ \alpha_i^* g_i(w^*) &= 0, i = 1, 2, \dots, k \\ g_i(w^*) &\leq 0, i = 1, 2, \dots, k \end{aligned}$$

SVM: Case 1: Optimisation Problem

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle - \sum_{i=1}^N \alpha_i [y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1] \quad (5)$$

where $\alpha_i \geq 0$ are the Lagrange multipliers.

The corresponding dual is found by differentiating with respect to \mathbf{w} and b , imposing stationarity,

$$\frac{\partial L(\mathbf{w}, b, \alpha)}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^N y_i \alpha_i \mathbf{x}_i = 0 \quad (6)$$

$$\frac{\partial L(\mathbf{w}, b, \alpha)}{\partial b} = \sum_{i=1}^N y_i \alpha_i = 0 \quad (7)$$

From (6)

$$\mathbf{w} = \sum_{i=1}^N y_i \alpha_i \mathbf{x}_i \quad (8)$$

$$f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle = \sum_{i=1}^N y_i \alpha_i \langle \mathbf{x}_i, \mathbf{x} \rangle$$

From the above equation it is clear that the hypothesis can be written as a linear combination of the training points [which also follows from the representer theorem].

From (7)

$$0 = \sum_{i=1}^N y_i \alpha_i \quad (9)$$

KKT Complimentary Condition

Using Karush-Kuhn-Tucker complementarity conditions,

$$\alpha_i^* [y_i(\langle \mathbf{w}, x_i \rangle + b) - 1] = 0, i = 1, 2, \dots, N \quad (10)$$

Therefore two cases

- $\alpha_i^* = 0$
 - $y_i(\langle \mathbf{w}, x_i \rangle + b) - 1 \geq 0$
 - x_i lies either above H_1 or on H_1 or above H_2 or on H_2
- $\alpha_i^* \neq 0$
 - $y_i(\langle \mathbf{w}, x_i \rangle + b) - 1 = 0$
 - x_i lies either on H_1 or H_2

Support Vectors

- The x_i 's corresponding to nonzero α_i 's are called support-vectors. Thus the only inputs that appear in the expression of w , and hence in the expression of hyperplane are the support vectors.
- Support vectors lie closest to the hyperplane.
- In the expression for the weight vector only support vectors are involved. It is for this reason that they are called support vectors.

Dual Formulation

Sub (8) and (9) into (5) to form the dual

$$\begin{aligned}\|w\|^2 &= \langle w, w \rangle \\ &= \left\langle \sum_i \alpha_i y_i x_i, \sum_i \alpha_i y_i x_i \right\rangle \\ &= \sum_i \sum_j \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle\end{aligned}$$

$$\begin{aligned}\sum_{i=1}^N \alpha_i [y_i (\langle w, x_i \rangle)] &= \sum_{i=1}^N \alpha_i [y_i (\langle \sum_{j=1}^N y_j \alpha_j x_j, x_i \rangle)] \\ &= \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle\end{aligned}$$

Collecting the terms of b ,

$$b \sum_i \alpha_i y_i = 0$$

Dual Formulation

Sub (8) and (9) into (5) to form the dual,

$$\begin{aligned} W(\alpha) &= \frac{1}{2} \sum_{i,j=1}^N y_i y_j \alpha_i \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle - \sum_{i,j=1}^N y_i y_j \alpha_i \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^N \alpha_i \\ &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N y_i y_j \alpha_i \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \end{aligned}$$

Hence the dual optimisation problem is

$$\text{maximise } W(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N y_i y_j \alpha_i \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

$$\text{subject to } \sum_{i=1}^N y_i \alpha_i = 0$$

$$\alpha_i \geq 0, i = 1, 2, \dots, N$$

$L(\mathbf{w}, \alpha, b)$ and $W(\alpha)$ arise from the same objective function but with different constraints and the solution is found by minimizing $L(\mathbf{w}, \alpha, b)$ or by maximizing $W(\alpha)$.

