

The background of the slide features a large, light gray watermark of the ShanghaiTech University logo. The logo is circular, with the university's name in Chinese characters '上海科技大学' at the top and 'SHANGHAITECH UNIVERSITY' at the bottom. In the center is a stylized graphic of a DNA double helix and a building. The year '2013' is also visible.

School of Information Science and Technology  
ShanghaiTech University

# Kernels Methods and Support Vector Machines

Shuhao Xia

2013 May 4, 2020



## Kernels Methods

- Understanding Kernels
- Constructing Kernels
- Other forms of Kernel Functions

## Support Vector Machines (SVM)

- Maximizing Margin Classifiers
- SVM Optimization
- Soft Margin

## *Memory-based method*

- ▶ Keep and use training samples,
- ▶ Measure the similarity of training and test samples,
- ▶ Fast to 'train' but slow at making predictions for test samples

E.g., Nearest neighbors.



Now we need a metric to measure such a similarity. Typically, we use inner product, and *kernels function* is therefore defined as

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}'),$$

where  $\phi(\cdot)$  is a fixed nonlinear *feature space* mapping.

## Mercer condition

$k(\cdot)$  is a valid kernel if and only if the Gram matrix  $\mathbf{K} \triangleq [k(\mathbf{x}_n, \mathbf{x}_m)]$  satisfies the following conditions,

- ▶  $\mathbf{K}$  is symmetric,
- ▶  $\mathbf{K}$  is positive semi-definite, e.g.,  $\mathbf{a}^\top \mathbf{K} \mathbf{a} \geq 0 \forall \mathbf{a}$ .

Mercer condition is often used to test whether a function is a valid kernel without explicit  $\phi(\mathbf{X})$ .

- Conceptually, to apply kernel function, we
  1. find  $\phi(\cdot)$  and calculate  $\phi(\mathbf{x})$ ,
  2. measure the similarity by some kernel  $k(\mathbf{x}, \mathbf{x}')$ .

However,  $\mathbf{x}$  is typically mapped to a high-dimensional space by  $\phi(\cdot)$  so that the computational complexity for inner product is huge.

- In practice, we can therefore work directly in terms of kernels, which allows us implicitly to use feature spaces of high, even infinite dimensionality.

Given valid kernels  $k_1(\mathbf{x}, \mathbf{x}')$ ,  $k_2(\mathbf{x}, \mathbf{x}')$ , the following new kernels will also valid:

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \quad (1)$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \quad (2)$$

$$k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}')) \quad (3)$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}')) \quad (4)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}') \quad (5)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}') \quad (6)$$

$$k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}')) \quad (7)$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{A} \mathbf{x}' \quad (8)$$

where  $c > 0$  is a constant,  $f(\cdot)$  is any function,  $q(\cdot)$  is a polynomial with nonnegative coefficients,  $k_3(\cdot)$  is a valid kernel and  $\mathbf{A}$  is a symmetric positive semidefinite matrix.

# Example: Gaussian Kernel



Gaussian Kernel is defined as follows,

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\|\mathbf{x} - \mathbf{x}'\|^2 / 2\sigma^2\right).$$

## Proof.

Expanding the square term,

$$\|\mathbf{x} - \mathbf{x}'\|^2 = \mathbf{x}^T \mathbf{x} + (\mathbf{x}')^T \mathbf{x}' - 2\mathbf{x}^T \mathbf{x}',$$

to give

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\mathbf{x}^T \mathbf{x} / 2\sigma^2) \exp(\mathbf{x}^T \mathbf{x}' / \sigma^2) \exp(-(\mathbf{x}')^T \mathbf{x}' / 2\sigma^2),$$

and then making use of (2) and (4), together with the validity of the linear kernel  $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$ . □





- ▶ *stationary* kernels  $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x} - \mathbf{x}')$ , invariant to translations in input space,
- ▶ *homogeneous* kernels  $k(\mathbf{x}, \mathbf{x}') = k(\|\mathbf{x} - \mathbf{x}'\|)$ , depend only on the magnitude of the distance.



Consider two-class classification problem using linear model of the form

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b. \quad (9)$$

If the training data set is linearly separable,

- ▶  $\geq 1$  *hyperplanes* to correctly classify all the training samples,
- ▶ PLA can efficiently find one of them dependent on initialization.

However, how to find the best one of these hyperplanes?

## Definition

The smallest distance between the decision boundary and any of the samples.

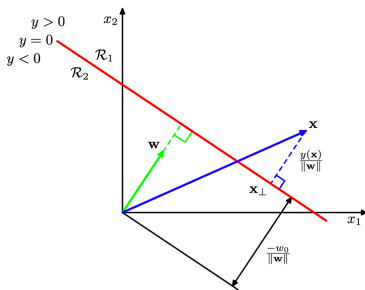


Figure: Linear Model

For all correct solutions, we have  $t_n y(x_n) > 0 \forall n$ . As shown in the figure, the distance of a point  $\mathbf{x}_n$  to the decision surface is given by

$$\frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|} = \frac{t_n (\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|}.$$

The margin is given by the smallest distance, and then we want to maximize the margin. The whole optimization problem can be written as

$$\arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n (\mathbf{w}^T \mathbf{x}_n + b)] \right\}.$$

