

Lecture 12: Support Vector Machine

Shukai Gong

Given a dataset $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, where $\mathbf{x}_n \in \mathbb{R}^D$ and $y_n \in \{-1, 1\}$, the goal of SVM is to find a hyperplane $\mathbf{w}^\top \mathbf{x} + b = 0$ with largest separation such that the discriminative power is maximized and the error risk is minimized. As is shown in Figure 1, **the desired hyperplane has the largest distance to the nearest training data points of any class** so that the hyperplane is less likely to be influenced by out-of-sample data points and more robust.

The hyperplane is shifted along two opposite directions $\mathbf{w}^\top \mathbf{x} = b + d$ and $\mathbf{w}^\top \mathbf{x} = b - d$ to form a margin of width $\frac{2d}{\|\mathbf{w}\|}$. The data points that are on the margin are called **support vectors**. Eventually, the SVM model can be written by

$$f^*(\mathbf{x}) = \text{sign}((\mathbf{w}^*)^\top \mathbf{x} + b^*)$$

where \mathbf{w}^* and b^* are optimal parameters.

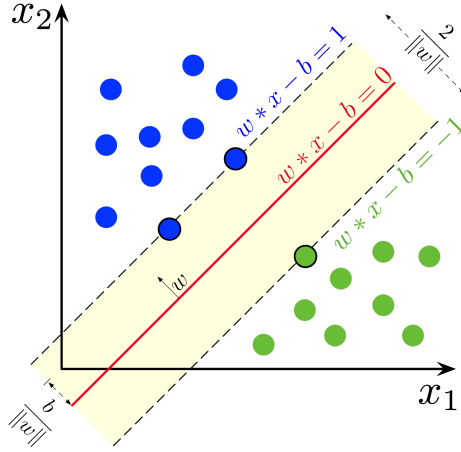


Figure 1: Hard-Margin Support Vector Machine

1 Hard-Margin SVM

1.1 Formulation of Optimization Goal

For a set of training data points $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $y_i \in \{-1, +1\}$ that are **2-class linear separable**, the hard-margin SVM problem can be formulated as

$$\begin{aligned} & \max \text{margin}(\mathbf{w}, b), \text{ s.t. } \begin{cases} \mathbf{w}^\top \mathbf{x}_i + b \geq 0, & \text{if } y_i = 1, \\ \mathbf{w}^\top \mathbf{x}_i + b < 0, & \text{if } y_i = -1. \end{cases}, i = 1, \dots, N \\ \iff & \max \text{margin}(\mathbf{w}, b), \text{ s.t. } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 0, i = 1, \dots, N \end{aligned}$$

Margin can be defined as the shortest distance from the hyperplane to the nearest data point of any class, i.e.

$$\text{margin}(\mathbf{w}, b) = \min_{i=1, \dots, N} \frac{|\mathbf{w}^\top \mathbf{x}_i + b|}{\|\mathbf{w}\|}$$

Therefore our optimization goal can be rewritten as

$$\begin{aligned}
& \max_{\mathbf{w}, b} \min_{i=1, \dots, N} \frac{|\mathbf{w}^\top \mathbf{x}_i + b|}{\|\mathbf{w}\|}, & \text{s.t. } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 0, i = 1, \dots, N \\
\iff & \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} \min_{i=1, \dots, N} y_i(\mathbf{w}^\top \mathbf{x}_i + b), & \text{s.t. } \exists r > 0, \min_{i=1, \dots, N} y_i(\mathbf{w}^\top \mathbf{x}_i + b) = r, i = 1, \dots, N \\
\iff & \max_{\mathbf{w}, b} \frac{r}{\|\mathbf{w}\|}, & \text{s.t. } \exists r > 0, \min_{i=1, \dots, N} y_i(\mathbf{w}^\top \mathbf{x}_i + b) = r, i = 1, \dots, N
\end{aligned}$$

Without loss of generality, we can set $r = 1$ since it is just a scaling factor and can be absorbed into \mathbf{w} and b . Therefore, the optimization problem can be rewritten as:

$$\begin{aligned}
& \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|}, & \text{s.t. } \min_{i=1, \dots, N} y_i(\mathbf{w}^\top \mathbf{x}_i + b) = 1, i = 1, \dots, N \\
\iff & \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2, & \text{s.t. } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, i = 1, \dots, N
\end{aligned}$$