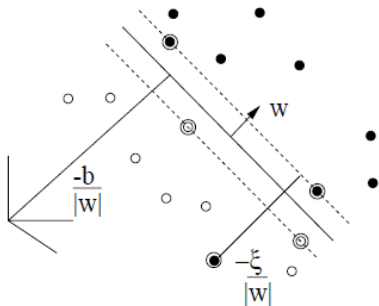
Bias

To find b , choose x_i for which $\alpha_i \neq 0$ and use (10) [take the average of b 's over all such x_i 's].

Test Phase

- The class of the testing point x is determined by $\text{sgn}(w^T x + b)$, that is,
 - x belongs to the positive class if $\text{sgn}(w^T x + b) \geq 0$
 - Else to the negative class

Case 2: Hyperplane and $V \neq 0$



$$\text{minimise}_{w \in \mathbb{R}^n, b \in \mathbb{R}} \frac{\|w\|^2}{2} + C \sum_{i=1}^N \max(0, 1 - y_i(w^T x_i + b))$$

$$V(y, \tilde{f}(x_i)) = \max(0, 1 - y_i \tilde{f}(x_i)) = \xi_i, i = 1, 2, \dots, N$$

Optimization

$$\text{minimise}_{w \in \mathbb{R}^n, b \in \mathbb{R}} \frac{\|w\|^2}{2} + C \sum_{i=1}^N \xi_i$$

where $C > 0$ is the regularization parameter, subject to the constraints

$$y_i(w^T x_i + b) \geq +1 - \xi_i \text{ for } i = 1, 2, \dots, N \quad (11)$$

$$\xi_i \geq 0 \text{ } i = 1, 2, \dots, N \quad (12)$$

Lagrangian Formulation

$$L(\mathbf{w}, \mathbf{b}, \xi, \alpha) = \frac{\|\mathbf{w}\|^2}{2} + C \sum_{i=1}^N \xi_i - \sum_i \alpha_i [y_i(\mathbf{w}^T x_i + \mathbf{b}) - 1 + \xi_i] - \sum_i \mu_i \xi_i \quad (13)$$

Applying the KKT conditions,
For $i = 1, 2, \dots, N$

$$\frac{\partial L(\mathbf{w}, \mathbf{b}, \alpha)}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^N y_i \alpha_i \mathbf{x}_i = 0 \quad (14)$$

Therefore,

$$\mathbf{w} = \sum_{i=1}^N y_i \alpha_i \mathbf{x}_i$$
$$\frac{\partial L(\mathbf{w}, \mathbf{b}, \alpha)}{\partial \mathbf{b}} = \sum_{i=1}^N y_i \alpha_i = 0 \quad (15)$$

For $i=1,2,\dots, N$,

$$\frac{\partial L(\mathbf{w}, \mathbf{b}, \alpha)}{\partial \xi_i} = C - \alpha_i - \mu_i = 0 \quad (16)$$

$$\alpha_i \geq 0 \quad (17)$$

$$\mu_i \geq 0 \quad (18)$$

KKT Complimentary Conditions

The KKT complimentary conditions are,

$$\alpha_i [y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + \mathbf{b}) - 1 + \xi_i] = 0 \quad (19)$$

$$\mu_i \xi_i = 0 \quad (20)$$

From (16) and(18)

$$C - \alpha_j = \mu_j \geq 0 \quad (21)$$

$$0 \leq \alpha_j \leq C, i = 1, 2, \dots N \quad (22)$$

Three cases

- $\alpha_j = 0$
- $0 < \alpha_j < C$
- $\alpha_j = C$

- When $\alpha_j = 0$, from (21), $\mu_j = C > 0$, therefore from (20), $\xi_j = 0$. Hence from (19), $y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1$.
- If $0 < \alpha_j < C$, from (21), $\mu_j > 0$. Hence from (20), $\xi_j = 0$. Therefore from (19), $y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) = 1$. Hence, points for which $0 < \alpha_j < C$ lie at the target distance of $\frac{1}{\|\mathbf{w}\|}$ from the hyperplane.
- If $\alpha_j = C$, from (21), $\mu_j = 0$. Hence from (20), $\xi_j \geq 0$. Hence (19), $y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) = 1 - \xi_j$. That is, $y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \leq 1$.

Support Vectors

Support vectors are those points for which $y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \leq 1$.

