



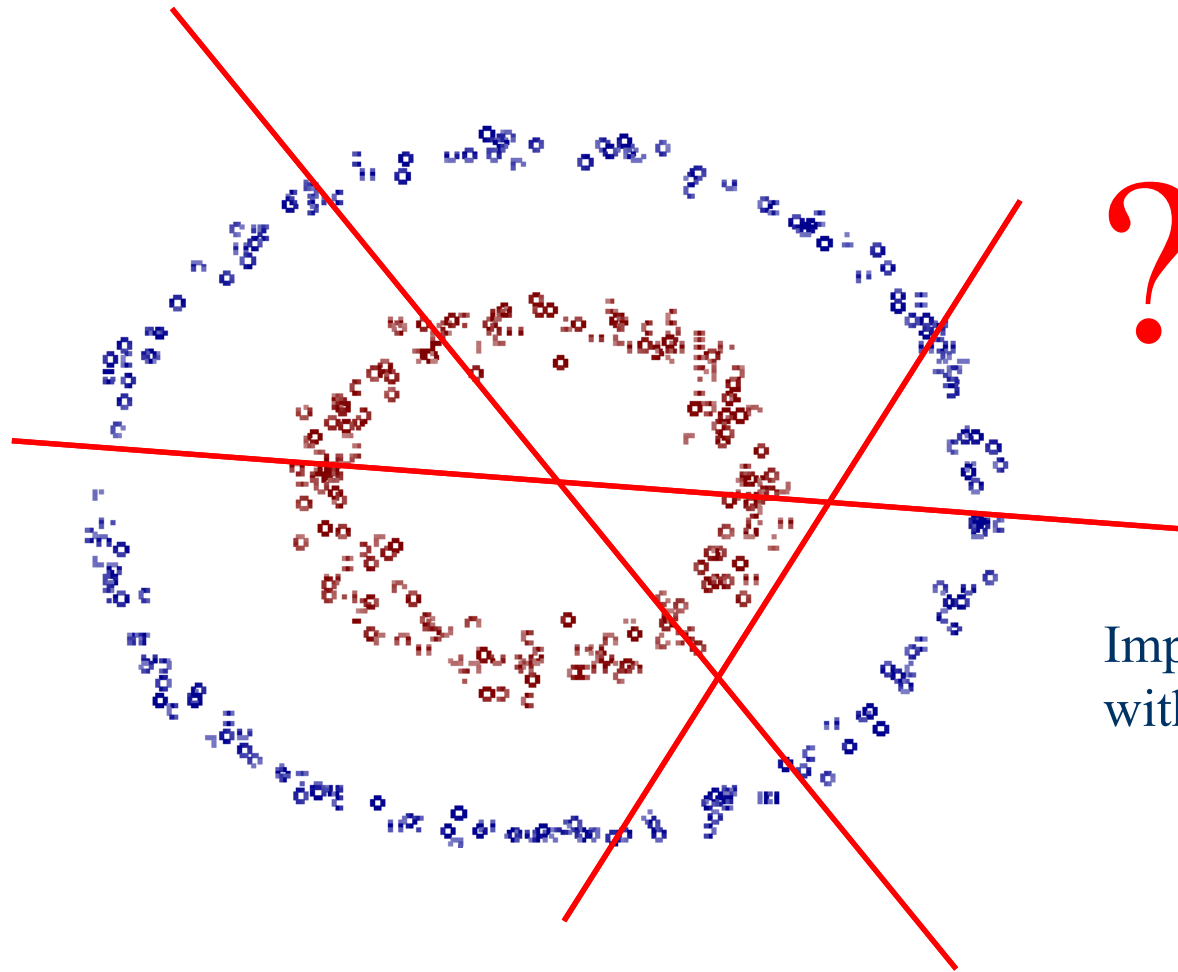
## So Far...

