

# IFI 9000 Analytics Methods

## Support Vector Machines

by **Houping Xiao**

Spring 2021



# Introduction

# What is a separating hyperplane?

- It is a flat affine subspace of dimension  $p - 1$  in the  $p$  dimensional space
- It takes the following equation

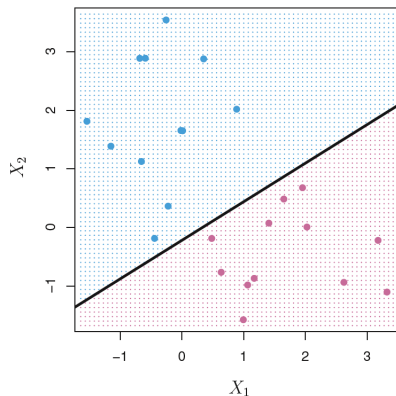
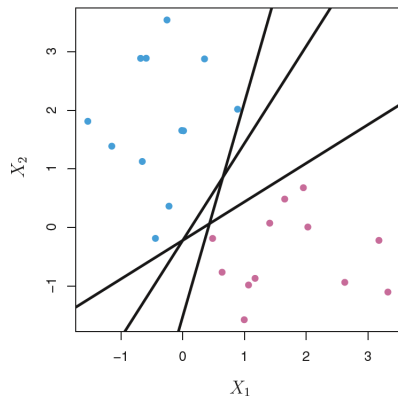
$$\beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p = 0.$$

- The vector  $(\beta_1, \cdots, \beta_p)$  is the normal to the plane
- The sign of

$$f(\mathbf{x}) = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p$$

can be used as a class indicator for classification problems

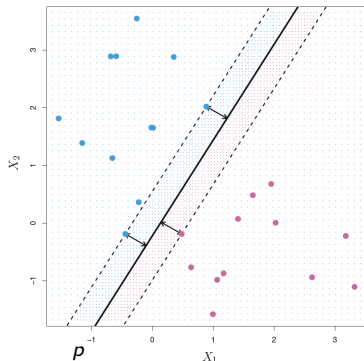
# What is a separating hyperplane?



- Coding the two classes by color, we have  $f(\mathbf{x}) > 0$  for one color and  $f(\mathbf{x}) \leq 0$  for the other class. If we assign  $y_i = 1$  or  $y_i = -1$  to the classes, then a separating hyperplane will satisfy  $y_i f(\mathbf{x}_i) > 0$ , and the decision boundary is identified by  $f(\mathbf{x}) = 0$ .

# Maximum Margin Classifier

- In selection of all possible separating hyperplanes, picking the ones that makes the **maximum symmetric gap** between the two classes (**maximum margin**)



$$\max_{\beta_0, \dots, \beta_p} M \quad \text{s.t.} \quad \sum_{j=1}^p \beta_j^2 = 1,$$

$$y_i(\beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p}) \geq M, \quad i = 1, \dots, N$$

# Convex Formulation of Maximum Margin Classifier

- The original problem

$$\max_{\beta_0, \dots, \beta_p} M \quad \text{s.t.} \quad \sum_{j=1}^p \beta_j^2 = 1,$$

$$y_i(\beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p}) \geq M, \quad i = 1, \dots, N$$

can be written in vector notation as

$$\max_{\beta_0, \beta} M \quad \text{s.t.} \quad \|\beta\|^2 = 1,$$

$$y_i(\beta_0 + \beta^\top \mathbf{x}_i) \geq M, \quad i = 1, \dots, N$$

- To write the program in convex form, we can replace it with

$$\max_{\beta_0, \beta} M \quad \text{s.t.} \quad \frac{y_i}{\|\beta\|}(\beta_0 + \beta^\top \mathbf{x}_i) \geq M, \quad i = 1, \dots, N$$

# Convex Formulation of Maximum Margin Classifier

$$\max_{\beta_0, \beta} M \quad \text{s.t.} \quad \frac{y_i}{\|\beta\|} (\beta_0 + \beta^\top \mathbf{x}_i) \geq M, \quad i = 1, \dots, N$$

in turn yields

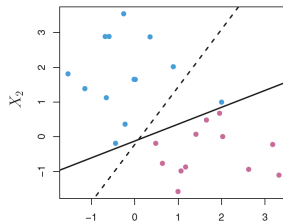
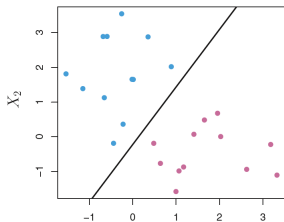
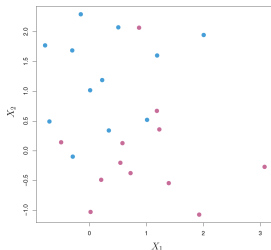
$$\max_{\beta_0, \beta} M \quad \text{s.t.} \quad y_i (\beta_0 + \beta^\top \mathbf{x}_i) \geq M \|\beta\|, \quad i = 1, \dots, N$$