- $y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) = 1$, \mathbf{x}_i on H_1 or on H_2
- $0 \leq y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) < 1$ Between H_1 and H or H_2 and H or on H
- $y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) < 0$ Incorrect classification

Dual Formulation

- The terms that involve w : $\frac{1}{2}\|w\|^2 - \sum_i \alpha_i [y_i (w^T x_i)]$
- The terms that involve b : $\sum_i \alpha_i y_i b$
- The terms that involve ξ_i : $\sum_{i=1}^N \xi_i (C - \alpha_i - \mu_i)$
- The remaining terms are: $\sum_i \alpha_i$

Dual form is

$$\text{maximise } W(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N y_i y_j \alpha_i \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

$$\text{subject to } \sum_{i=1}^N y_i \alpha_i = 0$$

$$C \geq \alpha_i \geq 0, i = 1, 2, \dots, N.$$

General Form

$$V(y_i, \tilde{f}(x_i)) = \max(0, 1 - y_i(\tilde{f}(x_i))) = \xi_i$$

$$\tilde{f}(x) = f(x) + b$$

$$\min_{f \in \mathcal{F}, b \in \mathbb{R}} \frac{1}{2} \|f\|^2 + C \sum_{i=1}^N \xi_i \tag{23}$$
$$y_i(f(x_i) + b) \geq 1 - \xi_i, \quad i = 1, \dots, N$$
$$\xi_i \geq 0.$$

Hyperplane Models: RKHS

If \tilde{f} is the unknown function of a datamodeling problem, then by kernel theory

$$\tilde{f}(x) = f(x) + b = \langle f, k_x \rangle + b, \forall x \in \mathcal{X}; f, k_x \in \mathcal{F}, b \in \mathbb{R} \quad (24)$$

Consider $\tilde{H}_{f,b} : \mathcal{F} \rightarrow \mathbb{R}$ where

$$\tilde{H}_{f,b}(g) = \langle f, g \rangle + b, \forall g \in \mathcal{F} \quad (25)$$

(25) is the equation of a hyperplane in RKHS.
Comparing (24) and (25),

$$\tilde{f}(x) = \tilde{H}_{f,b}(k_x), x \in \mathcal{X}, k_x \in \mathcal{F}$$

Thus finding \tilde{f} in input space is equivalent in finding $\tilde{H}_{f,b}$ in RKHS.

Optimisation Problem

$$\begin{aligned} \min_{f \in \mathcal{F}, b \in \mathbb{R}} \quad & \frac{1}{2} \|f\|^2 + C \sum_{i=1}^N \xi_i \\ \text{subject to} \quad & y_i(\langle f, k_{x_i} \rangle + b) - 1 + \xi_i \geq 0, \quad i = 1, \dots, N \\ & \xi_i \geq 0. \end{aligned} \tag{26}$$

- The objective function is quadratic and all the constraints affine.

Lagrangian Formulation

The primal Lagrangian is,

$$L(f, b, \xi, \alpha) = \frac{\|f\|^2}{2} + C \sum_{i=1}^N \xi_i - \sum_i \alpha_i [y_i (\langle f, k_{x_i} \rangle + b) - 1 + \xi_i] - \sum_i \mu_i \xi_i \quad (27)$$

KKT conditions

$$\frac{\partial L}{\partial f} = f - \sum_{i=1}^N \alpha_i y_i k_{x_i} = 0$$

$$f = \sum_{i=1}^N \alpha_i y_i k_{x_i} \quad (28)$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^N \alpha_i y_i = 0 \quad (29)$$

For $i=1,2,\dots, N$,

$$\frac{\partial L(\mathbf{w}, \mathbf{b}, \alpha)}{\partial \xi_i} = C - \alpha_i - \mu_i = 0 \quad (30)$$

$$\alpha_i \geq 0 \quad (31)$$

$$\mu_i \geq 0 \quad (32)$$

KKT Complimentary Conditions

The KKT complimentary conditions are,

$$\alpha_i [y_i (\langle f, k_{x_i} \rangle + b) - 1 + \xi_i] = 0 \quad (33)$$

$$\mu_i \xi_i = 0 \quad (34)$$

From (30) and(32)

