# Support Vector Machines (SVM)

## Linear separation of a feature space

A hyper plane in an n-D feature space can be represented by the following equation:

$$f(\mathbf{x}) = \mathbf{x}^T\mathbf{w} + b = \sum_{i=1}^{n} x_i w_i + b = 0$$

Dividing by $||\mathbf{w}||$, we get

$$\frac{\mathbf{x}^T\mathbf{w}}{||\mathbf{w}||} = P_\mathbf{w}(\mathbf{x}) = -\frac{b}{||\mathbf{w}||}$$

indicating that the projection of any point $\mathbf{x}$ on the plane onto the vector $\mathbf{w}$ is always $-b/||\mathbf{w}||$, i.e., $\mathbf{w}$ is the normal direction of the plane, and $|b|/||\mathbf{w}||$ is the distance from the origin to the plane. Note that the equation of the hyper plane is not unique. $c\,f(\mathbf{x}) = 0$ represents the same plane for any $c$.

The n-D space is partitioned into two regions by the plane. Specifically, we define a mapping function $y = sign(f(\mathbf{x})) \in \{1, -1\}$,

$$f(\mathbf{x}) = \mathbf{x}^T\mathbf{w} + b = \begin{cases} > 0, & y = sign(f(\mathbf{x})) = 1, \ \mathbf{x} \in P \\ < 0, & y = sign(f(\mathbf{x})) = -1, \ \mathbf{x} \in N \end{cases}$$

Any point $\mathbf{x} \in P$ on the positive side of the plane is mapped to 1, while any point $\mathbf{x} \in N$ on the negative side is mapped to -1. A point $\mathbf{x}$ of unknown class will be classified to P if $f(\mathbf{x}) > 0$, or N if $f(\mathbf{x}) < 0$.

**Example:**

A straight line in 2D space $\mathbf{x} = [x_1, x_2]^T$ described by the following equation:

$$f(\mathbf{x}) = \mathbf{x}^T\mathbf{w} + b = [x_1, x_2]\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + b = [x_1, x_2]\begin{bmatrix} 1 \\ 2 \end{bmatrix} - 1 = x_1 + 2x_2 - 1 = 0$$

devides the 2D plane into two halves. The distance between the origin and the line is

$$\frac{|b|}{||\mathbf{w}||} = \frac{1}{\sqrt{w_1^2 + w_2^2}} = \frac{1}{\sqrt{5}} = 0.447$$

Consider three points:

- $\mathbf{x}_0 = [0.5,\ 0.25]^T$, $f(\mathbf{x}_0) = 0.5 + 2 \times 0.25 - 1 = 0$, i.e., $\mathbf{x}_0$ is on the plane;

- $\mathbf{x}_1 = [1,\ 0.25]^T$, $f(\mathbf{x}_1) = 1 + 2 \times 0.25 - 1 = 0.5 > 0$, i.e., $\mathbf{x}_1$ is above the straight line;

- $\mathbf{x}_2 = [0.5,\ 0]^T$, $f(\mathbf{x}_2) = 0.5 + 2 \times 0 - 1 = -0.5 < 0$, i.e., $\mathbf{x}_2$ is below the straight line.

## The learning problem

Given a set $K$ training samples from two linearly separable classes P and N:

$$\{(\mathbf{x}_k, y_k), k = 1, \cdots, K\}$$

where $y_k \in \{1, -1\}$ labels $\mathbf{x}_k$ to belong to either of the two classes. we want to find a hyper-plane in terms of $\mathbf{w}$ and $b$, that linearly separates the two classes.

Before the classifier is properly trained, the actual output $y' = sign(f(\mathbf{x}))$ may not be the same as the desired output $y$. There are four possible cases:

|   | Input $(\mathbf{x}, y)$ | Output $y' = sign(f(\mathbf{x}))$ | result |
|---|---|---|---|
| 1 | $(\mathbf{x}, y = 1)$ | $y' = 1 = y$ | corrrect |
| 2 | $(\mathbf{x}, y = -1)$ | $y' = 1 \neq y$ | incorrect |
| 3 | $(\mathbf{x}, y = 1)$ | $y' = -1 \neq y$ | incorrect |
| 4 | $(\mathbf{x}, y = -1)$ | $y' = -1 = y$ | corrrect |

The weight vector $\mathbf{w}$ is updated whenever the result is incorrect (mistake driven):

- If $(\mathbf{x}, y = -1)$ but $y' = 1 \neq y$ (case 2 above), then

$$\mathbf{x}^{new} = \mathbf{w}^{old} + \eta y \mathbf{x} = \mathbf{w}^{old} - \eta \mathbf{x}$$

When the same $\mathbf{x}$ is presented again, we have

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w}^{new} + b = \mathbf{x}^T \mathbf{w}^{old} - \eta \mathbf{x}^T \mathbf{x} + b < \mathbf{x}^T \mathbf{w}^{old} + b$$

The output $y' = sign(f(\mathbf{x}))$ is more likely to be $y = -1$ as desired. Here $0 < \eta < 1$ is the learning rate.