1.2 Dual Problem Optimization

The primal problem can be solved by the Lagrange multiplier method, and the Lagrangian function is a convex optimization problem with N linear constraints,

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2, \text{ s.t. } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, i = 1, \dots, N$$

It's **equivalent** to an unconstrained primal problem with N Lagrange multipliers $\lambda_i \geq 0$,

$$\begin{aligned}
\mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) &= \min_{\mathbf{w}, b} \max_{\boldsymbol{\lambda}} \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \underbrace{\lambda_i}_{\geq 0} \underbrace{(1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b))}_{\leq 0} \\
&\Rightarrow \begin{cases} \min_{\mathbf{w}, b} \max_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) \\ \text{s.t. } \lambda_i \geq 0, i = 1, \dots, N \end{cases}
\end{aligned}$$

Non-negative Lagrangian Multipliers λ_i

For arbitrary data point (\mathbf{x}_i, y_i) ,

- If $1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b) > 0$, then $\max_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) = \frac{1}{2} \|\mathbf{w}\|^2 + \infty = \infty$. (Meaningless!)
- If $1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b) \leq 0$, then $\max_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) = \frac{1}{2} \|\mathbf{w}\|^2$. (Set all $\lambda_i = 0$ maximizes the non-positive term $\lambda_i(1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b))$)

So

$$\min_{\mathbf{w}, b} \max_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) = \min_{\mathbf{w}, b} (\infty, \frac{1}{2} \|\mathbf{w}\|^2) = \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

By setting the Lagrangian multipliers $\lambda_i \geq 0$ for all $i = 1, \dots, N$, we made an equivalent optimization problem to the primal problem.

It can be proved that the dual problem is **equivalent** to the primal problem when the optimization goal is convex and the constraints are linear (**strong duality**), i.e.

$$\begin{cases} \min_{\mathbf{w}, b} \max_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) \\ \text{s.t. } \lambda_i \geq 0, i = 1, \dots, N \end{cases} = \begin{cases} \max_{\boldsymbol{\lambda}} \min_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) \\ \text{s.t. } \lambda_i \geq 0, i = 1, \dots, N \end{cases}$$

First, we take the derivative of $\mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda})$ w.r.t \mathbf{w} and b and set them to zero,

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i = 0 \\ \frac{\partial \mathcal{L}}{\partial b} = - \sum_{i=1}^N \lambda_i y_i = 0 \end{cases} \Rightarrow \begin{cases} \mathbf{w} = \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i \\ \sum_{i=1}^N \lambda_i y_i = 0 \end{cases}$$

Plugging these back to $\mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda})$, we have

$$\begin{aligned} \mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) &= \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \lambda_i (1 - y_i (\mathbf{w}^\top \mathbf{x}_i + b)) \\ &= \frac{1}{2} \mathbf{w}^\top \mathbf{w} + \sum_{i=1}^N \lambda_i - \sum_{i=1}^N \lambda_i y_i (\mathbf{w}^\top \mathbf{x}_i) + b \sum_{i=1}^N \lambda_i y_i \\ &= \frac{1}{2} \left(\sum_{i=1}^N \lambda_i y_i \mathbf{x}_i \right)^\top \left(\sum_{i=1}^N \lambda_i y_i \mathbf{x}_i \right) + \sum_{i=1}^N \lambda_i - \sum_{i=1}^N \lambda_i y_i \left(\sum_{j=1}^N \lambda_j y_j \mathbf{x}_j \right)^\top \mathbf{x}_i \\ &= \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j + \sum_{i=1}^N \lambda_i - \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \mathbf{x}_j^\top \mathbf{x}_i \\ &= \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \end{aligned}$$

So the duality problem is simplified to

$$\begin{cases} \max_{\boldsymbol{\lambda}} & \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \\ \text{s.t.} & \lambda_i \geq 0, i = 1, \dots, N \\ & \sum_{i=1}^N \lambda_i y_i = 0 \end{cases} = \begin{cases} \min_{\boldsymbol{\lambda}} & \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j - \sum_{i=1}^N \lambda_i \\ \text{s.t.} & \lambda_i \geq 0, i = 1, \dots, N \\ & \sum_{i=1}^N \lambda_i y_i = 0 \end{cases}$$

