

# Kernel Methods

- Radial Basis Function (RBF) Networks
- Support Vector Machines (SVM)
- Kernel PCA
- Kernel k-Means

# Revisiting Class Separability

A 2-class problem is linearly separable if a decision hyperplane (defined by  $\mathbf{w}^*$ ) exists such that

$$(\mathbf{w}^*)^T \mathbf{x} > 0 \quad \forall \mathbf{x} \in \omega_1$$

$$(\mathbf{w}^*)^T \mathbf{x} < 0 \quad \forall \mathbf{x} \in \omega_2$$

Now, in a more general case, let us consider a mapping from the original input space to a space of a different dimension:

$$\varphi: \mathbb{R}^{m_0} \rightarrow \mathbb{R}^{m_1} \quad \text{where} \quad \varphi(\mathbf{x}) = [\varphi_1(\mathbf{x}) \ \varphi_2(\mathbf{x}) \ \dots \ \varphi_{m_1}(\mathbf{x})]^T$$

Then we can define a  $\varphi$ -separability: A 2-class problem is  $\varphi$ -linearly separable if there exists a  $\mathbf{w}^*$  such that

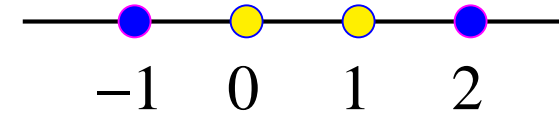
$$(\mathbf{w}^*)^T \varphi(\mathbf{x}) > 0 \quad \forall \mathbf{x} \in \omega_1$$

$$(\mathbf{w}^*)^T \varphi(\mathbf{x}) < 0 \quad \forall \mathbf{x} \in \omega_2$$

# Cover's Theorem

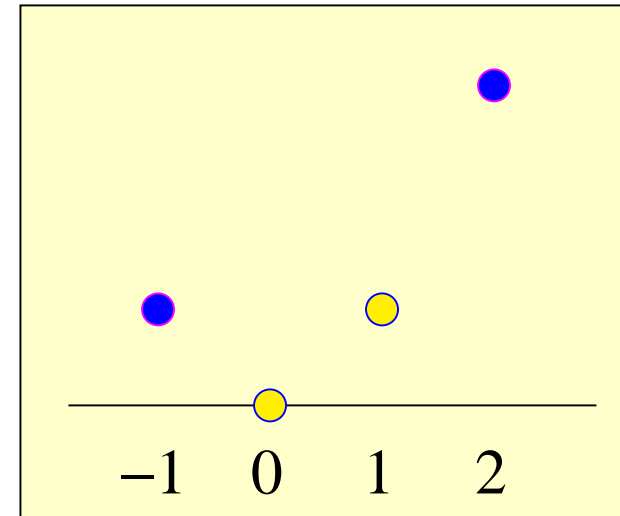
*A complex pattern recognition problem cast in a high-dimensional space nonlinearly is more likely to be linearly separable than in a low-dimensional space. (1965)*

Example: Non-separable 1-D data:



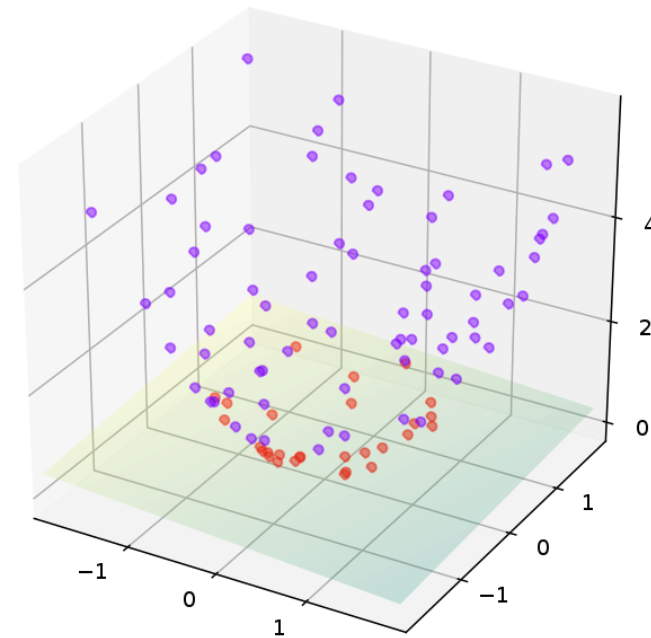
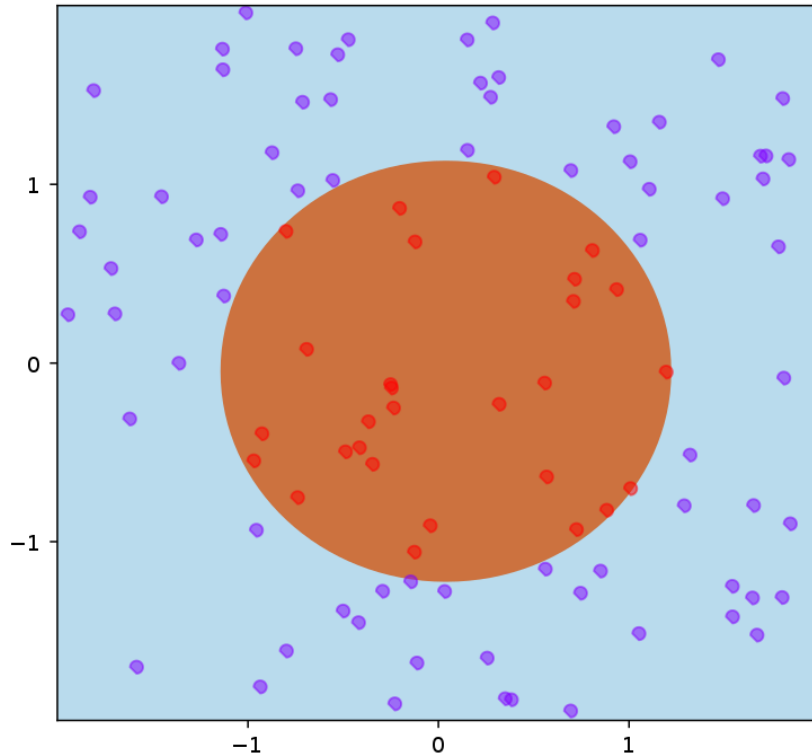
Mapped to 2-D by  $\phi(x) = [x \ x^2]^T$ :

Now linearly separable



# Cover's Theorem

An example 2-D to 3-D mapping to make the samples of two classes linearly separable:



[https://en.wikipedia.org/wiki/Cover%27s\\_theorem](https://en.wikipedia.org/wiki/Cover%27s_theorem)

# The Radial Basis Functions (RBF)

- Each radial basis function  $\varphi(\mathbf{x})$  has a center  $\mathbf{c}$ .
- $\varphi(\mathbf{x})$  is determined by  $\|\mathbf{x} - \mathbf{c}\|$ , i.e., the distance between  $\mathbf{x}$  and  $\mathbf{c}$ . (Euclidean distance is common.)
- A few forms of  $\varphi(\mathbf{x})$  :

- Gaussian: 
$$\varphi(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}\|^2}{2\sigma^2}\right)$$

- Multiquadrics: 
$$\varphi(\mathbf{x}) = \left(\|\mathbf{x} - \mathbf{c}\|^2 + a^2\right)^{1/2}$$

- Inverse multiquadrics: 
$$\varphi(\mathbf{x}) = \left(\|\mathbf{x} - \mathbf{c}\|^2 + a^2\right)^{-1/2}$$

# RBF Networks

Radial Basis Function (RBF) networks are neural networks where the activation functions of hidden neurons are radial basis functions.