- We have a scaling ambiguity (if  $\beta_0, \beta$  are solutions, so are their scaled forms) we can pick  $\|\beta\| = \frac{1}{M}$ , which reduces the optimization to the convex form

$$\max_{\beta_0, \beta} \|\beta\| \quad \text{s.t.} \quad y_i (\beta_0 + \beta^\top \mathbf{x}_i) \geq 1, \quad i = 1, \dots, N$$

# Non-ideal Cases

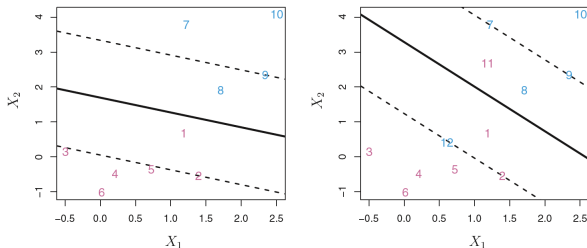
- In the majority of cases, we have non-separable or noisy data issues (here the model drastically changes due to the addition of a single point)





# Support Vector Classifier

- Support vector classifier deals with the non-ideal issues by allowing some misclassification



- We find a **soft margin** that would include some misclassified points

$$\max_{\beta_0, \beta} M \quad \text{s.t.} \quad \|\beta\|^2 = 1$$

$$y_i(\beta_0 + \beta^T \mathbf{x}_i) \geq M(1 - \epsilon_i), \quad i = 1, \dots, N$$

$$\epsilon \geq 0, \quad \sum_{i=1}^N \epsilon_i \leq C$$

# Support Vector Classifier

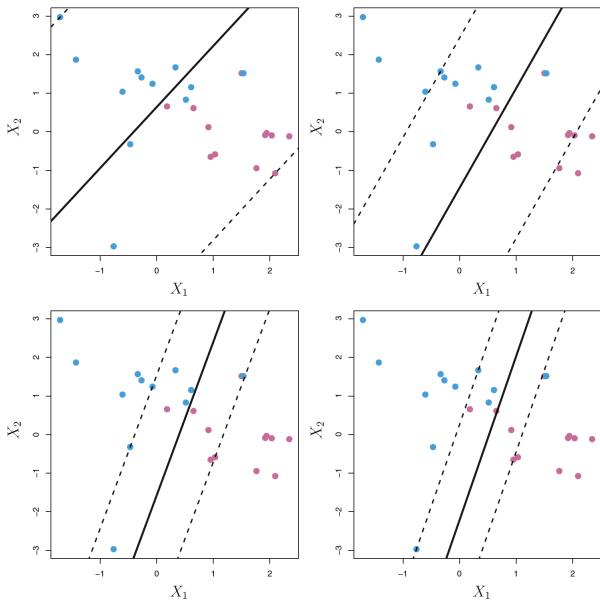
$$\max_{\beta_0, \beta, \epsilon_1, \dots, \epsilon_N, M} M \quad \text{s.t.} \quad \|\beta\|^2 = 1$$

$$y_i(\beta_0 + \beta^\top \mathbf{x}_i) \geq M(1 - \epsilon_i), \quad i = 1, \dots, N$$

$$\epsilon_i \geq 0, \quad \sum_{i=1}^N \epsilon_i \leq C$$

The slack variables  $\epsilon_1, \dots, \epsilon_N$  would allow individual observations to be on the wrong side of the margin. The parameter  $C$  controls the width of the soft margin, larger  $C$  increases the tolerance for more misclassified points.