- If $(\mathbf{x}, y = 1)$ but $y' = -1 \neq y$ (case 3 above), then

$$\mathbf{w}^{new} = \mathbf{w}^{old} + \eta y \mathbf{x} = \mathbf{w}^{old} + \eta \mathbf{x}$$

When the same $\mathbf{x}$ is presented again, we have

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w}^{new} + b = \mathbf{x}^T \mathbf{w}^{old} + \eta \mathbf{x}^T \mathbf{x} + b > \mathbf{x}^T \mathbf{w}^{old} + b$$

The output $y' = sign(f(\mathbf{x}))$ is more likely to be $y = 1$ as desired.

Summarizing the two cases:

$$\text{if } \ yf(\mathbf{x}) = y(\mathbf{x}^T \mathbf{w}^{old} + b) < 0, \ \text{ then } \ \mathbf{w}^{new} = \mathbf{w}^{old} + \eta y \mathbf{x}$$

The two correct cases (cases 1 and 4) can also be summarized as

$$yf(\mathbf{x}) = y(\mathbf{x}^T \mathbf{w} + b) \geq 0$$

which is the condition a successful classifier should satisfy.

We assume initially $\mathbf{w} = 0$, and the $K$ training samples are presented repeatedly, the training will yield:

$$\mathbf{w} = \sum_{i=1}^{K} \alpha_i y_i \mathbf{x}_i$$

where $\alpha_i > 0$. Note that $\mathbf{w}$ is expressed as a linear combination of the training samples. After receiving a new sample $(\mathbf{x}_i, y_i)$, vector $\mathbf{w}$ is updated by

$$\text{if } \ y_i f(\mathbf{x}_i) = y_i(\mathbf{x}_i^T \mathbf{w}^{old} + b) = y_i \left( \sum_{j=1}^{m} \alpha_j y_j (\mathbf{x}_i^T \mathbf{x}_j) + b \right) < 0,$$

$$\text{then } \ \mathbf{w}^{new} = \mathbf{w}^{old} + \eta y_i \mathbf{x}_i = \sum_{j=1}^{m} \alpha_j y_j \mathbf{x}_j + \eta y_i \mathbf{x}_i, \quad \text{i.e.} \quad \alpha_i^{new} = \alpha_i^{old} + \eta$$

Now both the decision function

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w} + b = \sum_{j=1}^{m} \alpha_j y_j (\mathbf{x}^T \mathbf{x}_j) + b$$

and the learning law

$$\text{if } \ y_i \left( \sum_{j=1}^{m} \alpha_j y_j (\mathbf{x}_i^T \mathbf{x}_j) + b \right) < 0, \quad \text{then } \ \alpha_i^{new} = \alpha_i^{old} + \eta$$

are expressed in terms of the inner production of input vectors.

## SVM Dereivations

For a decision hyper-plane $\mathbf{x}^T\mathbf{w} + b = 0$ to separate the two classes $P = \{(\mathbf{x}_i, 1)\}$ and $N = \{(\mathbf{x}_i, -1)\}$, it has to satisfy

$$y_i(\mathbf{x}_i^T\mathbf{w} + b) \geq 0$$

for both $\mathbf{x}_i \in P$ and $\mathbf{x}_i \in N$. Among all such planes satisfying this condition, we want to find the optimal one $H_0$ that separates the two classes with the maximal margin (the distance between the decision plane and the closest sample points).