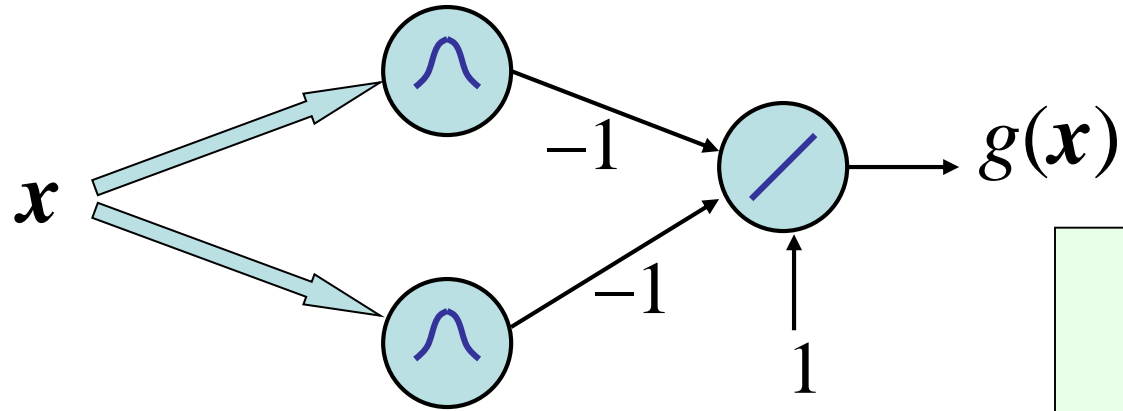
- A typical RBF network has one hidden layer.
- A typical RBF network uses linear activation for output neurons. Therefore, the output of a RBF network has the form of

$$F(\mathbf{x}) = \sum_{i=1}^{m_1} w_i \varphi(\|\mathbf{x} - \mathbf{c}_i\|)$$

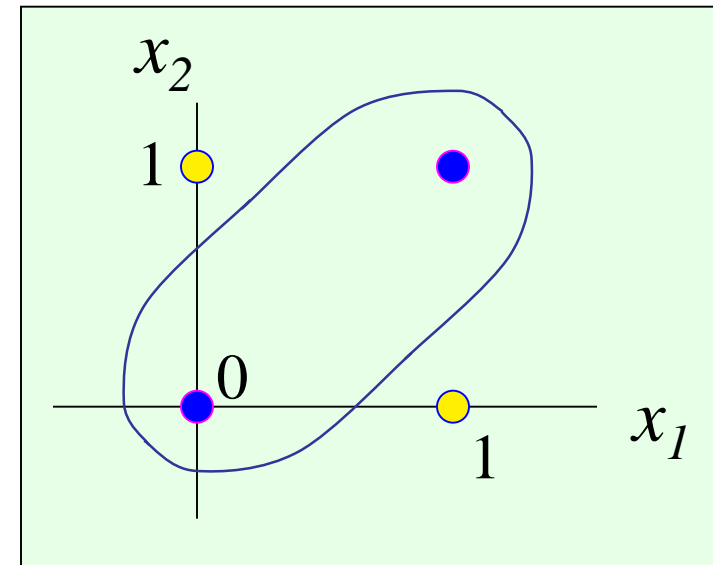
- RBF networks are universal approximators. (This can be easily understood by considering the case when we have one hidden neuron for each  $x_i$ .)
- The hidden neurons perform the mapping, as described for the Cover's theorem.

# Revisiting the XOR Problem

Let us consider the following RBF network with two RBF units centered at  $(0,0)$  and  $(1,1)$  and  $\sigma^2=0.5$ .



The decision boundary ( $g(x)=0$ ):



The final output is set to 1 for  $g(x)>0$  and 0 otherwise, to be consistent with Boolean logic.

# Training RBF Networks

Compared with MLPs, in addition to the weights, there are two more types of parameters that can be trained in RBF networks: the centers  $c_i$  and the standard deviations  $\sigma_i$ . This makes learning more complicated. The options are:

- Training all  $w$ ,  $c_i$  and  $\sigma_i$  together using gradient descent. This is computationally more expensive.
- Two-stage training: Derive the RBF parameters ( $c_i$  and  $\sigma_i$ ) first, and then the weights.



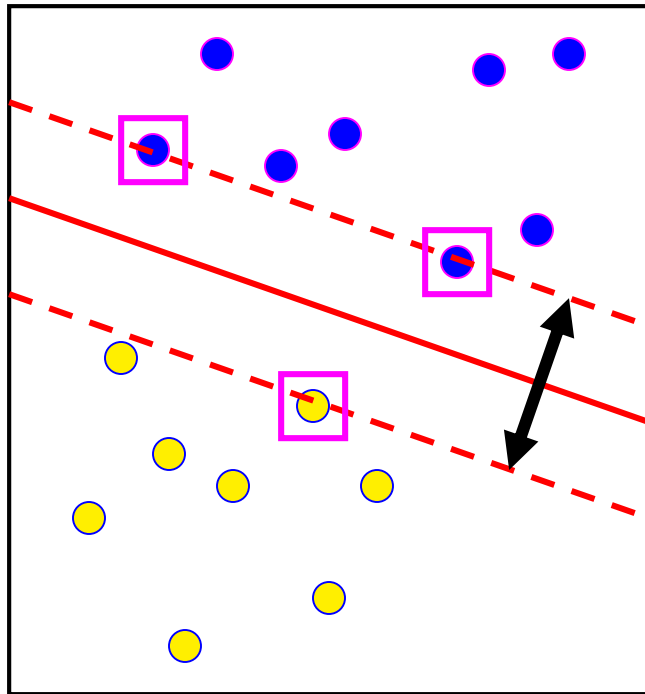
# Support Vector Machines (SVM)

Idea (SVM as a linear classifier):

Choose the hyperplane as far as possible from any  $x$

⇒ New samples near those  $x$  in the training set are more likely to be classified correctly.

⇒ Good generalization ability.



For linearly separable cases:

The **margin** is twice the distance from the hyperplane to the closest sample. The **support vectors** are the samples that are closest to the hyperplane.

# Margin in SVM

The linear discriminant function is  $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$

The distance between  $\mathbf{x}$  and the decision boundary ( $g(\mathbf{x})=0$ ) is given by  $\frac{|\mathbf{w}^T \mathbf{x} + w_0|}{\|\mathbf{w}\|}$

To make the values of  $\mathbf{w}$  and  $w_0$  unique, we add the requirement:  $|\mathbf{w}^T \mathbf{x} + w_0| = 1$  if  $\mathbf{x}$  is a support vector.

The margin is now given by  $2/\|\mathbf{w}\|$ . Maximizing the margin is the same as minimizing  $\|\mathbf{w}\|$ .