# Support Vector Classifier



# Convex Formulation

- If we follow a similar procedure as before we end up with the following convex formulation

$$\max_{\beta_0, \boldsymbol{\beta}, \epsilon_1, \dots, \epsilon_N} \|\boldsymbol{\beta}\| \quad s.t. \quad y_i(\beta_0 + \boldsymbol{\beta}^\top \mathbf{x}_i) \geq 1 - \epsilon_i, i = 1, \dots, N$$

$$\epsilon_i \geq 0, \sum_{i=1}^N \epsilon_i \leq C$$

or in a regularized form as

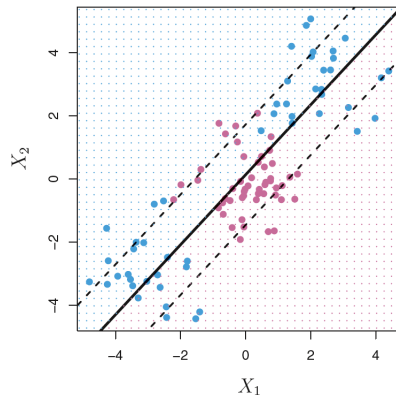
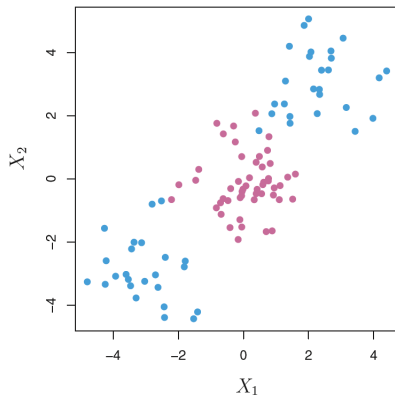
$$\max_{\beta_0, \boldsymbol{\beta}, \epsilon_1, \dots, \epsilon_N} \|\boldsymbol{\beta}\| + \lambda \sum_{i=1}^N \epsilon_i \quad s.t. \quad y_i(\beta_0 + \boldsymbol{\beta}^\top \mathbf{x}_i) \geq 1 - \epsilon_i, i = 1, \dots, N$$

$$\epsilon_i \geq 0, i = 1, \dots, N$$

Try some coding demos!

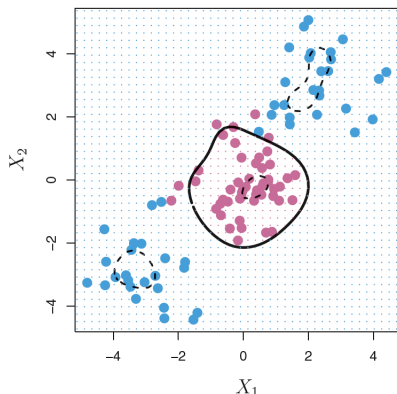
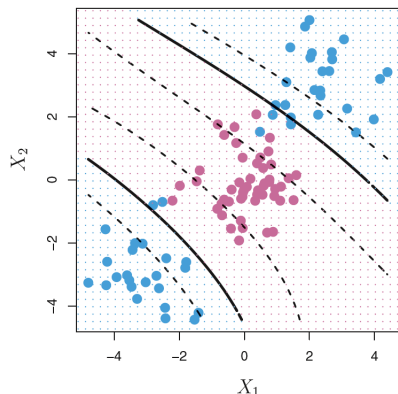
# Generalization to Nonlinear Boundary

- Different classes cannot always be separated via a separating hyperplane



# Feature Extension

- One remedy to handle nonlinear decision boundaries is adding nonlinear features to the variable list, e.g.,  $x_1, x_2, x_1^2, x_2^2, x_1x_2$ , etc



- Features:  $x_1, x_2, x_1^2, x_2^2, x_1x_2, x_1^3, x_2^3, x_1^2x_2, x_1x_2^2$

- Use of kernel is a more systematic way of introducing nonlinearities in SVM
- Generalization to kernels is easier to apply when working with the dual of the SVM convex program, however we skip such detail here
- Instead we look into an alternative way of representing the hyperplanes

$$f(\mathbf{x}) = \beta_0 + \boldsymbol{\beta}^\top \mathbf{x}$$

or to incorporate all data points we can rewrite it as

$$f(\mathbf{x}) = \beta_0 + \sum_{i=1}^N \alpha_i \mathbf{x}_i^\top \mathbf{x}$$

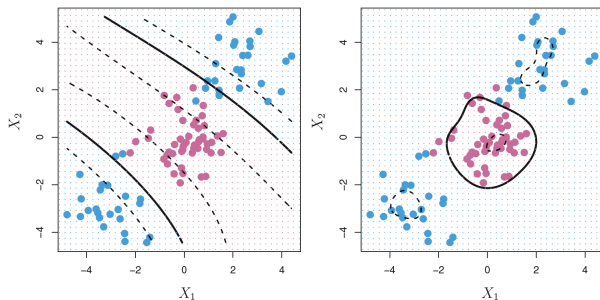
- In kernel methods, we replace  $\mathbf{x}_i^\top \mathbf{x}$  with some (symmetric) function  $k(\mathbf{x}_i, \mathbf{x})$ :

$$f(\mathbf{x}) = \beta_0 + \sum_{i=1}^N \alpha_i \cdot k(\mathbf{x}_i, \mathbf{x})$$