The optimal plane should be in the middle of the two classes, so that the distance from the plane to the closest point on either side is the same. We define two additional planes $H_+$ and $H_-$ that are parallel to $H_0$ and go through the point closest to the plane on either side:

$$\mathbf{x}^T\mathbf{w} + b = 1, \quad \text{and} \quad \mathbf{x}^T\mathbf{w} + b = -1$$

All points $\mathbf{x}_i \in P$ on the positive side should satisfy

$$\mathbf{x}_i^T\mathbf{w} + b \geq 1, \quad y_i = 1$$

and all points $\mathbf{x}_i \in N$ on the negative side should satisfy

$$\mathbf{x}_i^T\mathbf{w} + b \leq -1, \quad y_i = -1$$

These can be combined into one inequality:

$$y_i(\mathbf{x}_i^T\mathbf{w} + b) \geq 1, \quad (i = 1, \cdots, m)$$

The equality holds for those points on the planes $H_+$ or $H_-$. Such points are called *support vectors*, for which

$$\mathbf{x}_i^T\mathbf{w} + b = y_i$$

i.e., the following holds for all support vectors:

$$b = y_i - \mathbf{x}_i^T\mathbf{w} = y_i - \sum_{j=1}^m \alpha_j y_j(\mathbf{x}_i^T\mathbf{x}_j)$$

Moreover, the distances from the origin to the three parallel planes $H_-$, $H_0$ and $H_+$ are, respectively, $|b-1|/||\mathbf{w}||$, $|b|/||\mathbf{w}||$, and $|b+1|/||\mathbf{w}||$, and the distance between planes $H_-$ and $H_+$ is $2/||\mathbf{w}||$.

Our goal is to maximize this distance, or, equivalantly, to minimize the norm $||\mathbf{w}||$. Now the problem of finding the optimal decision plane in terms of $\mathbf{w}$ and $b$ can be formulated as:

$$\text{minimize} \quad \frac{1}{2}\mathbf{w}^T\mathbf{w} = \frac{1}{2}||\mathbf{w}||^2 \quad \text{(objective function)}$$

$$\text{subject to} \quad y_i(\mathbf{x}_i^T\mathbf{w} + b) \geq 1, \quad \text{or} \quad 1 - y_i(\mathbf{x}_i^T\mathbf{w} + b) \leq 0, \quad (i = 1, \cdots, m)$$

Since the objective function is quadratic, this constrained optimization problem is called a quadratic program (QP) problem. (If the objective function is linear instead, the problem is a linear program (LP) problem). This QP problem can be solved by Lagrange multipliers method to minimize the following

$$L_p(\mathbf{w}, b, \alpha) = \frac{1}{2}||\mathbf{w}||^2 + \sum_{i=1}^m \alpha_i(1 - y_i(\mathbf{x}_i^T\mathbf{w} + b))$$

with respect to $\mathbf{w}$, $b$ and the Lagrange coefficients $\alpha_i \geq 0$ $(i = 1, \cdots, \alpha_m)$. We let

$$\frac{\partial}{\partial W}L_p(\mathbf{w}, b) = 0, \quad \frac{\partial}{\partial b}L_p(\mathbf{w}, b) = 0$$

These lead, respectively, to

$$\mathbf{w} = \sum_{j=1}^m \alpha_j y_j \mathbf{x}_j, \quad \text{and} \quad \sum_{i=1}^m \alpha_i y_i = 0$$

Substituting these two equations back into the expression of $L(\mathbf{w}, b)$, we get the *dual problem* (with respect to $\alpha_i$) of the above *primal problem*:

$$\text{maximize} \quad L_d(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2}\sum_{i=1}^m\sum_{j=1}^m \alpha_i\alpha_j y_i y_j \mathbf{x}_i^T, \mathbf{x}_j$$

$$\text{subject to} \quad \alpha_i \geq 0, \quad \sum_{i=1}^m \alpha_i y_i = 0$$

The dual problem is related to the primal problem by:

$$L_d(\alpha) = inf_{(\mathbf{w},b)}L_p(\mathbf{w}, b, \alpha)$$

5

i.e., $L_d$ is the greatest lower bound (infimum) of $L_p$ for all $\mathbf{w}$ and $b$.

Solving this dual problem (an easier problem than the primal one), we get $\alpha_i$, from which $\mathbf{w}$ of the optimal plane can be found.

Those points $\mathbf{x}_i$ on either of the two planes $H_+$ and $H_-$ (for which the equality $y_i(\mathbf{w}^T\mathbf{x}_i + b) = 1$ holds) are called *support vectors* and they correspond to positive Lagrange multipliers $\alpha_i > 0$. The training depends only on the support vectors, while all other samples away from the planes $H_+$ and $H_-$ are not important.