# The Optimization Problem

We can phrase the task of finding the decision boundary in SVM as an optimization problem:

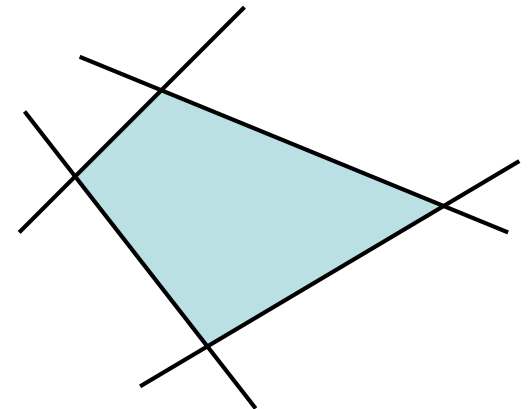
To minimize  $J = \frac{1}{2} \| \mathbf{w} \|^2$

Subject to the constraints  $y_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1, \forall i$

Here  $y_i = +1 \quad \forall \mathbf{x}_i \in \omega_1$   
 $y_i = -1 \quad \forall \mathbf{x}_i \in \omega_2$

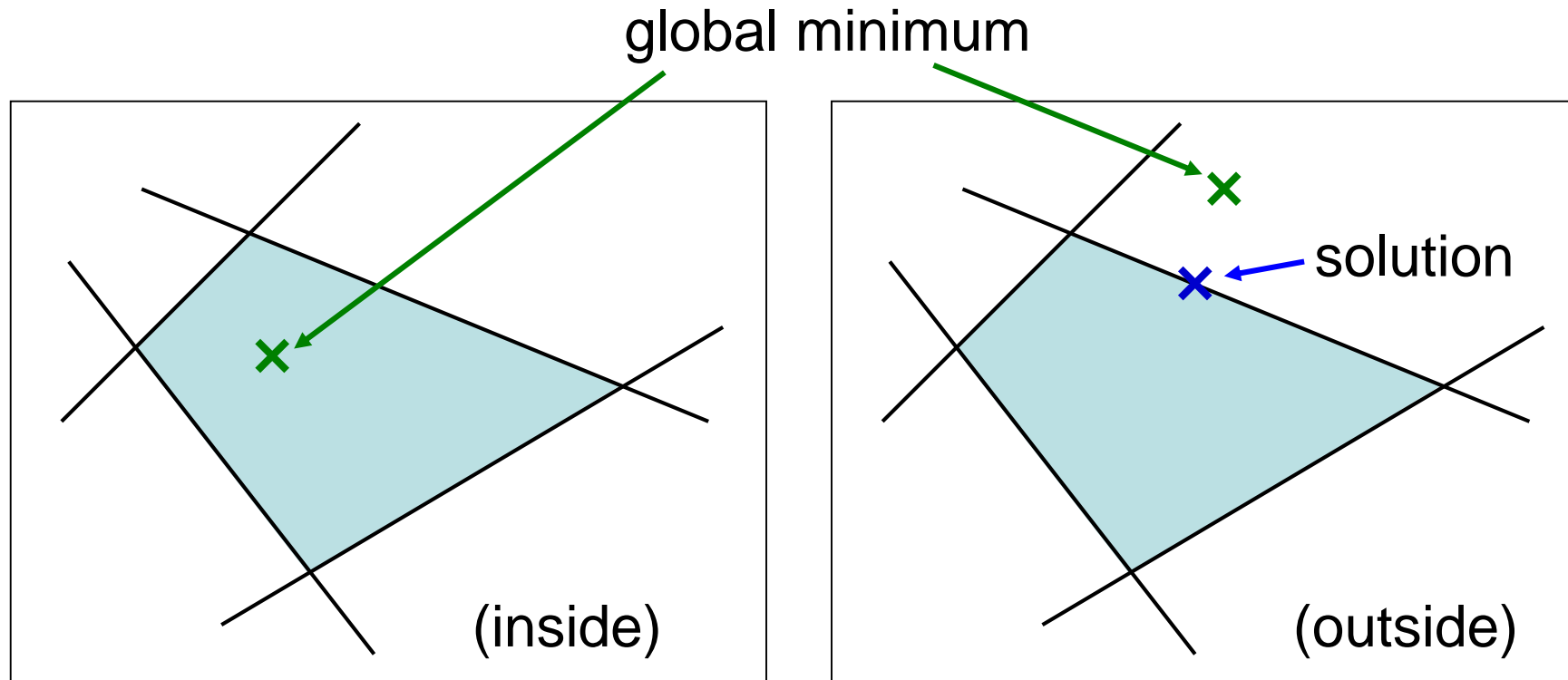
This is a **primal problem**, with a convex cost function (which has a single global minimum) and linear constraints. (Each training sample leads to a constraint.)

The (inequality) constraints define a feasible region of the solution (in the space of the weights/bias). Example:



# The Optimization Problem

The global minimum may fall inside or outside of the feasible region. (The constrained minimum is unique.)



There are active and inactive constraints.

# The Optimization Problem: Example

Consider the 1-D 2-class example:

$x_1=1, x_2=2$  in  $\omega_1$ , and  $x_3=-1, x_4=-2$  in  $\omega_2$ .

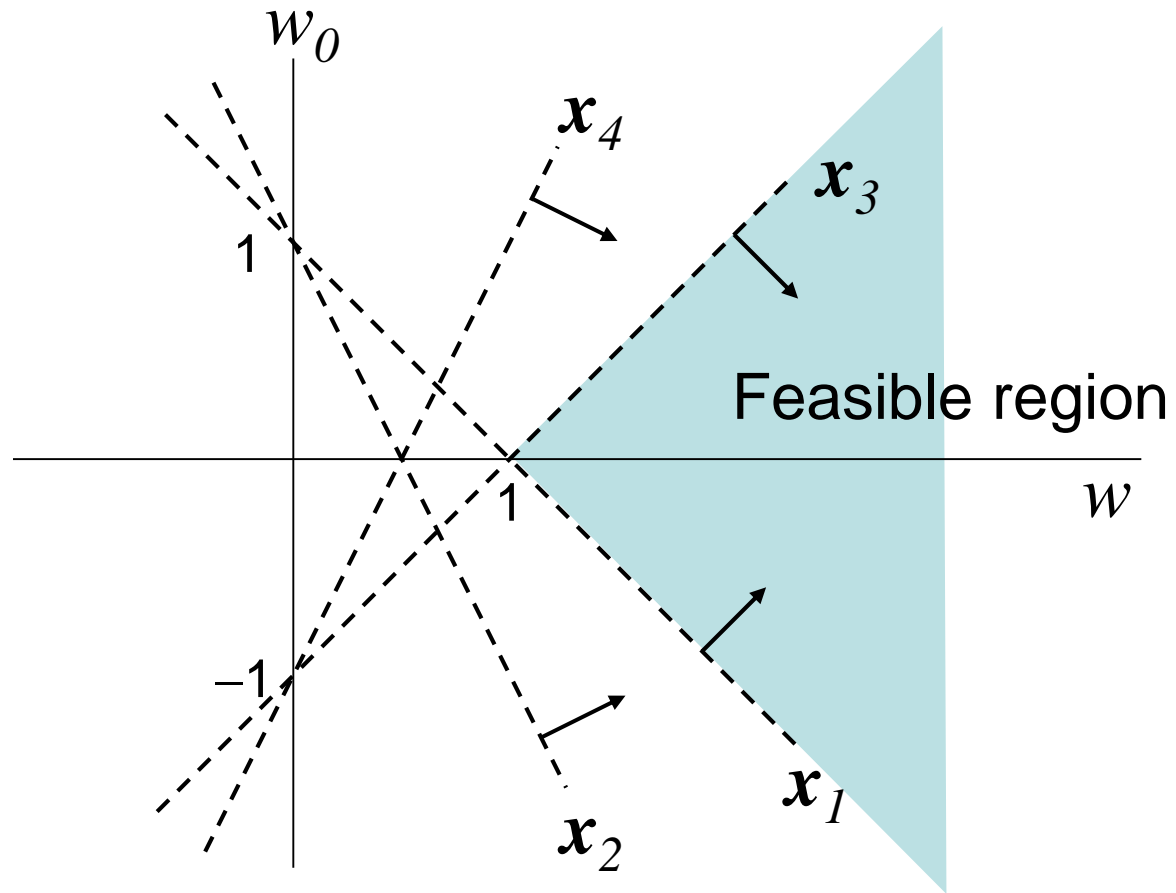
The constraints:

$$x_1 : (w + w_0) \geq 1$$

$$x_2 : (2w + w_0) \geq 1$$

$$x_3 : -(-w + w_0) \geq 1$$

$$x_4 : -(-2w + w_0) \geq 1$$



# Finding the Decision Boundary

Using Lagrange multipliers to include the constraints:

$$L = \frac{1}{2} \| \mathbf{w} \|^2 - \sum_{i=1}^N \lambda_i \left[ y_i (\mathbf{w}^T \mathbf{x}_i + w_0) - 1 \right]$$

Minimizing  $L$  relative to  $\mathbf{w}$  and  $w_0$  by setting the partial derivatives to zero results in

$$\mathbf{w} = \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i \quad \text{and} \quad \sum_{i=1}^N \lambda_i y_i = 0$$

These are the solutions for minimizing  $J$  with the constraints. The problem now is to determine the Lagrange multipliers.

# Finding the Decision Boundary

We also get the following regarding the solution from the Karush-Kuhn-Tucker (KKT) conditions:

$$\lambda_i \left[ y_i (\mathbf{w}^T \mathbf{x}_i + w_0) - 1 \right] = 0, \quad \forall i \quad \text{and} \quad \lambda_i \geq 0, \quad \forall i$$

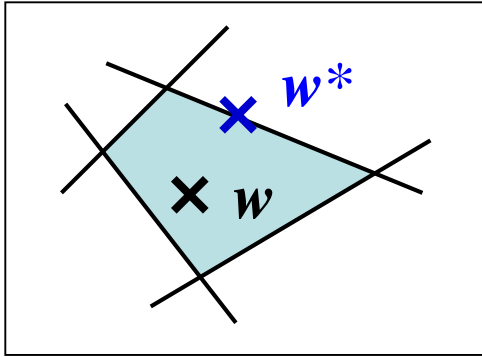
A Lagrange multiplier  $\lambda_i$  can be non-zero only if

$$y_i (\mathbf{w}^T \mathbf{x}_i + w_0) = 1$$

that is, when  $\mathbf{x}_i$  is a support vector.

We can understand this using the concept of active and inactive constraints. The Lagrange multipliers are zero for the inactive constraints.

# Finding the Decision Boundary



To understand the requirement that all  $\lambda_i \geq 0$ , consider the case when the solution (called  $\mathbf{w}^*$ ) is on the boundary given by the  $i^{\text{th}}$  constraint, and  $\mathbf{w}$  is a point inside the feasible region:

$$(\partial L / \partial \mathbf{w})_{\mathbf{w}^*} = 0 = (\partial J / \partial \mathbf{w})_{\mathbf{w}^*} - \lambda_i y_i \mathbf{x}_i$$

$$(J(\mathbf{w}) > J(\mathbf{w}^*)) \quad (\mathbf{w} - \mathbf{w}^*)^T ((\partial J / \partial \mathbf{w})_{\mathbf{w}^*}) = \lambda_i y_i (\mathbf{w} - \mathbf{w}^*)^T \mathbf{x}_i \geq 0$$

$$\text{(KKT) assuming } \lambda_i \neq 0 \quad y_i (\mathbf{w}^{*T} \mathbf{x}_i + w_0) = 1$$

$$\text{original constraint} \quad \Rightarrow \lambda_i \left[ y_i (\mathbf{w}^T \mathbf{x}_i + w_0) - 1 \right] \geq 0 \quad \Rightarrow \lambda_i \geq 0$$

This is still true when  $\mathbf{w}^*$  falls on multiple constraint boundaries.



# Finding the Decision Boundary

Now how do we obtain the Lagrange multipliers  $\lambda_i$ ? After incorporating expressions for optimal  $\mathbf{w}$ , the problem now is to maximize

$$L = \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \quad (\text{Now } L \text{ is a concave function with respect to } \lambda.)$$

Subject to the constraints 
$$\sum_{i=1}^N \lambda_i y_i = 0 \quad \text{and} \quad \lambda_i \geq 0, \quad \forall i$$

This is easier to solve than the original optimization problem, as now we have equality instead of inequality constraints. Once the optimal Lagrange Multipliers are obtained through some optimization procedure, we can compute  $\mathbf{w}$  using

$$\mathbf{w} = \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i$$

Using any support vector  $\mathbf{x}_i$ ,  $w_0$  can be computed using 
$$w_0 = y_i^{-1} - \mathbf{w}^T \mathbf{x}_i$$

# Example Continued

Continued example:  $x_1=1, x_2=2$  in  $\omega_1$ , and  $x_3=-1, x_4=-2$  in  $\omega_2$ .

The process of determining the Lagrange multipliers is not trivial in general. Based on what we know about the problem, we will set  $\lambda_2=\lambda_4=0$  first. We then use

$$\frac{\partial L}{\partial \lambda_1} = \frac{\partial L}{\partial \lambda_3} = 0$$

This is combined with the constraint  $\sum_{i=1}^N \lambda_i y_i = 0$

to yield the solution:  $\lambda_1=\lambda_3=0.5$ .

We can then determine  $w$  and  $w_0$ .

# Nonseparable Cases

In most cases, the two classes are not linearly separable. This means that the previous constraints ( $y_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1, \forall i$ ) are not valid.