# Radial Kernels

- Radial kernels are among the most robust and promising kernels used in SVM applications

$$k(\mathbf{u}, \mathbf{v}) = \exp \left( -\gamma \sum_{j=1}^p (u_j - v_j)^2 \right)$$





# SVM: More Than Two Classes

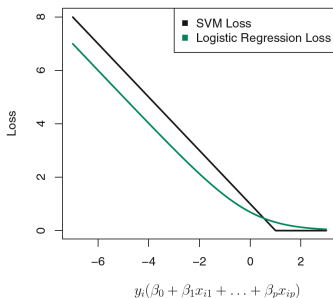
Here are two approaches we can use to classify  $K > 2$  classes

- **One vs All (OVA)**: Take one class, take remaining as the alternative class, repeat this for all  $k$  classes to get the classifiers  $f_k(\mathbf{x})$ . As the prediction for a sample  $\mathbf{x}_t$  pick the class with largest  $f_k(\mathbf{x}_t)$
- **One vs One (OVO)**: Apply SVM to all combinations of two from the  $K$  classes, for each test point  $\mathbf{x}_t$  count the number of times it lands on each class, then pick the max as the predicted class

OVO is usually more reliable, however for large number of classes OVA is computationally more desirable

# SVM vs Logistic Regression

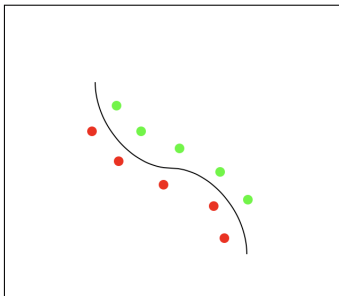
- The loss functions have some similarities when plotted as a function of  $y_i \beta^\top \mathbf{x}_i$



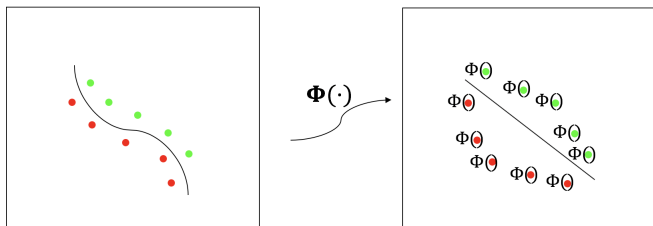
- LR becomes ill-posed when the classes are (almost) separable, LDA and SVM are fine
- LR would provide class probabilities
- SVMs easily generalize to nonlinear boundaries, LR and LDA are computationally harder

# Why Kernel Methods are Needed?

- When the relations between dependent and independent variables are **nonlinear** in the original space, i.e., none hyperplane can be found to separate the data for classification problem in SVM



# Why Kernel Methods are Needed?



- Kernel methods allow us to map the data (no separable in the original space) into a new space (possibly high-dimensional space) where data are linear separable
- Commonly used mapping:
  - $x \longrightarrow \{x, x^2\}$
  - $\{x_1, x_2\} \longrightarrow \{x_1^2, \sqrt{2}x_1x_2, x_2^2\}$

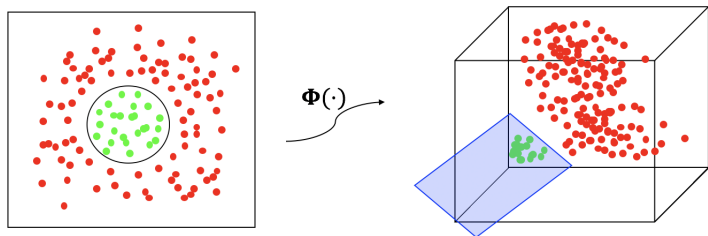
- **Computationally difficulty:** More features are in demand to capture the nonlinearity (i.e., higher dimensional representation)
- **Kernel functions:** Given  $\mathbf{x}_i, \mathbf{x}_j$  in the original space  $\mathbb{R}$ ,  $\phi$  is a kernel function if

$$k(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$

- $\phi(\mathbf{x}_i)$  is the coordinate representation in new space.
- The kernel function is a similarity measure; not need to calculate the  $\phi(\mathbf{x}_i)$  and  $\phi(\mathbf{x}_j)$  individually. We are more interested of their inner product