- ▶ Our goal (supervised learning):
  - To learn a set of discriminant functions
- ▶ Bayesian framework
  - We could design an optimal classifier if we knew:
    - $P(\omega_i)$  : priors and  $P(x | \omega_i)$  : class-conditional densities
    - Using training data to estimate  $P(\omega_i)$  and  $P(x | \omega_i)$
- ▶ Directly learning discriminant functions from the training data
  - We only know the form of the discriminant functions
  - Linear Regression
  - Logistic Regression
  - SVM

Linear



# Nonlinear Distributed Data



Impossible to separate  
with a hyperplane

# Generalized Linear Function & Kernel Methods

**Deng Cai (蔡登)**

College of Computer Science  
Zhejiang University

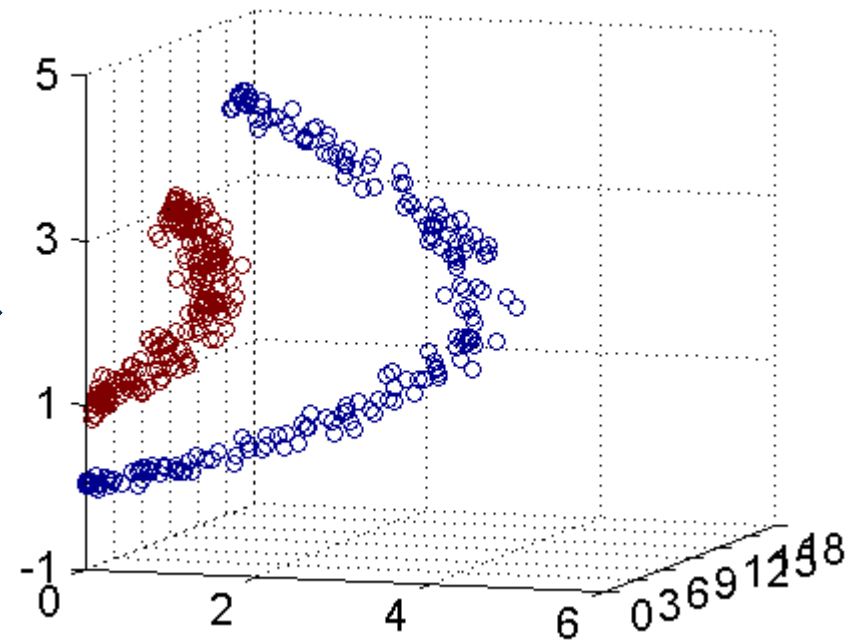
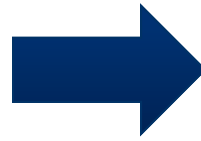
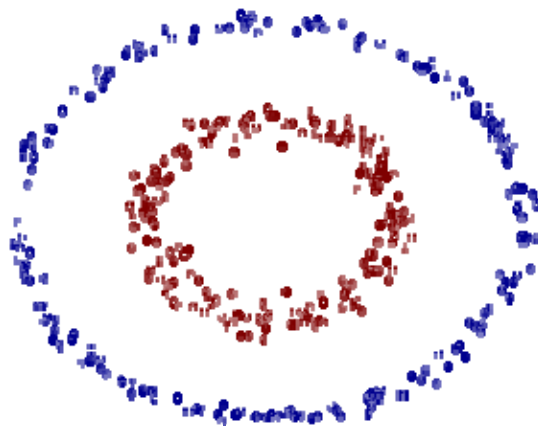
dengcai@gmail.com





# A Circle from 2D to 3D

- Here is an example of mapping a (special case) circle in 2D to 3D (the result is linear separable):





# Generalized Linear Discriminant Functions

- ▶ Recall the Linear Discriminant Function

$$g(\mathbf{x}) = \sum_{i=1}^d w_i x_i + w_0$$

- $g(\mathbf{x})$  positive implies class 1
- $g(\mathbf{x})$  negative implies class 2

- ▶ Generalized Linear Discriminant

- Add additional terms involving the products of features
- For example,
  - Given:  $[x_1, x_2, x_3]$
  - Make it:  $[x_1, x_2, x_3, x_1x_2, x_2x_3, x_1x_2x_3]$  by adding products of features.
- Learn a discriminant function that is linear in the new feature space



