

Soft Margin SVM

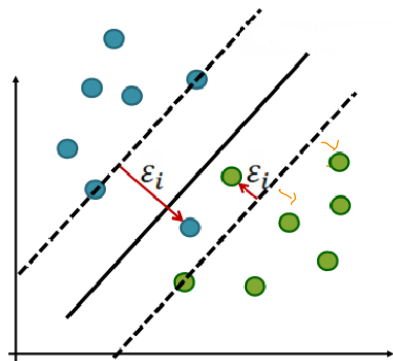
Functional margin $1 - \xi_i \leq 1$:

$$\min_{w, b, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$$

← slack variable.

$$\text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1 - \xi_i$$
$$\xi_i \geq 0, i = 1, \dots, m$$

- ▶ C: relative weight on the regularizer
- ▶ L_1 regularization let most $\xi_i = 0$, such that their functional margins $1 - \xi_i = 1$



Soft Margin SVM

The generalized Lagrangian function:

$$L(w, b, \xi, \alpha, r) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i - \sum_i^m \alpha_i \left[y^{(i)} (w^T x^{(i)} + b) - 1 + \xi_i \right] - \sum_{i=1}^m r_i \xi_i$$

Soft Margin SVM

The generalized Lagrangian function:

$$L(w, b, \xi, \alpha, r) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i - \sum_i^m \alpha_i \left[y^{(i)} (w^T x^{(i)} + b) - 1 + \xi_i \right] - \sum_{i=1}^m r_i \xi_i$$

Dual problem:

Soft Margin SVM



hard margin: separable

$$y^i (w^T x^i + b) \geq 1 - \xi_i$$

soft margin: allows unseparable samples.

The generalized Lagrangian function:

$$L(w, b, \xi, \alpha, r) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y^{(i)}(w^T x^{(i)} + b) - 1 + \xi_i] - \sum_{i=1}^m r_i \xi_i$$

Dual problem:

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

$$\text{s.t. } 0 \leq \alpha_i \leq C, i = 1, \dots, m$$

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0$$

w^* is the same as the non-regularizing case, but b^* has changed.

Soft Margin SVM

Dual problem:

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

$$\text{s.t. } 0 \leq \alpha_i \leq C, i = 1, \dots, m$$

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0$$

By the KKT dual-complementary conditions, for all i , $\alpha_i^* g_i(w^*) = 0$

$$\alpha_i = 0 \quad \iff$$

$$\alpha_i = C \quad \iff$$

$$0 < \alpha_i < C \quad \iff$$

Soft Margin SVM

Dual problem:

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

$$\text{s.t. } 0 \leq \alpha_i \leq C, i = 1, \dots, m$$

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0$$

By the KKT dual-complementary conditions, for all i , $\alpha_i^* g_i(w^*) = 0$

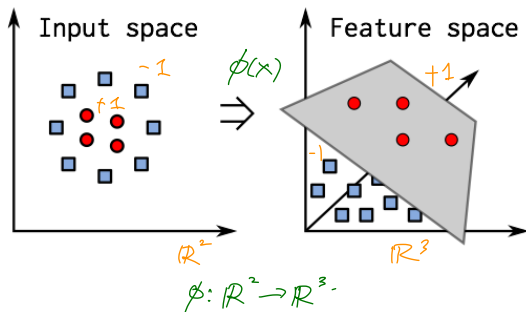
support vector $\alpha_i \neq 0$	<u>$\alpha_i = 0$</u>	\iff	<u>$y^{(i)}(w^T x^{(i)} + b) \geq 1$</u>	correct side of margin
	<u>$\alpha_i = C$</u>	\iff	<u>$y^{(i)}(w^T x^{(i)} + b) \leq 1$</u>	wrong side of margin
	<u>$0 < \alpha_i < C$</u>	\iff	<u>$y^{(i)}(w^T x^{(i)} + b) = 1$</u>	<u>at margin</u>

w^*

Kernel SVM

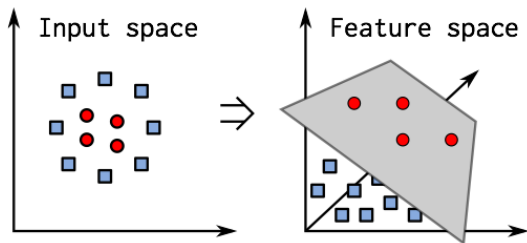
Non-linear SVM

For non-separable data, we can use the **kernel trick**: Map input values $x \in \mathbb{R}^d$ to a higher dimension $\phi(x) \in \mathbb{R}^D$, such that the data becomes separable.