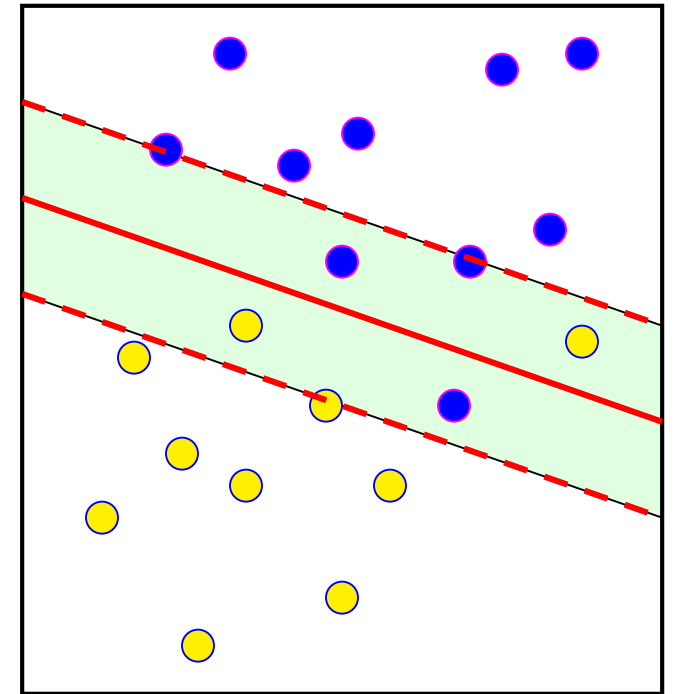
The constraints become

$$y_i(\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1 - \xi_i$$

where

$$\xi_i \geq 0$$

Cases with  $\xi_i > 0$  correspond to those  $\mathbf{x}_i$  that are misclassified or are within the "separation band".



# Nonseparable Cases

Now we have two goals:

- To minimize the number of  $x$  that are either misclassified or fall in the separation band.
- To maximize the margin.

This leads to the cost function

$$J = \frac{1}{2} \| \mathbf{w} \|^2 + C \sum_{i=1}^N I(\xi_i) \quad \text{where} \quad I(\xi_i) = \begin{cases} 1 & \text{for } \xi_i > 0 \\ 0 & \text{for } \xi_i = 0 \end{cases}$$

This is discontinuous and too difficult to optimize. Instead, we use

$$J = \frac{1}{2} \| \mathbf{w} \|^2 + C \sum_{i=1}^N \xi_i \quad (C > 0)$$

# Nonseparable Cases

Using Lagrange multipliers to include the constraints:

$$L = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \lambda_i \left[ y_i (\mathbf{w}^T \mathbf{x}_i + w_0) - 1 + \xi_i \right] - \sum_{i=1}^N \alpha_i \xi_i$$

Minimizing  $L$  relative to  $\mathbf{w}$ ,  $w_0$ , and  $\xi_i$  by setting the partial derivatives to zero results in two conditions as in linearly separable cases:

$$\mathbf{w} = \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i \quad \text{and} \quad \sum_{i=1}^N \lambda_i y_i = 0$$

As well as

$$\lambda_i = C - \alpha_i$$

# Nonseparable Cases

We also get the following from the KKT conditions:

$$\lambda_i \left[ y_i (\mathbf{w}^T \mathbf{x}_i + w_0) - 1 + \xi_i \right] = 0, \quad \forall i$$

$$\text{and} \quad \alpha_i \xi_i = 0, \quad \forall i$$

A Lagrange multiplier  $\lambda_i$  can be non-zero only if

$$y_i (\mathbf{w}^T \mathbf{x}_i + w_0) = 1 - \xi_i$$

that is, when  $\mathbf{x}_i$  is a support vector (inside or on the boundaries of the separation band, or misclassified).

In addition, each support vector  $\mathbf{x}_i$  that is not on the boundaries of the separation band has  $\lambda_i = C$ .

# Nonseparable Cases

Using the same techniques of the linearly separable cases, we end up with a very similar optimization problem: To maximize

$$L = \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

Subject to the constraints

$$\sum_{i=1}^N \lambda_i y_i = 0 \quad \text{and} \quad 0 \leq \lambda_i \leq C, \quad \forall i$$

Once we've found the optimal  $\lambda_i$ , we can compute  $\mathbf{w}$  using

$$\mathbf{w} = \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i$$

The value  $w_0$  can be computed using  $w_0 = y_i^{-1} - \mathbf{w}^T \mathbf{x}_i$

where  $\mathbf{x}_i$  is a support vector on the boundary of the separation band, i.e.,  $\xi_i=0$ .

# Nonseparable Cases

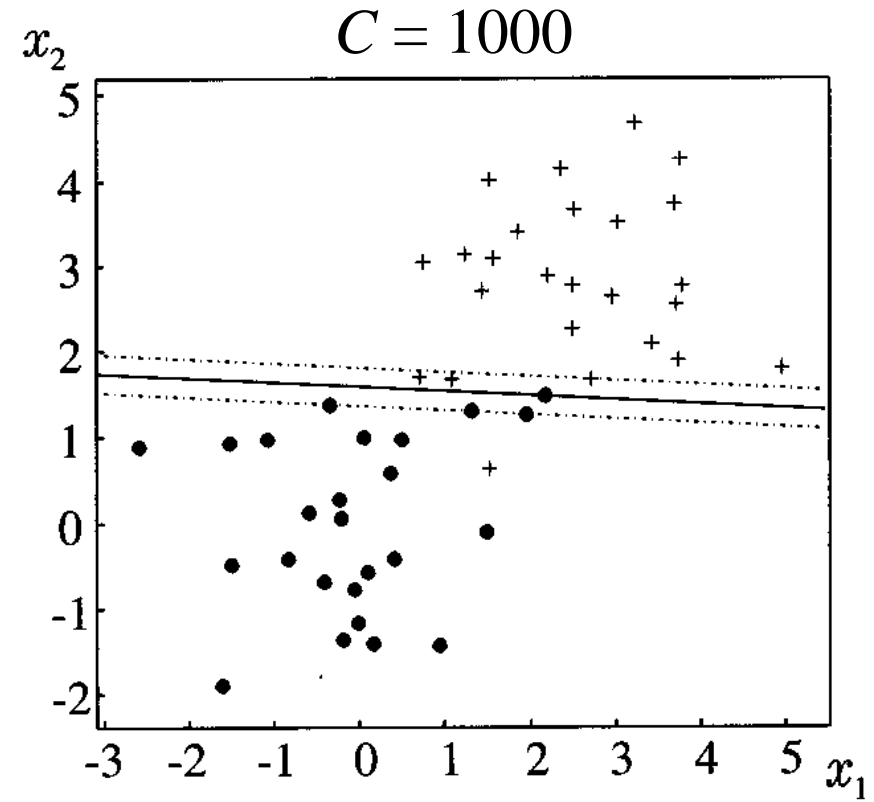
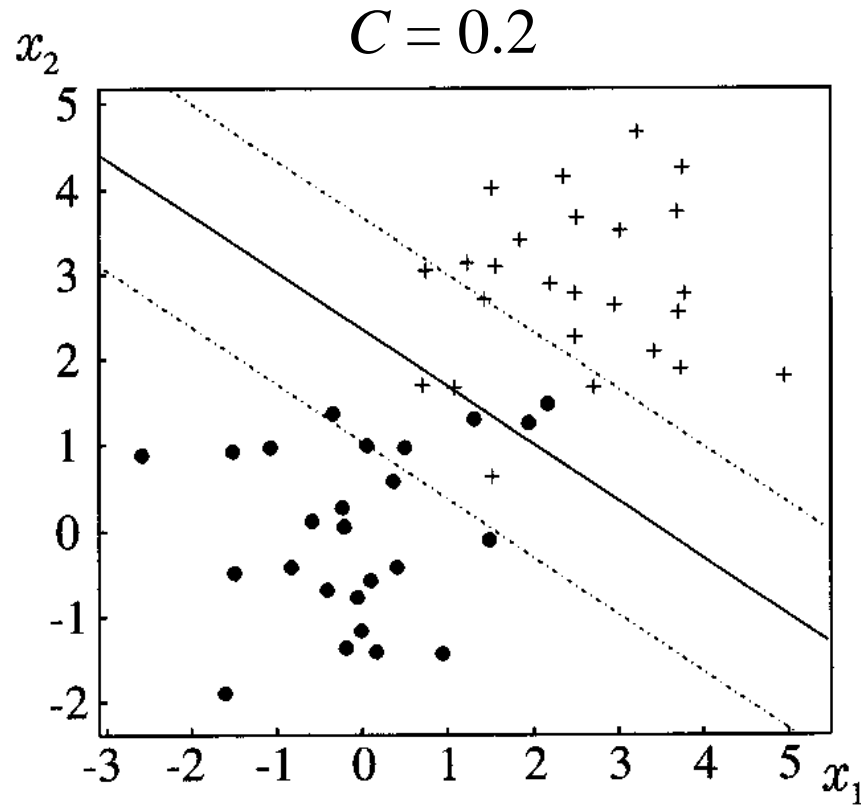
Some points about the parameter  $C$ :