A sufficient and necessary condition for **strong duality** (we use without proof) is that the primal problem satisfies the **Karush-Kuhn-Tucker (KKT) conditions**. The KKT conditions for this specific problem are

$$\begin{cases} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) = 0 & \Rightarrow \mathbf{w}^* = \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i \\ \nabla_b \mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) = 0 & \Rightarrow \sum_{i=1}^N \lambda_i y_i = 0 \\ \lambda_i (1 - y_i (\mathbf{w}^\top \mathbf{x}_i + b)) = 0, & i = 1, \dots, N \\ \lambda_i \geq 0, & i = 1, \dots, N \\ 1 - y_i (\mathbf{w}^\top \mathbf{x}_i + b) \leq 0, & i = 1, \dots, N \end{cases}$$

Let's focus on the third condition (**complementary slackness condition**). For arbitrary data point (\mathbf{x}_i, y_i) ,

- If the data point is a support vector, i.e. $1 - y_i (\mathbf{w}^\top \mathbf{x}_i + b) = 0$, then $\lambda_i > 0$.
- If the data point is beyond the margin, i.e. $1 - y_i (\mathbf{w}^\top \mathbf{x}_i + b) < 0$, then $\lambda_i = 0$.

which means that non-support-vector data points won't show up in $\mathbf{w}^* = \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i$ since their $\lambda_i = 0$. Taking advantage of this, by picking one support vector (\mathbf{x}_k, y_k) , we can derive b by

$$1 - y_k (\mathbf{w}^\top \mathbf{x}_k + b) = 0 \Rightarrow y_k^2 (\mathbf{w}^\top \mathbf{x}_k + b) = y_k \Rightarrow b^* = y_k - \mathbf{w}^\top \mathbf{x}_k$$

The hard-margin SVM is therefore

$$f^*(\mathbf{x}) = \text{sign}((\mathbf{w}^*)^\top \mathbf{x} + b^*), \quad \begin{cases} \mathbf{w}^* = \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i \\ b^* = y_k - \mathbf{w}^\top \mathbf{x}_k \end{cases}$$

which only **depends on the support vectors** and the corresponding Lagrange multipliers.

2 Soft-Margin SVM

Hard-margin SVM is too idealized about the separability of the data points. In practice, the data points are usually either not linearly separable, or subject to noise. As is shown in Figure 2, Soft-margin SVM **allows for some data points to be misclassified** by introducing a slack variable $\xi_i \geq 0$ for each data point (\mathbf{x}_i, y_i) to measure the degree of misclassification.

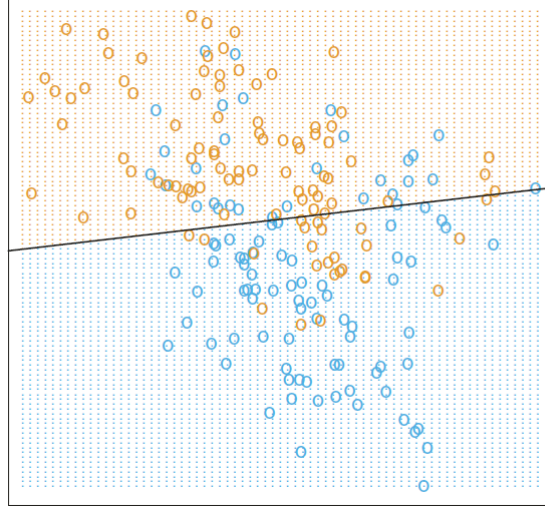


Figure 2: Soft-Margin Support Vector Machine

Recall that $y_i(\mathbf{w}^\top \mathbf{x} + b) \geq 1$ indicates that the data point is correctly classified. The misclassification is presented by $y_i(\mathbf{w}^\top \mathbf{x} + b) < 1$. Some methods of defining slackness ξ_i is presented below:

- **Counts:** $\xi_i = \mathbf{I}(y_i(\mathbf{w}^\top \mathbf{x}_i + b) < 1)$ (**Discontinuous!**)
- **Hinge Loss:** $\xi_i = \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b))$ (**Recommended**)