# An Example of Kernel Methods

- $\phi : \mathbf{x} \in \mathbb{R}^2 \rightarrow \mathbb{R}^3$ ,  $\Phi(\mathbf{x}_1, \mathbf{x}_2) \rightarrow (x_1^2, \sqrt{2}x_1x_2, x_2^2)$
- Easy to check  $\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = \langle \mathbf{x}_i, \mathbf{x}_j \rangle^2 = f(\langle \mathbf{x}_i, \mathbf{x}_j \rangle) = k(\mathbf{x}_i, \mathbf{x}_j)$



- $\phi(\mathbf{x}_1, \mathbf{x}_2) \rightarrow (x_1^2, \sqrt{2}x_1x_2, x_2^2)$ . Left: input space, Right: feature space
- Different feature transformation can share the same kernel function
- $\psi(\mathbf{x}_1, \mathbf{x}_2) = (x_1^2, x_2^2, x_1x_2, x_2x_1)$ , then

$$\langle \psi(\mathbf{x}_i), \psi(\mathbf{x}_j) \rangle = \langle \mathbf{x}_i, \mathbf{x}_j \rangle^2 = k(\mathbf{x}_i, \mathbf{x}_j)$$

# Kernel Methods in Ridge Regression

- Ridge regression

$$\min_{\beta} \sum_{i=1}^N \|y_i - \beta^{\top} \mathbf{x}_i\|^2 + \lambda \sum_{i=1}^N \beta_i^2$$

- Prediction at  $\mathbf{x}'$ :

$$\begin{aligned} y^* &= \beta^* \mathbf{x}' \\ &= ((\lambda I_D + \mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{y})^{\top} \mathbf{x}' \\ &= \mathbf{y}^{\top} (\lambda I_N + \mathbf{X} \mathbf{X}^{\top})^{-1} \mathbf{X} \mathbf{x}' \end{aligned}$$

# Kernel Methods in Ridge Regression

- $y^* = \mathbf{y}^\top (\lambda I_N + \mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{X}\mathbf{x}'$

$\mathbf{X}\mathbf{X}^\top$

$$\mathbf{X}\mathbf{X}^\top = \begin{pmatrix} \langle \mathbf{x}_1, \mathbf{x}_1 \rangle & \langle \mathbf{x}_1, \mathbf{x}_2 \rangle & \cdots & \langle \mathbf{x}_1, \mathbf{x}_N \rangle \\ \langle \mathbf{x}_2, \mathbf{x}_1 \rangle & \langle \mathbf{x}_2, \mathbf{x}_2 \rangle & \cdots & \langle \mathbf{x}_2, \mathbf{x}_N \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle \mathbf{x}_N, \mathbf{x}_1 \rangle & \langle \mathbf{x}_N, \mathbf{x}_2 \rangle & \cdots & \langle \mathbf{x}_N, \mathbf{x}_N \rangle \end{pmatrix}$$

$\mathbf{X}\mathbf{x}'$

$$\mathbf{X}\mathbf{x}' = \begin{pmatrix} \langle \mathbf{x}_1, \mathbf{x}' \rangle \\ \langle \mathbf{x}_2, \mathbf{x}' \rangle \\ \vdots \\ \langle \mathbf{x}_N, \mathbf{x}' \rangle \end{pmatrix}$$



# Generalizing to Non-Linear Regression

- **Ridge regression:** consider the feature transformation  $\phi = (\phi_1, \dots, \phi_P)$  and

$$\Phi = \begin{pmatrix} \phi_1(\mathbf{x}_1) & \phi_2(\mathbf{x}_1) & \cdots & \phi_P(\mathbf{x}_1) \\ \phi_1(\mathbf{x}_2) & \phi_2(\mathbf{x}_2) & \cdots & \phi_P(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_1(\mathbf{x}_N) & \phi_2(\mathbf{x}_N) & \cdots & \phi_P(\mathbf{x}_N) \end{pmatrix}$$