# Quadratic Discriminant Function

- ▶ Quadratic Discriminant Function
  - Obtained by adding pair-wise products of features

$$g(\mathbf{x}) = w_0 + \underbrace{\sum_{i=1}^d w_i x_i}_{\text{Linear Part}} + \underbrace{\sum_{i=1}^d \sum_{j=1}^d w_{ij} x_i x_j}_{\text{Quadratic part, } d(d+1)/2 \text{ additional parameters}}$$

Linear Part  
(d+1) parameters

Quadratic part,  $d(d+1)/2$   
additional parameters

- $g(x)$  positive implies class 1;  $g(x)$  negative implies class 2
- ▶  $g(x) = 0$ , represents a hyperquadric (hyperparaboloid, hyperellipsoid, hyperhyperboloids), as opposed to hyperplanes in linear discriminant case.
- ▶ Adding more terms such as  $w_{ijk} x_i x_j x_k$  results in polynomial discriminant functions.



# Quadratic Discriminant Function

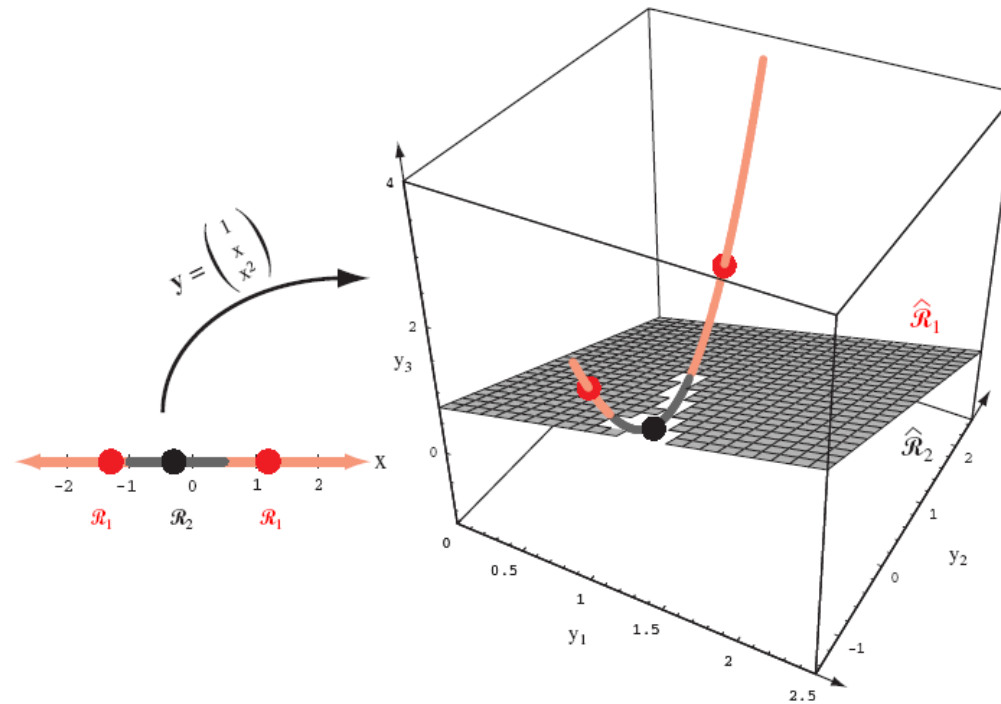


Figure 5.5: The mapping  $\mathbf{y} = (1, x, x^2)^t$  takes a line and transforms it to a parabola in three dimensions. A plane splits the resulting  $\mathbf{y}$  space into regions corresponding to two categories, and this in turn gives a non-simply connected decision region in the one-dimensional  $x$  space.