Non-linear SVM

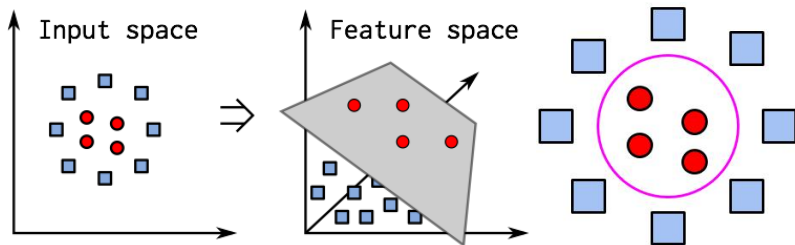
For non-separable data, we can use the **kernel trick**: Map input values $x \in \mathbb{R}^d$ to a higher dimension $\phi(x) \in \mathbb{R}^D$, such that the data becomes separable.



- ▶ ϕ is called a **feature mapping**.

Non-linear SVM

For non-separable data, we can use the **kernel trick**: Map input values $x \in \mathbb{R}^d$ to a higher dimension $\phi(x) \in \mathbb{R}^D$, such that the data becomes separable.



- ▶ ϕ is called a **feature mapping**.
- ▶ The classification function $w^T x + b$ becomes nonlinear: $\underline{w}^T \underline{\phi(x)} + \underline{b}$

Kernel Function

Given a feature mapping ϕ , we define the **kernel function** to be

$$\underline{K(x, z)} = \underline{\phi(x)}^T \underline{\phi(z)}$$

Kernel Function

Given a feature mapping ϕ , we define the kernel function to be

$$K(x, z) = \underline{\phi(x)}^T \underline{\phi(z)}$$

Some kernel functions are easier to compute than $\phi(x)$, e.g.

$$\chi: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\begin{aligned} \underline{K(x, z)} &= \underline{(x^T z)^2} = \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \right)^2 = (x_1 z_1 + x_2 z_2)(x_1 z_1 + x_2 z_2) \\ &= (x_1 z_1)^2 + (x_1 z_1)(x_1 z_2) + (x_2 z_2)(x_1 z_1) + (x_2 z_2)^2 \\ &= \underline{(x_1^2 z_1^2)} + \underline{(x_1 x_2)(z_1 z_2)} + \underline{(x_2 x_1)(z_2 z_1)} + \underline{(x_2^2 z_2^2)} \end{aligned}$$

$$x, z \in \mathbb{R}^2$$

$$\phi \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) = \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ x_2 x_1 \\ x_2^2 \end{bmatrix}$$

$$= \left\langle \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ x_2 x_1 \\ x_2^2 \end{bmatrix}, \begin{bmatrix} z_1^2 \\ z_1 z_2 \\ z_2 z_1 \\ z_2^2 \end{bmatrix} \right\rangle$$

$$= \underline{\langle \phi(x), \phi(z) \rangle}.$$

Kernel Function

Given a feature mapping ϕ , we define the **kernel function** to be

$$K(x, z) = \phi(x)^T \phi(z)$$

Some kernel functions are easier to compute than $\phi(x)$, e.g.

$$\begin{aligned} x, z \in \mathbb{R}^n. \\ K(x, z) &= \underbrace{(x^T z)}^2 = \underbrace{\left(\sum_{i=1}^n x_i z_i \right)}_{\phi(x)^T} \underbrace{\left(\sum_{j=1}^n x_j z_j \right)}_{\phi(z)} = \sum_{i=1}^n \sum_{j=1}^n x_i x_j z_i z_j \\ &= \phi(x)^T \phi(z) \end{aligned}$$

Kernel Function

Given a feature mapping ϕ , we define the **kernel function** to be

$$K(x, z) = \phi(x)^T \phi(z)$$

Some kernel functions are easier to compute than $\phi(x)$, e.g.