- A trade-off between large margin and minimum classification error (in training data)
- Has to be selected separately (before SVM training). One common approach is to use cross-validation.
- For linearly separable cases, we can obtain the same results as obtained previously by setting  $C \rightarrow \infty$ .



# Nonseparable Cases

A computer experiment result for using different  $C$  values in a nonseparable case:



# Nonlinear SVM

In SVM, the input vectors are used only in inner products. This includes  $(\mathbf{x}_i^T \mathbf{x}_j)$  in  $L$ , and in the discriminant function:

$$L = \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i^T \mathbf{x} + w_0$$

This allows a "kernel trick": The inner-product of the high-dimensional mapping of two vectors can be expressed as a function of these two vectors.

# Kernel Trick

- Some methods use explicit mappings of data into high-dimensional spaces. There are several challenges:
  - The need to choose the dimensionality and the mapping functions.
  - Increased computational cost.
- For many problems where the problem can be expressed using inner products of the data points (feature vectors), we do not need to do explicit mapping.
  - Use kernel functions in the original space to compute the inner products of high-dimensional vectors.
  - No need to select the dimensionality and compute the mapped high-dimensional vectors.
  - Support vector machines are the most well-known example.

# Mercer's Theorem

Given any mapping  $x \rightarrow \varphi(x) \in H$  ,  
there exists a symmetric  $K(x,x')$  that satisfies (A) and (B)

Given any symmetric  $K(x,x')$  that satisfies (A),  
there exists a mapping  $x \rightarrow \varphi(x) \in H$  that satisfies (B)

$$(A) \quad \iint K(x, x') f(x) f(x') dx dx' \geq 0$$

for any  $f(x)$  that satisfies  $\int f(x)^2 dx < +\infty$

$$(B) \quad K(x, x') = \sum_r \varphi_r(x) \varphi_r(x')$$

meaning  $K(x,x')$  is the inner product of  $\varphi(x)$  and  $\varphi(x')$ .

# SVM Tasks

The optimization task, with  $\mathbf{x}$  now mapped to high-dimensional space, becomes the maximization of

$$\begin{aligned} L &= \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \boldsymbol{\varphi}(\mathbf{x}_i)^T \boldsymbol{\varphi}(\mathbf{x}_j) \\ &= \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) \end{aligned}$$

Subject to the constraints

$$\sum_{i=1}^N \lambda_i y_i = 0 \quad \text{and} \quad 0 \leq \lambda_i \leq C, \quad \forall i$$

# SVM Tasks

The linear discriminant function becomes

$$\begin{aligned} g(\mathbf{x}) &= \mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}) + w_0 = \sum_{i=1}^N \lambda_i y_i \boldsymbol{\varphi}(\mathbf{x}_i)^T \boldsymbol{\varphi}(\mathbf{x}) + w_0 \\ &= \sum_{i=1}^N \lambda_i y_i K(\mathbf{x}_i, \mathbf{x}) + w_0 \end{aligned}$$

As a result, we can select the function  $K$  (the inner-product kernel), and then do all the computation without calculating (or even knowing)  $\boldsymbol{\varphi}(\mathbf{x})$ !!

This "kernel trick" avoids the increased computational cost involved in working with high-dimensional data.

# Types of Inner-Product Kernels

Different common types of  $K(\mathbf{x}, \mathbf{x}')$ :

## Polynomial

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + 1)^q \quad (q > 0)$$

- ◆  $q$  is specified by the user

## RBF

$$K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right)$$

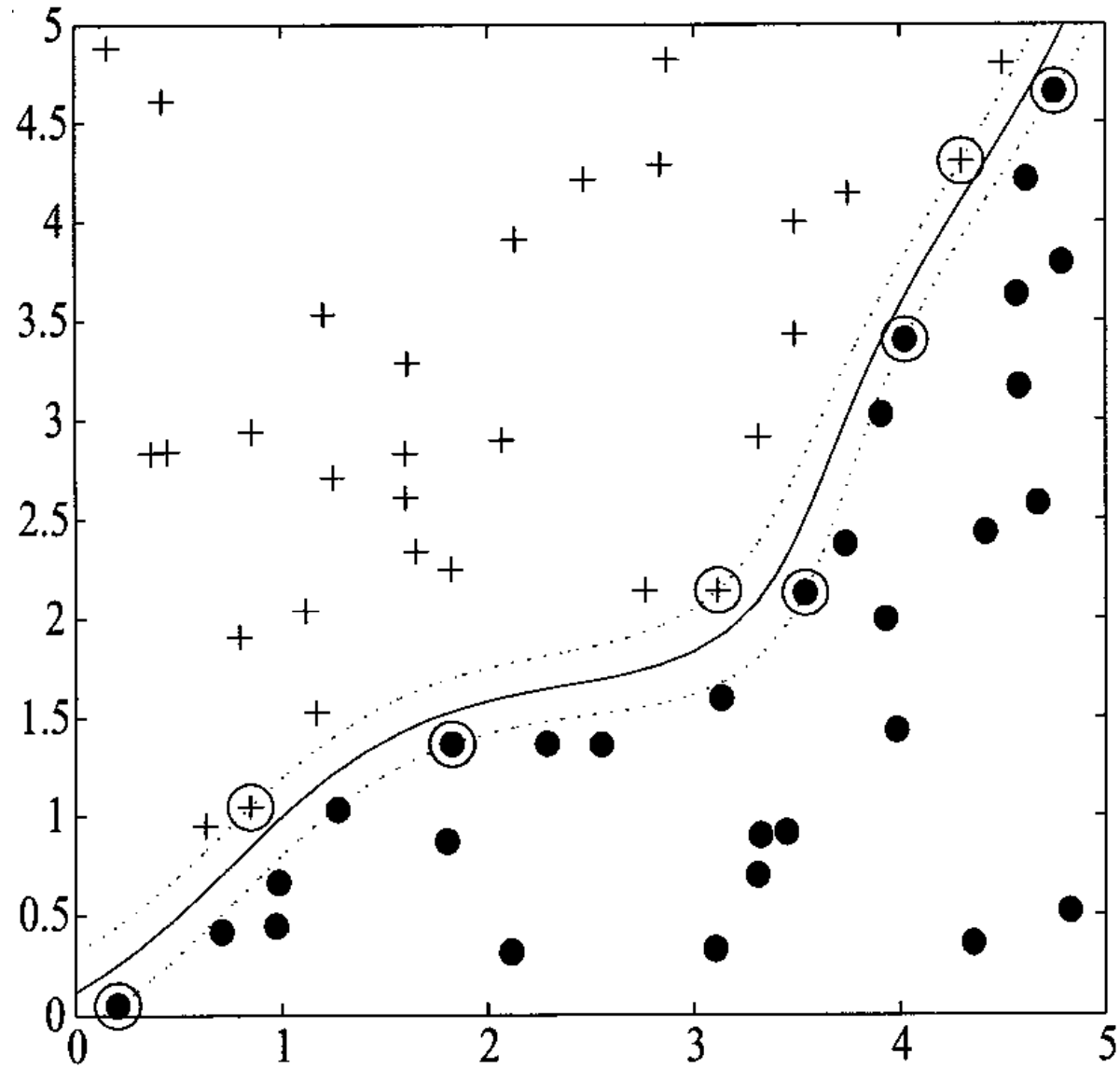
- ◆  $\sigma$  is specified by the user and common to all kernels
- ◆ RBF network