We can formulate the optimization problem as

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)), & \text{s.t. } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i, i = 1, \dots, N \\ \iff \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i, & \text{s.t. } \begin{cases} y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i, & i = 1, \dots, N \\ \xi_i \geq 0, & i = 1, \dots, N \end{cases} \end{aligned}$$

Similar to the hard-margin SVM, we construct the Lagrangian function and derive its dual problem.

$$\begin{aligned}
& \min_{\mathbf{w}, b, \boldsymbol{\xi}} \max_{\boldsymbol{\lambda}, \boldsymbol{\mu}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i + \sum_{i=1}^N \lambda_i (1 - \xi_i - y_i(\mathbf{w}^\top \mathbf{x}_i + b)) - \sum_{i=1}^N \mu_i \xi_i \\
& = \max_{\boldsymbol{\lambda}, \boldsymbol{\mu}} \min_{\mathbf{w}, b, \boldsymbol{\xi}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i + \sum_{i=1}^N \lambda_i (1 - \xi_i - y_i(\mathbf{w}^\top \mathbf{x}_i + b)) - \sum_{i=1}^N \mu_i \xi_i \\
& \Rightarrow \begin{cases} \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i = 0 \\ \frac{\partial \mathcal{L}}{\partial b} = - \sum_{i=1}^N \lambda_i y_i = 0 \\ \frac{\partial \mathcal{L}}{\partial \xi_i} = C - \lambda_i - \mu_i = 0 \end{cases} \Rightarrow \begin{cases} \mathbf{w} = \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i \\ \sum_{i=1}^N \lambda_i y_i = 0 \\ \lambda_i + \mu_i = C \end{cases}
\end{aligned}$$

and the dual problem of the primal soft-margin SVM is

$$\begin{aligned}
\begin{cases} \max_{\boldsymbol{\lambda}} & \mathcal{L}(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \\ \text{s.t.} & \lambda_i \geq 0, i = 1, \dots, N \\ & \mu_i \geq 0, i = 1, \dots, N \end{cases} = \begin{cases} \min_{\boldsymbol{\lambda}} & \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j - \sum_{i=1}^N \lambda_i \\ \text{s.t.} & \lambda_i \geq 0, \mu_i \geq 0, i = 1, \dots, N \\ & \sum_{i=1}^N \lambda_i y_i = 0 \\ & \lambda_i + \mu_i = C, i = 1, \dots, N \end{cases} \\
& = \begin{cases} \min_{\boldsymbol{\lambda}} & \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j - \sum_{i=1}^N \lambda_i \\ \text{s.t.} & 0 \leq \lambda_i \leq C, i = 1, \dots, N \\ & \sum_{i=1}^N \lambda_i y_i = 0 \end{cases}
\end{aligned}$$

Compare to the hard-margin SVM case, the soft-margin SVM has an additional constraint $\lambda_i \leq C$ and $\mu_i \geq 0$ to control the trade-off between the margin and the misclassification. The KKT condition for the primal soft-margin SVM is

$$\begin{cases} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = 0 & \Rightarrow \mathbf{w}^* = \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i \\ \nabla_b \mathcal{L}(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = 0 & \Rightarrow \sum_{i=1}^N \lambda_i y_i = 0 \\ \nabla_{\xi_i} \mathcal{L}(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = 0 & \Rightarrow \lambda_i + \mu_i = C, i = 1, \dots, N \\ \lambda_i (1 - \xi_i - y_i(\mathbf{w}^\top \mathbf{x}_i + b)) = 0, & i = 1, \dots, N \\ \mu_i \xi_i = 0, & i = 1, \dots, N \\ \lambda_i \geq 0, \mu_i \geq 0, \xi_i \geq 0, & i = 1, \dots, N \\ 1 - \xi_i - y_i(\mathbf{w}^\top \mathbf{x}_i + b) \leq 0, & i = 1, \dots, N \end{cases}$$

Let's break down this huge KKT conditions: For arbitrary data point (\mathbf{x}_i, y_i) , either $\lambda_i = 0$ or $1 - \xi_i - y_i(\mathbf{w}^\top \mathbf{x}_i + b) = 0$.