A kernel example:

$$\begin{aligned} \underline{K(x, z)} &= \underline{(x^T z)^2} = \left(\sum_{i=1}^n x_i z_i \right) \left(\sum_{j=1}^n x_j z_j \right) = \sum_{i=1}^n \sum_{j=1}^n x_i x_j z_i z_j \\ &= \underline{\phi(x)^T} \underline{\phi(z)} \end{aligned}$$

where $\underline{\phi(x)} = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_n \\ x_2 x_1 \\ x_n x_{n-1} \\ x_n x_n \end{bmatrix}$ takes $\underline{O(n^2)}$ operations to compute, while $(x^T z)^2$ only takes $\underline{O(n)}$

Kernel SVM

In the dual problem, replace $\langle x_i, x_j \rangle$ with $\langle \phi(x_i), \phi(y_i) \rangle = K(x_i, x_j)$

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \underbrace{K(x_i, x_j)}_{\langle \phi(x_i), \phi(x_j) \rangle}$$

$$\text{s.t. } 0 \leq \alpha_i \leq C, i = 1, \dots, m$$

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0$$

Kernel SVM

In the dual problem, replace $\langle x_i, y_j \rangle$ with $\langle \phi(x_i), \phi(y_j) \rangle = K(x_i, x_j)$

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \underline{K(x_i, x_j)}$$

$$\text{s.t. } 0 \leq \alpha_i \leq C, i = 1, \dots, m$$

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0$$

No need to compute $w^* = \sum_{i=1}^m \alpha_i^* y^{(i)} \phi(x^{(i)})$ explicitly since

Given new sample x .

$$\underline{f(x)} = \underline{w^T} \phi(x) + b = \left(\sum_{i=1}^m \underline{\alpha_i y^{(i)}} \phi(x^{(i)}) \right)^T \cdot \underline{\phi(x)} + b$$

$$= \sum_{i=1}^m \alpha_i y^{(i)} \langle \phi(x^{(i)}), \phi(x) \rangle + b$$

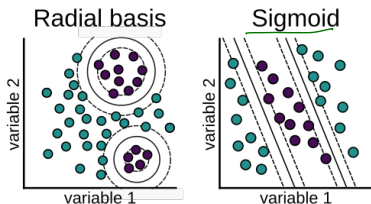
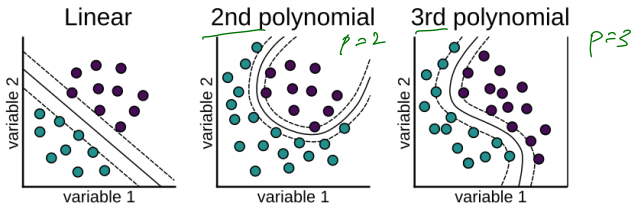
$$= \sum_{i=1}^m \alpha_i y^{(i)} \underline{K(x^{(i)}, x)} + b$$

$\sum_{i=1}^m \alpha_i \neq 0$

Kernel Matrix

kernel functions measure the similarity between samples x, z , e.g.

- ▶ Linear kernel: $K(x, z) = \underline{x^T z}$
- ▶ Polynomial kernel: $K(x, z) = \underline{(x^T z + 1)^p}$
- ▶ Gaussian / radial basis function (RBF) kernel:
 $K(x, z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)$



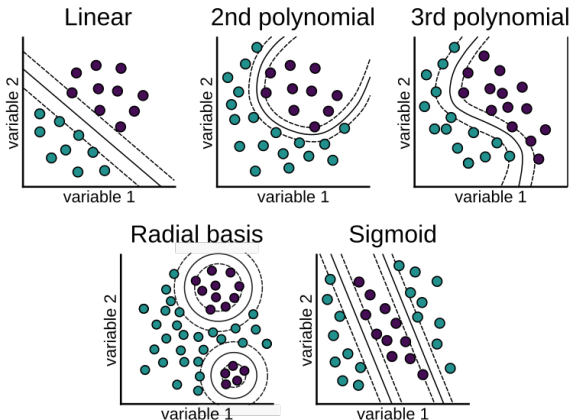
Kernel Matrix

$$K(x_1, x_2)$$

$$\forall x, z \in \mathbb{R}^n.$$

kernel functions measure the similarity between samples x, z , e.g.

- ▶ Linear kernel: $K(x, z) = (x^T z)$
- ▶ Polynomial kernel: $K(x, z) = (x^T z + 1)^p$
- ▶ Gaussian / radial basis function (RBF) kernel:
$$K(x, z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)$$

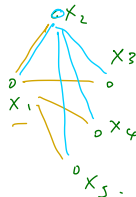


Can any function $K(x, y)$ be a kernel function?

$$\| \langle \phi(x), \phi(y) \rangle \text{ for some } \phi ?$$

Kernel Matrix

Given m training samples $(\underline{x_1}, \underline{y_1}), \dots, (\underline{x_m}, \underline{y_m})$