For a support vector $\mathbf{x}_i$ (on the $H_-$ or $H_+$ plane), the constraining condition is

$$y_i\left(\mathbf{x}_i^T\mathbf{w} + b\right) = 1 \quad (i \in sv)$$

here $sv$ is a set of all indices of support vectors $\mathbf{x}_i$ (corresponding to $\alpha_i > 0$). Substituting

$$\mathbf{w} = \sum_{j=1}^{m} \alpha_j y_j \mathbf{x}_j = \sum_{j \in sv} \alpha_j y_j \mathbf{x}_j$$

we get

$$y_i\left(\sum_{j \in sv} \alpha_j y_j \mathbf{x}_i^T \mathbf{x}_j + b\right) = 1$$

Note that the summation only contains terms corresponding to those support vectors $\mathbf{x}_j$ with $\alpha_j > 0$, i.e.

$$y_i \sum_{j \in sv} \alpha_j y_j \mathbf{x}_i^T \mathbf{x}_j = 1 - y_i b$$

For the optimal weight vector $\mathbf{w}$ and optimal $b$, we have:

$$
\begin{aligned}
||\mathbf{w}||^2 &= \mathbf{w}^T\mathbf{w} = \sum_{i \in sv} \alpha_i y_i \mathbf{x}_i^T \sum_{j \in sv} \alpha_j y_j \mathbf{x}_j = \sum_{i \in sv} \alpha_i y_i \sum_{j \in sv} \alpha_j y_j \mathbf{x}_i^T \mathbf{x}_j \\
&= \sum_{i \in sv} \alpha_i(1 - y_i b) = \sum_{i \in sv} \alpha_i - b \sum_{i \in sv} \alpha_i y_i \\
&= \sum_{i \in sv} \alpha_i
\end{aligned}
$$

The last equality is due to $\sum_{i=1}^{m} \alpha_i y_i = 0$ shown above. Recall that the distance between the two margin planes $H_+$ and $H_-$ is $2/||\mathbf{w}||$, and the margin, the distance between $H_+$ (or $H_-$) and the optimal decision plane $H_0$, is

$$\frac{1}{||\mathbf{w}||} = \left(\sum_{i \in sv} \alpha_i\right)^{-1/2}$$

## Soft Margin SVM

When the two classes are not linearly separable (e.g., due to noise), the condition for the optimal hyper-plane can be relaxed by including an extra term:

$$y_i(\mathbf{x}_i^T \mathbf{w} + b) \geq 1 - \xi_i, \quad (i = 1, \cdots, m)$$

For minimum error, $\xi_i \geq 0$ should be minimized as well as $||\mathbf{w}||$, and the objective function becomes:

$$\text{minimize} \quad \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^{m} \xi_i^k$$
$$\text{subject to} \quad y_i(\mathbf{x}_i^T \mathbf{w} + b) \geq 1 - \xi_i, \quad \text{and} \quad \xi_i \geq 0; \quad (i = 1, \cdots, m)$$

Here $C$ is a regularization parameter that controls the trade-off between maximizing the margin and minimizing the training error. Small C tends to emphasize the margin while ignoring the outliers in the training data, while large C may tend to overfit the training data.

When $k = 2$, it is called 2-norm soft margin problem:

$$\text{minimize} \quad \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^{m} \xi_i^2$$
$$\text{subject to} \quad y_i(\mathbf{x}_i^T \mathbf{w} + b) \geq 1 - \xi_i, \quad (i = 1, \cdots, m)$$

$$(1)$$

Note that the condition $\xi_i \geq 0$ is dropped, as if $\xi_i < 0$, we can set it to zero and the objective function is further reduced.) Alternatively, if we let $k = 1$, the problem can be formulated as

$$\text{minimize} \quad \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^{m} \xi_i$$
$$\text{subject to} \quad y_i(\mathbf{x}_i^T \mathbf{w} + b) \geq 1 - \xi_i \quad \text{and} \quad \xi_i \geq 0; \quad (i = 1, \cdots, m)$$

$$(2)$$

This is called 1-norm soft margin problem. The algorithm based on 1-norm setup, when compared to 2-norm algorithm, is less sensitive to outliers in training data. When the data is noisy, 1-norm method should be used to ignore the outliers.