$$\begin{cases} \text{If } \lambda_i = 0 \Rightarrow \text{No effect on } \mathbf{w}^* \\ \text{If } \lambda_i > 0 \Rightarrow (\mathbf{x}_i, y_i) \text{ is a SV} \end{cases} \begin{cases} \text{If } \lambda_i < C \Rightarrow \mu_i > 0 \Rightarrow \xi_i = 0 \Rightarrow b^* = y_i - \mathbf{w}^\top \mathbf{x}_i \\ \text{If } \lambda_i = C \Rightarrow \mu_i = 0 \begin{cases} \text{If } \xi_i \leq 1 \Rightarrow \text{Tolerated Misclassification} \\ \text{If } \xi_i > 1 \Rightarrow \text{Misclassification} \end{cases} \end{cases}$$

So we need to find a support vector (\mathbf{x}_k, y_k) with $0 < \lambda_k < C$ to derive b^* . Still, the soft-margin SVM is

$$f^*(\mathbf{x}) = \text{sign}((\mathbf{w}^*)^\top \mathbf{x} + b^*), \quad \begin{cases} \mathbf{w}^* = \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i \\ b^* = y_k - \mathbf{w}^\top \mathbf{x}_k \end{cases}$$

3 Kernel SVM

We have assumed in hard and soft SVM that data points are linearly separable. For non-linearly separable data points, we can use the **kernel trick to map the data points to a higher-dimensional feature space where they are linearly separable.**

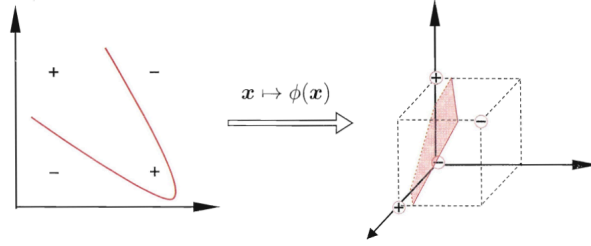


Figure 3: Kernel Trick in SVM

Denote $\phi(\mathbf{x})$ as the feature vector of \mathbf{x} in the higher-dimensional space, the SVM model corresponding to the hyperplane in the feature space is $\mathbf{w}^\top \phi(\mathbf{x}) + b = 0$ and the optimization problem can be formulated as (take soft-margin SVM as an example)

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i, \quad \text{s.t.} \quad \begin{cases} y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b) \geq 1 - \xi_i, & i = 1, \dots, N \\ \xi_i \geq 0, & i = 1, \dots, N \end{cases}$$

The dual problem of the primal kernel soft-margin SVM is

$$\begin{cases} \min_{\lambda} & \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) - \sum_{i=1}^N \lambda_i \\ \text{s.t.} & 0 \leq \lambda_i \leq C, i = 1, \dots, N \\ & \sum_{i=1}^N \lambda_i y_i = 0 \end{cases}$$

We can observe that the only difference between the kernel SVM and the linear SVM is that the dot product $\mathbf{x}_i^\top \mathbf{x}_j$ is replaced by the kernel function $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$. As an analogy, $\mathbf{w}^* = \sum_{i=1}^N \lambda_i y_i \phi(\mathbf{x}_i)$ and $b^* = y_k - \mathbf{w}^\top \phi(\mathbf{x}_k)$, and the kernel SVM model is written by

$$\begin{aligned} f^*(\mathbf{x}) &= \text{sign}((\mathbf{w}^*)^\top \phi(\mathbf{x}) + b^*) \\ &= \text{sign} \left(\sum_{i=1}^N \lambda_i y_i \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}) + b^* \right) \\ &= \text{sign} \left(\sum_{i=1}^N \lambda_i y_i \kappa(\mathbf{x}_i, \mathbf{x}) + b^* \right) \end{aligned}$$

Similar to other kernel methods, we don't have to find ϕ explicitly, but only need to introduce the kernel function κ to derive the hyperplane.

4 Noticable Facts about SVM

- Both primal and dual problems are quadratic programming problems, which can be solved by numerical methods. Ultimately we need to know certain λ_k to derive the hyperplane.
- When $C \rightarrow \infty$, the soft-margin SVM "behaves like" the hard-margin SVM. (From $0 < \lambda_k < C$ to $0 < \lambda_k$)
- Linear SVM is a special case of kernel SVM with a linear kernel: $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$.

References

- [Machine Learning, Zhi-Hua Zhou](#)
- [Support Vector Machine](#)