Represent kernel function as a matrix $\underline{K} \in \underline{\mathbb{R}^{m \times m}}$ where $K_{i,j} = K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$.

$$\underline{K} = \begin{bmatrix} \phi(x_1)^T \phi(x_1) & \phi(x_1)^T \phi(x_2) & \dots & \phi(x_1)^T \phi(x_m) \\ \phi(x_2)^T \phi(x_1) & \dots & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots \\ \phi(x_m)^T \phi(x_1) & \dots & \dots & \phi(x_m)^T \phi(x_m) \end{bmatrix}$$

$x_1 \rightarrow \begin{bmatrix} \phi(x_1)^T \phi(x_1) \\ \phi(x_1)^T \phi(x_2) \\ \vdots \\ \phi(x_1)^T \phi(x_m) \end{bmatrix}$

$\begin{bmatrix} K(x_1, x_1) \\ K(x_1, x_2) \\ \vdots \\ K(x_1, x_m) \end{bmatrix}$

Kernel Matrix

How to prove function K is a valid kernel?

e.g. $K(x, z) = (x^T z)^2$. $\xrightarrow{(1)}$ $\langle \phi(x), \phi(z) \rangle$
 $\xrightarrow{(2)}$ applying Mercer theorem

Represent kernel function as a matrix $K \in \mathbb{R}^{m \times m}$ where let $a \in \mathbb{R}^m$

$$K_{i,j} = K(x_i, x_j) = \phi(x_i)^T \phi(x_j).$$

Show $a^T K a \geq 0$.

\Rightarrow K is a valid kernel

$(x_i^T x_j + c)^p$
Theorem (Mercer)

Let $K : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$ Then K is a valid (Mercer) kernel if and only if for any finite training set $\{x^{(1)}, \dots, x^{(m)}\}$, K is symmetric positive semi-definite.

i.e. $K_{i,j} = K_{j,i}$ and $x^T K x \geq 0$ for all $x \in \mathbb{R}^m$

$$K = K^T$$

Kernel SVM Summary

- ▶ Input: m training samples $(x^{(i)}, y^{(i)})$, $y^i \in \{-1, 1\}$, kernel function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, constant $C > 0$
- ▶ Output: non-linear decision function $f(x)$
- ▶ Step 1: solve the dual optimization problem for α^*

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j K(x^{(i)}, x^{(j)})$$

$$s.t. 0 \leq \alpha_i \leq C, \sum_{i=1}^m \alpha_i y^{(i)} = 0, i = 1, \dots, m$$

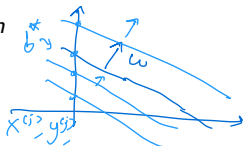
- ▶ Step 2: compute the optimal decision function

$$b^* = y^{(j)} - \sum_{i=1}^m \alpha_i^* y^{(i)} K(x^{(i)}, x^{(j)}) \text{ for some } 0 < \alpha_j < C$$

for some $x^{(j)}, y^{(j)}$

$$f(x) = \sum_{i=1}^m \alpha_i y^{(i)} K(x^{(i)}, x) + b^*$$

In practice, it's more efficient to compute kernel matrix K in advance. ($m \times m$)



SVM in Practice

(SMO) \rightarrow , coordinate ascent α_i, α_j

Sequential Minimal Optimization: a fast algorithm for training soft margin kernel SVM

- ▶ Break a large SVM problem into smaller chunks, update two α_i 's at a time
- ▶ Implemented by most SVM libraries.

SVM in Practice

Sequential Minimal Optimization: a fast algorithm for training soft margin kernel SVM

- ▶ Break a large SVM problem into smaller chunks, update two α_i 's at a time
- ▶ Implemented by most SVM libraries.

Other related algorithms

- ▶ Support Vector Regression (SVR)
- ▶ Multi-class SVM (Koby Crammer and Yoram Singer. 2002. *On the algorithmic implementation of multiclass kernel-based vector machines*. J. Mach. Learn. Res. 2 (March 2002), 265-292.)

- least square SVM

SVM in Practice

Sequential Minimal Optimization: a fast algorithm for training soft margin kernel SVM

- ▶ Break a large SVM problem into smaller chunks, update two α_i 's at a time
- ▶ Implemented by most SVM libraries.

Other related algorithms

- ▶ Support Vector Regression (SVR)
- ▶ Multi-class SVM (Koby Crammer and Yoram Singer. 2002. *On the algorithmic implementation of multiclass kernel-based vector machines*. J. Mach. Learn. Res. 2 (March 2002), 265-292.)