$$C - \alpha_j = \mu_j \geq 0 \quad (35)$$

$$0 \leq \alpha_j \leq C, i = 1, 2, \dots N \quad (36)$$

Three cases

- $\alpha_j = 0$
- $0 < \alpha_j < C$
- $\alpha_j = C$

- When $\alpha_j = 0$, from (35), $\mu_j = C > 0$, therefore from (34), $\xi_j = 0$. Hence from (33), $y_i(\langle f, k_{x_i} \rangle + b) \geq 1$.
- If $0 < \alpha_j < C$, from (35), $\mu_j > 0$. Hence from (34), $\xi_j = 0$. Therefore from (33), $y_i(\langle f, k_{x_i} \rangle + b) = 1$. Hence, points for which $0 < \alpha_j < C$ lie at the target distance of $\frac{1}{\|f\|}$ from the hyperplane.
- If $\alpha_j = C$, from (35), $\mu_j = 0$. Hence from (34), $\xi_j \geq 0$. Hence (33), $y_i(\langle f, k_{x_i} \rangle + b) = 1 - \xi_j$. That is, $y_i(\langle f, k_{x_i} \rangle + b) \leq 1$.

Support Vectors

Support vectors are those points for which $y_i(\langle f, k_{x_i} \rangle + b) \leq 1$.

- $y_i(\langle f, k_{x_i} \rangle + b) = 1$, k_{x_i} on H_1 or on H_2
- $0 \leq y_i(\langle f, k_{x_i} \rangle + b) < 1$ Between H_1 and H or H_2 and H or on H
- $y_i(\langle f, k_{x_i} \rangle + b) < 0$ Incorrect classification

Dual Formulation

- The terms that involve f : $\frac{1}{2}\|f\|^2 - \sum_j \alpha_j [y_j \langle f, k_{x_j} \rangle]$
- The terms that involve b : $\sum_j \alpha_j y_j b$
- The terms that involve ξ_j : $\sum_{i=1}^N \xi_i (C - \alpha_i - \mu_i)$
- The remaining terms are: $\sum_j \alpha_j$

$$\|f\|^2 = \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j k(x_i, x_j)$$

$$\|f\|^2 = \beta^T K \beta$$

where K is the kernel matrix and $\beta = (\alpha_1 y_1, \alpha_2 y_2, \dots, \alpha_N y_N)^T$

- $\|f\|^2 \geq 0$
 - k should be positive semi-definite
 - K should be positive semi-definite

Dual Optimisation

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N y_i y_j \alpha_i \alpha_j k(x_i, x_j)$$

$$\text{subject to } \sum_{i=1}^N y_i \alpha_i = 0$$

$$C \geq \alpha_i \geq 0, i = 1, 2, \dots, N.$$

SVM Model

$$f = \sum_{i=1}^{sv} y_i \alpha_i^* k_{x_i} \quad (37)$$

where sv is the number of support vectors. Then

$$\begin{aligned} \tilde{f}(x) &= \langle f, k_x \rangle + b \\ &= \sum_{i=1}^{sv} y_i \alpha_i^* \langle k_{x_i}, k_x \rangle + b \\ &= \sum_{i=1}^{sv} y_i \alpha_i^* k(x_i, x) + b \end{aligned}$$

Characteristics of SVM Model

- Sparse algorithm
- In SVM classification, the optimal hyperplane is found in RKHS by minimizing the training error and maximizing the margin.

RKHS and Reproducing Kernel

- There exists a bijection between the set of all reproducing kernel Hilbert spaces and the set of all positive kernel functions.
- Regularization attempts to provide well-posed solutions to a learning task, specifically ERM, by constraining the capacity of the hypothesis space through the elimination of complex functions that are unlikely to generalize, thereby isolating a unique and stable solution.

Hyperparameters

Training a model requires the choice of few relevant quantities:

- The kernel function, that determines the shape of the decision surface;
- Parameter in the kernel function (eg: for Gaussian kernel: variance of the Gaussian, for polynomial kernel: degree of the polynomial)
- The regularization parameter

Non Parametric Algorithm

In kernel methods, the number of basis functions depends on the number of data and hence should not be fixed a priori. Hence they are non parametric methods.