## Setting

$$t_n (\mathbf{w}^T \mathbf{x}_n + b) = 1$$

for the point that is closest to the surface. Therefore, all samples will satisfy the constraints

$$t_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1, \quad n = 1, \dots, N.$$

The original problem can be converted to the following optimization problem

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s. t.} \quad & t_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1, \quad n = 1, \dots, N. \end{aligned}$$

Introducing Lagrange multipliers  $a_n \geq 0$ , the Lagrangian function is given by

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N a_n \{t_n (\mathbf{w}^T \mathbf{x}_n + b) - 1\}.$$

Hence, we want to optimize the Lagrangian function as follows

$$\begin{aligned} \min_{\mathbf{w}, b} \max_{\mathbf{a}} L(\mathbf{w}, b, \mathbf{a}) \\ \text{s. t. } a_n \geq 0 \end{aligned}$$

Using strong duality, the problem becomes

$$\begin{aligned} \max_{\mathbf{a}} \min_{\mathbf{w}, b} L(\mathbf{w}, b, \mathbf{a}) \\ \text{s. t. } a_n \geq 0 \end{aligned}$$

# Dual Representation



Setting the derivatives of  $L(\mathbf{w}, b, \mathbf{a})$  with respect to  $\mathbf{w}, b$  to zero, we obtain the following two conditions

$$\begin{aligned}\mathbf{w} &= \sum_{n=1}^N a_n t_n \mathbf{x}_n \\ 0 &= \sum_{n=1}^N a_n t_n\end{aligned}$$

Eliminating  $\mathbf{w}$  and  $b$  from  $L(\mathbf{w}, b, \mathbf{a})$ , we arrive the dual form

$$\begin{aligned}\max_{\mathbf{a}} \tilde{L}(\mathbf{a}) &= \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m) \\ \text{s. t. } a_n &\geq 0, \quad n = 1, \dots, N \\ \sum_{n=1}^N a_n t_n &= 0\end{aligned}$$

The K.K.T condition can be derived as follows

$$a_n \geq 0$$

$$t_n y(\mathbf{x}_n) - 1 \geq 0$$

$$a_n \{t_n y(\mathbf{x}_n) - 1\} = 0$$

which can be interpreted as

- ▶  $a_n = 0, t_n y(\mathbf{x}_n) - 1 \geq 0$ ,  $n$ -th sample plays no role for predicting,
- ▶  $a_n \geq 0, t_n y(\mathbf{x}_n) - 1 = 0$ ,  $n$ -th sample is called support vectors.



In practice, the training samples are almost not linear separable. Hence, we may not want exact separation, which can also lead to poor generalization.

To tolerate "wrong side", we introduce *slack variables*  $\xi_n \geq 0$ , so the SVM optimization problem can be written as

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n \\ \text{s. t.} \quad & 1 - t_n y(\mathbf{x}_n) - \xi_n \leq 0, \forall n \\ & \xi_n \geq 0 \end{aligned}$$

►  $C \rightarrow \infty \implies \xi_n \rightarrow 0$

►  $C < \infty \implies \xi_n > 0$

The corresponding Lagrangian is given by

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N a_n \{t_n y(\mathbf{x}_n) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n,$$

and K.K.T condition is given by

$$\begin{aligned} a_n &\geq 0 \\ t_n y(\mathbf{x}_n) - 1 + \xi_n &\geq 0 \\ a_n (t_n y(\mathbf{x}_n) - 1 + \xi_n) &= 0 \\ \mu_n &\geq 0 \\ \xi_n &\geq 0 \\ \mu_n \xi_n &= 0 \end{aligned}$$

The optimization problem is

$$\begin{aligned} \max_{\mathbf{w}, b, \xi} \min_{\mathbf{a}, \mu} L(\mathbf{w}, b, \mathbf{a}, \mu) \\ \text{s. t. } a_n \geq 0, \mu_n \geq 0, \forall n. \end{aligned}$$

The dual form then is written as

$$\begin{aligned} \min_{\mathbf{a}, \mu} \max_{\mathbf{w}, b, \xi} L(\mathbf{w}, b, \mathbf{a}, \mu) \\ \text{s. t. } a_n \geq 0, \mu_n \geq 0, \forall n. \end{aligned}$$

After optimizing  $\mathbf{w}, b, \xi_n$ , we arrive at

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{n=1}^N a_n t_n \mathbf{x}_n$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{n=1}^N a_n t_n = 0$$

$$\frac{\partial L}{\partial \xi_n} = 0 \Rightarrow a_n = C - \mu_n$$

Dual Lagrangian

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m).$$

The optimization problem

$$\begin{aligned} \max_{\mathbf{a}} \quad & \tilde{L}(\mathbf{a}) \\ \text{s. t.} \quad & a_n \geq 0, \quad n = 1, \dots, N \\ & \sum_{n=1}^N a_n t_n = 0 \\ & C - a_n - \mu_n = 0. \end{aligned}$$

## Pros

- ▶ Rigorous mathematical theory, strong interpretability
- ▶ Able to solve the high-dimensional problem with kernels methods
- ▶ Computational complexity depends on support vectors

## Cons

- ▶ Long training time, e.g. SMO algorithm
- ▶ Large space if using kernels methods
- ▶ Predicting time depends on the number of support machines