## Hyperbolic

$$K(\mathbf{x}, \mathbf{x}') = \tanh(\beta \mathbf{x}^T \mathbf{x}' + \gamma)$$

- ◆ Only some values of  $\beta$  and  $\gamma$  can be used, e.g.,  $\beta=2$  and  $\gamma=1$ .
- ◆ Two-layer perceptron with tanh activation functions.

# SVM Experiments





# Kernel PCA

- When PCA is not effective for the original data distribution, mapping to a high-dimensional space might lead to more useful PCA results to subsequent tasks.
- Mapping:  $\mathbf{x} \leftarrow \varphi(\mathbf{x})$
- Assuming that the projection of data points to the high-dimensional space is centered (having zero mean), the covariance matrix becomes:

$$\Sigma = \frac{1}{N} \sum_{i=1}^N \varphi(\mathbf{x}_i) \varphi(\mathbf{x}_i)^T$$

- However, we never actually compute the high-dimensional vectors, so this covariance matrix cannot be used directly.

# Derivation of Kernel PCA

- As PCA involves the eigen-decomposition of the covariance matrix, let us consider this form ( $\mathbf{u}$  is an eigenvector here):

$$\Sigma \mathbf{u} = \lambda \mathbf{u}$$

- We can express  $\mathbf{u}$  as a linear combination of the projected high-dimensional data points, as the subspace spanned by the eigenvectors with nonzero eigenvalues is the same as the subspace spanned by the data points:

$$\mathbf{u} = \sum_{i=1}^N a_i \varphi(\mathbf{x}_i)$$

# Derivation of Kernel PCA

- Substituting  $u$  back to the eigen-decomposition of  $\Sigma$ :

$$\Rightarrow \left[ \frac{1}{N} \sum_{i=1}^N \varphi(\mathbf{x}_i) \varphi(\mathbf{x}_i)^T \right] \left[ \sum_{j=1}^N a_j \varphi(\mathbf{x}_j) \right] = \lambda \left[ \sum_{i=1}^N a_i \varphi(\mathbf{x}_i) \right]$$

- Exchange of terms:

$$\Rightarrow \frac{1}{N} \sum_{i=1}^N \varphi(\mathbf{x}_i) \left[ \sum_{j=1}^N a_j \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j) \right] = \lambda \sum_{i=1}^N a_i \varphi(\mathbf{x}_i)$$

- To convert everything involving the high-dimensional vectors to inner products:

$$\Rightarrow \frac{1}{N} \sum_{i=1}^N \varphi(\mathbf{x}_k)^T \varphi(\mathbf{x}_i) \left[ \sum_{j=1}^N a_j \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j) \right] = \lambda \sum_{i=1}^N a_i \varphi(\mathbf{x}_k)^T \varphi(\mathbf{x}_i)$$

# Derivation of Kernel PCA

- Now, express this with the inner-product kernels:

$$\Rightarrow \frac{1}{N} \sum_{i=1}^N K(\mathbf{x}_k, \mathbf{x}_i) \left[ \sum_{j=1}^N a_j K(\mathbf{x}_j, \mathbf{x}_i) \right] = \lambda \sum_{i=1}^N a_i K(\mathbf{x}_k, \mathbf{x}_i)$$

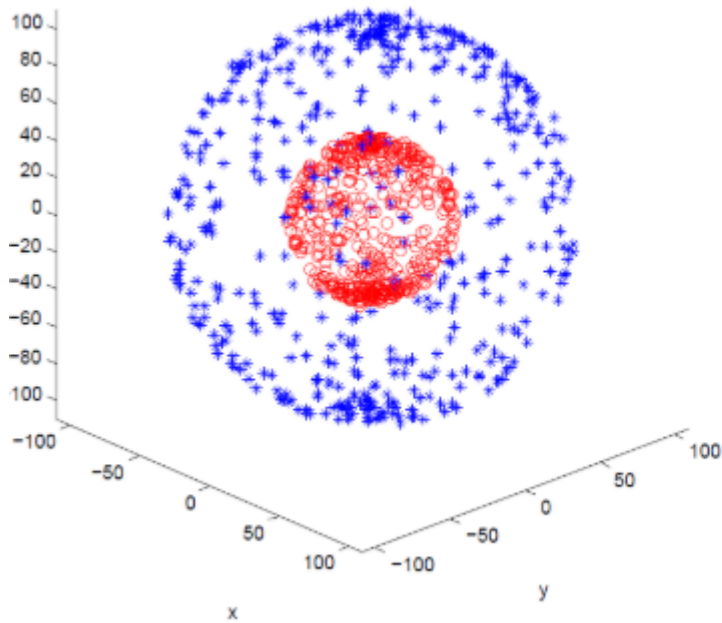
- Express this in matrix form:  $\mathbf{K}(\mathbf{K}\boldsymbol{\alpha}) = N\lambda\mathbf{K}\boldsymbol{\alpha}$
- Here  $\mathbf{K}$  is the kernel matrix, and  $\boldsymbol{\alpha}$  is the vector of the coefficients. We can remove one  $\mathbf{K}$  from both sides because we are not concerned with eigenvectors of zero eigenvalues:

$$\mathbf{K}\boldsymbol{\alpha} = N\lambda\boldsymbol{\alpha}$$

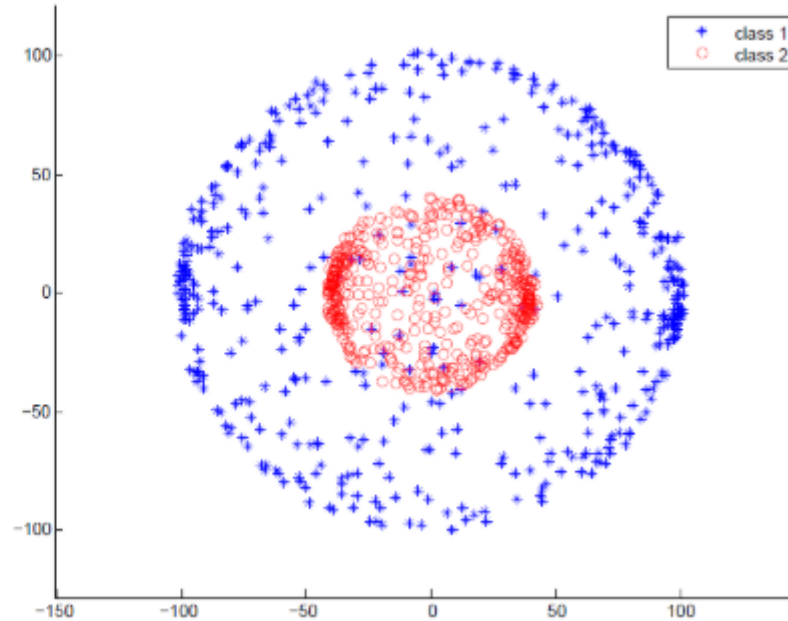
- Now we again have an eigen-decomposition problem. We can then project the data points to a low-dimension subspace spanned by a subset of the eigenvectors (with the largest eigenvalues) of  $\mathbf{K}$ .

# Kernel PCA Example

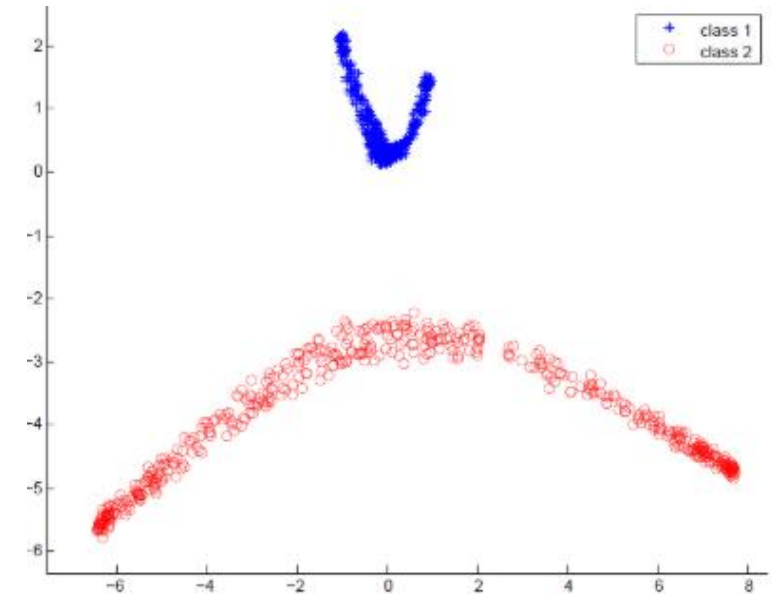
Original Data (3D)



Linear PCA (2D)

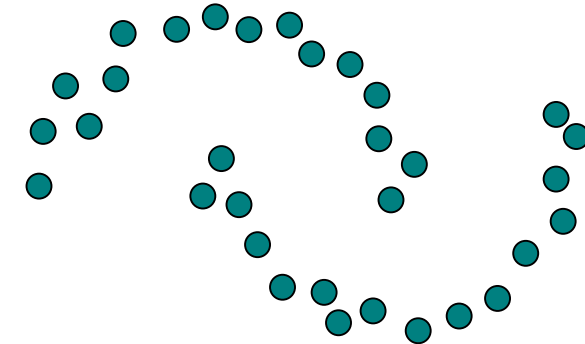
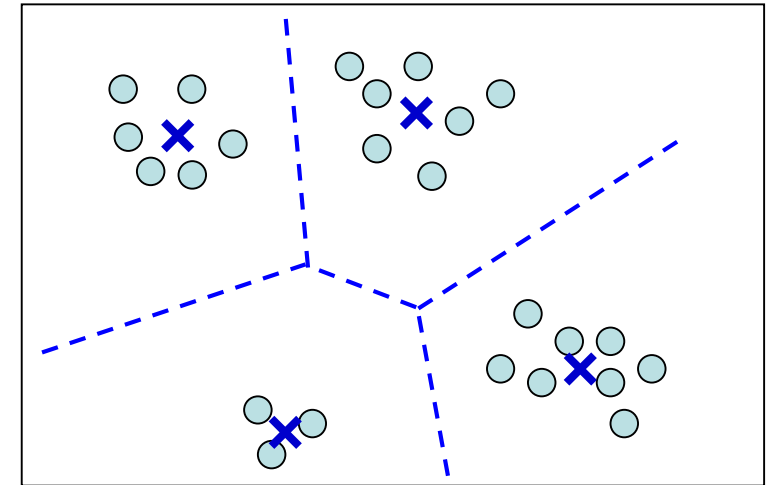


Kernel PCA (2D)  
with Gaussian Kernels



# Kernel k-Means: Motivation

- The “kernel” version of k-means aims to solve another limitation of k-means: Linear cluster boundaries.
  - In standard k-means, the prototypes define “cluster boundaries” that are hyperplanes in the feature space:
- Q: Can we separate the "clusters" in this example?



# Kernel Trick in Kernel k-Means

- Let us rewrite the cost function of k-means:

$$J = \sum_{j=1}^k \sum_{\mathbf{x} \in C_j} \|\mathbf{x} - \mathbf{m}_j\|^2$$

- $C_j$  : The set of samples in cluster #j.
  - $\mathbf{m}_j$  : Centroid (prototype) of cluster #j.
- Now after the mapping to a high-dimensional space:

$$\mathbf{x} \leftarrow \varphi(\mathbf{x})$$

$$\mathbf{m}_j \leftarrow \frac{1}{|C_j|} \sum_{\mathbf{x} \in C_j} \varphi(\mathbf{x})$$

# Kernel Trick in Kernel k-Means

The squared distance in the cost function can be rewritten as

$$\varphi(\mathbf{x})^T \varphi(\mathbf{x}) - \frac{2}{|C|} \sum_{\mathbf{x}_j \in C} \varphi(\mathbf{x})^T \varphi(\mathbf{x}_j) + \frac{1}{|C|^2} \sum_{\mathbf{x}_j, \mathbf{x}_k \in C} \varphi(\mathbf{x}_j)^T \varphi(\mathbf{x}_k)$$

assuming that  $\mathbf{x}$  belongs to cluster  $C$ .

We can see that everything is expressed in inner products of mapped vectors. These inner products can be replaced with the corresponding kernel function:

$$K(\mathbf{x}, \mathbf{x}) - \frac{2}{|C|} \sum_{\mathbf{x}_j \in C} K(\mathbf{x}, \mathbf{x}_j) + \frac{1}{|C|^2} \sum_{\mathbf{x}_j, \mathbf{x}_k \in C} K(\mathbf{x}_j, \mathbf{x}_k)$$



# Kernel k-Means

- Iterative update of cluster assignments only: In each iteration, each sample is assigned to the cluster with the smallest sample-cluster distance.
- No explicit computation of the prototypes. (They are in the high dimensional space.) → Not suitable if the actual prototypes are needed.
- When  $N$  is very large, processing of the  $N \times N$  kernel matrix can still be computationally expensive.