**L2-Norm Soft Margin**

The primal Lagrangian for 2-norm problem above is

$$L_p(\mathbf{w}, b, \xi, \alpha) = \frac{1}{2}\mathbf{w}^T\mathbf{w} + \frac{C}{2}\sum_{i=1}^{m}\xi_i^2 - \sum_{i=1}^{m}\alpha_i[y_i(\mathbf{w}^T\mathbf{x} + b) - 1 + \xi_i]$$

Substituting

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^{m}y_i\alpha_i\mathbf{x}_i = 0; \quad \frac{\partial L}{\partial \xi} = C\xi - \alpha = 0; \quad \frac{\partial L}{\partial b} = \sum_{i=1}^{m}y_i\alpha_i = 0$$

into the primal Lagrangian, we get the dual problem

$$\text{maximize} \quad L_d(\alpha) = \sum_{i=1}^{m}\alpha_i - \frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m}y_iy_j\alpha_i\alpha_j\mathbf{x}_j^T\mathbf{x}_i - \frac{1}{2C}\sum_{i=1}^{m}\alpha_i^2$$

$$= \sum_{i=1}^{m}\alpha_i - \frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m}y_iy_j\alpha_i\alpha_j(\mathbf{x}_j^T\mathbf{x}_i + \frac{1}{C}\delta_{ij})$$

$$\text{subject to} \quad \alpha_i \geq 0, \quad \sum_{i=1}^{m}\alpha_iy_i = 0$$

This QP program can be solved for $\alpha_i$. All support vectors $\mathbf{x}_i$ corresponding to $\alpha_i > 0$ satisfy:

$$y_i(\mathbf{x}_i^T\mathbf{w} + b) = 1 - \xi_i$$

Substituting $\mathbf{w} = \sum_{j \in sv}y_j\alpha_j\mathbf{x}_j$ into this equation, we get

$$y_i(\sum_{j \in sv}y_j\alpha_j(\mathbf{x}_i^T\mathbf{x}_j) + b) = 1 - \xi_i, \quad \text{i.e.,} \quad y_i\sum_{j \in sv}y_j\alpha_j(\mathbf{x}_i^T\mathbf{x}_j) = 1 - \xi_i - y_ib$$

For the optimal weight $\mathbf{w}$, we have

$$
\begin{aligned}
||\mathbf{w}||^2 &= \mathbf{w}^T\mathbf{w} = \sum_{i \in sv}\alpha_iy_i\mathbf{x}_i^T\sum_{j \in sv}\alpha_jy_j\mathbf{x}_j = \sum_{i \in sv}\alpha_iy_i\sum_{j \in sv}\alpha_jy_j\mathbf{x}_i^T\mathbf{x}_j \\
&= \sum_{i \in sv}\alpha_i(1 - \xi_i - y_ib) = \sum_{i \in sv}\alpha_i - \sum_{i \in sv}\alpha_i\xi_i - b\sum_{i \in sv}y_i\alpha_i \\
&= \sum_{i \in sv}\alpha_i - \sum_{i \in sv}\alpha_i\xi_i = \sum_{i \in sv}\alpha_i - \frac{1}{C}\sum_{i \in sv}\alpha_i^2
\end{aligned}
$$

The last equation is due to $\xi_i = \alpha_i/C$. The optimal margin is

$$1/||\mathbf{w}|| = (\sum_{i \in sv}\alpha_i - \frac{1}{C}\sum_{i \in sv}\alpha_i^2)^{-1/2}$$

**L1-Norm Soft Margin**

The primal Lagrangian for 1-norm problem above is

$$L_p(\mathbf{w}, b, \xi, \alpha, \gamma) = \frac{1}{2}\mathbf{w}^T\mathbf{w} + C\sum_{i=1}^{m}\xi_i - \sum_{i=1}^{m}\alpha_i[y_i(\mathbf{w}^T\mathbf{x} + b) - 1 + \xi_i] - \sum_{i=1}^{m}\gamma_i\xi_i$$

with $\alpha_i \geq 0$ and $\gamma_i \geq 0$. Substituting

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^{m}y_i\alpha_i\mathbf{x}_i = 0; \quad \frac{\partial L}{\partial \xi} = C - \alpha_i - \gamma_i = 0; \quad \frac{\partial L}{\partial b} = \sum_{i=1}^{m}y_i\alpha_i = 0$$

into the primal Lagrangian, we get the dual problem

$$\text{maximize} \quad L_d(\alpha, \gamma) = \sum_{i=1}^{m}\alpha_i - \frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m}y_iy_j\alpha_i\alpha_j\mathbf{x}_j^T\mathbf{x}_i - \sum_{i=1}^{m}\alpha_i\xi_i - \sum_{i=1}^{m}\gamma_i\xi_i + C\sum_{i=1}^{m}\xi_i$$

$$= \sum_{i=1}^{m}\alpha_i - \frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m}y_iy_j\alpha_i\alpha_j\mathbf{x}_j^T\mathbf{x}_i$$

$$\text{subject to} \quad 0 \leq \alpha_i \leq C, \quad \sum_{i=1}^{m}\alpha_iy_i = 0$$