- Making a prediction at

$$y^* = \mathbf{y}^\top (\lambda I_N + \Phi \Phi^\top)^{-1} \Phi \phi(\mathbf{x}')$$

where  $(\Phi \Phi^\top)_{ij} = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$

# What Is A Kernel Trick?

- Ridge Regression

$$y^* = \mathbf{y}^\top \left( \lambda I_N + \mathbf{X}\mathbf{X}^\top \right)^{-1} \mathbf{X}\mathbf{x}'$$

- Replace inner product with a kernel function  $k(\mathbf{x}_i, \mathbf{x}_j)$
- Replace  $\mathbf{X}\mathbf{X}^\top$  with the Gram matrix  $\mathbf{K}$  where  $\mathbf{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$
- Kernel Regression

$$y^* = \mathbf{y}^\top (\lambda I_N + \mathbf{K})^{-1} k(\mathbf{X}, \mathbf{x}')$$

- We want to map the patterns into a high-dimensional feature space  $F$  and compare them using a dot product
- To avoid working in the  $F$ , choose a map function mapping to feature space where the dot product can be directly evaluated using a nonlinear function in the input space
- this is called the kernel trick

$$k(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$

- Recall the solution of dual problem:

$$\hat{\beta} = \sum_{i=1}^N \hat{\alpha}_i y_i \mathbf{x}_i$$

- Decision rule:

$$\hat{y} = \text{sign} \left( \sum_{i \in SV} \hat{\alpha}_i y_i (\mathbf{x}_i^T \mathbf{x}) + \beta_0 \right)$$

- Using kernel would give us:

$$\hat{y} = \text{sign} \left( \sum_{i \in SV} \hat{\alpha}_i y_i \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle + \beta_0 \right)$$

# How to Construct a Kernel?

- The simplest kernel function we have is that

$$k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$$

- We start with basis functions

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

- Direct design

- Measure similarity between  $\mathbf{x}_i$  and  $\mathbf{x}_j$
- Gram matrix must be positive semi-definite
- $k$  should be symmetric

- Creating more complicated kernels based on the known kernels  $k_1$  and  $k_2$

- $k(\mathbf{x}_i, \mathbf{x}_j) = ck_1(\mathbf{x}_i, \mathbf{x}_j)$
- $k(\mathbf{x}_i, \mathbf{x}_j) = f(\mathbf{x}_i)k_1(\mathbf{x}_i, \mathbf{x}_j)f(\mathbf{x}_j)$
- $k(\mathbf{x}_i, \mathbf{x}_j) = q(k_1(\mathbf{x}_i, \mathbf{x}_j))$ ,  $q$  is polynomial
- $k(\mathbf{x}_i, \mathbf{x}_j) = k_1(\mathbf{x}_i, \mathbf{x}_j) + k_2(\mathbf{x}_i, \mathbf{x}_j)$ ,  $k_1(\mathbf{x}_i, \mathbf{x}_j)k_2(\mathbf{x}_i, \mathbf{x}_j)$

# Popular Kernels

- Mercer Kernel:  $\mathbf{K}$  is positive definite
- Gaussian Kernel or Radial Basis Function (RBF)

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{x}_i - \mathbf{x}_j\|^2\right)$$

Mapping inputs to an infinite dimensional space

- Cosine Similarity

$$k(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{x}_i^\top \mathbf{x}_j}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|}$$

- Probabilistic Kernel Functions

$$k(\mathbf{x}_i, \mathbf{x}_j) = Pr(\mathbf{x}_i|\boldsymbol{\theta})Pr(\mathbf{x}_j|\boldsymbol{\theta})$$

Two inputs are more similar if both have high probabilities

- Bayesian Kernel

$$k(\mathbf{x}_i, \mathbf{x}_j) = \int Pr(\mathbf{x}_i|\boldsymbol{\theta})Pr(\mathbf{x}_j|\boldsymbol{\theta})Pr(\boldsymbol{\theta})d\boldsymbol{\theta}$$

- We can use kernel function to generate new features. Apply the kernel function for each input and a set of  $P$  centroids

$$\Phi(\mathbf{x}) = [k(\mathbf{x}, \boldsymbol{\mu}_1), k(\mathbf{x}, \boldsymbol{\mu}_1), \dots, k(\mathbf{x}, \boldsymbol{\mu}_2)]$$

- How to choose centroids  $\boldsymbol{\mu}_i$ 
  - Randomly selection
  - Clustering
  - Using all available data points

The End