Note that interestingly the objective function of the dual problem is identical to that of the linearly separable problem discussed previously, due to the nice cancellation based on $C = \alpha_i + \gamma_i$. Also, since $\alpha_i \geq 0$ and $\gamma_i \geq 0$, we have $0 \leq \alpha_i \leq C$. Solving this QP problem for $\alpha_i$, we get the optimal decision plane $\mathbf{w}$ and $b$ with the margin

$$\left(\sum_{i\in sv}\sum_{j\in sv}\alpha_i\alpha_jy_iy_j\mathbf{x}_i^T\mathbf{x}_j\right)^{-1/2}$$

# Kernel Mapping

The algorithm above converges only for linearly separable data. If the data set is not linearly separable, we can map the samples $\mathbf{x}$ into a feature space of higher dimensions:

$$\mathbf{x} \longrightarrow \phi(\mathbf{x})$$

in which the classes can be linearly separated. The decision function in the new space becomes:

$$f(\mathbf{x}) = \phi(\mathbf{x})^T\mathbf{w} + b = \sum_{j=1}^{m}\alpha_jy_j(\phi(\mathbf{x})^T\phi(\mathbf{x}_j)) + b$$

where
$$\mathbf{w} = \sum_{j=1}^{m} \alpha_j y_j \phi(\mathbf{x}_j)$$

and $b$ are the parameters of the decision plane in the new space. As the vectors $\mathbf{x}_i$ appear only in inner products in both the decision function and the mapping function $\phi(\mathbf{x})$ does not need to be explicitly specified. Instead, all we need is the inner product of the vectors in the new space. The function $\phi(\mathbf{x})$ is a kernel-induced *implicit* mapping.

**Definition:** A kernel is a function that takes two vectors $\mathbf{x}_i$ and $\mathbf{x}_j$ as arguments and returns the value of the inner product of their images $\phi(\mathbf{x}_i)$ and $\phi(\mathbf{x}_j)$:
$$K(\mathbf{x}_1, \mathbf{x}_2) = \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2)$$

As only the inner product of the two vectors in the new space is returned, the dimensionality of the new space is not important.

The learning algorithm in the kernel space can be obtained by replacing all inner products in the learning algorithm in the original space with the kernels:
$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w} + b = \sum_{j=1}^{m} \alpha_j y_j K(\mathbf{x}, \mathbf{x}_j) + b$$

The parameter $b$ can be found from any support vectors $\mathbf{x}_i$:
$$b = y_i - \phi(\mathbf{x}_i)^T \mathbf{w} = y_i - \sum_{j=1}^{m} \alpha_j y_j (\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)) = y_i - \sum_{j=1}^{m} \alpha_j y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

**Example 0:** linear kernel
Assume $\mathbf{x} = [x_1, \cdots, x_n]^T$, $\mathbf{z} = [z_1, \cdots, z_n]^T$,
$$K(\mathbf{x}, \mathbf{z}) = \mathbf{x}^T \mathbf{z} = \sum_{i=1}^{n} x_1 z_1$$

**Example 1:** polynomial kernels
Assume $\mathbf{x} = [x_1, x_2]^T$, $\mathbf{z} = [z_1, z_2]^T$,
$$
\begin{aligned}
K(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2 &= (x_1 z_1 + x_2 z_2)^2 = x_1^2 z_1^2 + x_2^2 z_2^2 + 2 x_1 z_1 x_2 z_2 \\
&= \ <(x_1^2, x_2^2, \sqrt{2} x_1 x_2), (z_1^2, z_2^2, \sqrt{2} z_1 z_2) > = \phi(\mathbf{x})^T \phi(\mathbf{z})
\end{aligned}
$$

This is a mapping from a 2-D space to a 3-D space. The order can be changed from 2 to general d.

**Example 2:**
$$K(\mathbf{x}, \mathbf{z}) = e^{-||\mathbf{x}-\mathbf{z}||^2/2\sigma^2}$$

**Example 3:**

$$K(\mathbf{x}, \mathbf{z}) = K(\mathbf{x}, \mathbf{z})K(\mathbf{x}, \mathbf{x})^{-1/2}K(\mathbf{z}, \mathbf{z})^{-1/2}$$