# Quadratic Discriminant Functions

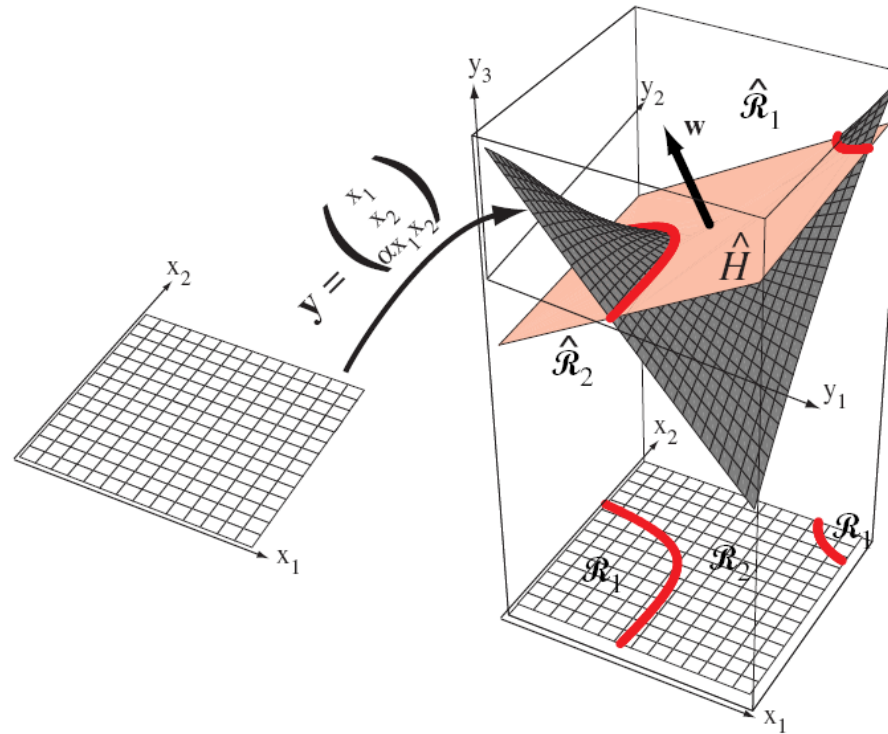


Figure 5.6: The two-dimensional input space  $\mathbf{x}$  is mapped through a polynomial function  $f$  to  $\mathbf{y}$ . Here the mapping is  $y_1 = x_1$ ,  $y_2 = x_2$  and  $y_3 \propto x_1x_2$ . A linear discriminant in this transformed space is a hyperplane, which cuts the surface. Points to the positive side of the hyperplane  $\hat{H}$  correspond to category  $\omega_1$ , and those beneath it  $\omega_2$ . Here, in terms of the  $\mathbf{x}$  space,  $\mathcal{R}_1$  is a not simply connected.





# Generalized Discriminant Function

- ▶ A *generalized linear discriminant* function can be written as,

Dimensionality of the augmented feature space.

$$g(\mathbf{x}) = \sum_{i=1}^{\hat{d}} a_i y_i(\mathbf{x})$$

Setting  $y_i(x)$  to be monomials results in polynomial discriminant functions

Weights in the augmented feature space. Note that the function is linear in  $a$ .

- ▶ Equivalently,

$$g(\mathbf{x}) = \mathbf{a}^t \mathbf{y}$$

$$\mathbf{a} = [a_1, a_2, \dots, a_{\hat{d}}]^t \quad \mathbf{y} = [y_1(x), y_2(x), \dots, y_{\hat{d}}(x)]^t$$

also called the **augmented feature vector**.



# Phi Function

- ▶ The discriminant function  $g(x)$  is not linear in  $x$ , but is linear in  $y$ .
- ▶ The mapping  $y = [y_1(x), y_2(x), \dots, y_{\hat{d}}(x)]^t$  is taking a  $d$ -dimensional vector  $x$  and mapping it to a  $\hat{d}$ -dimensional space. The mapping  $y$  is called the **phi-function**.
- ▶ When the input patterns  $x$  are non-linearly separable in the input space, mapping them using the *right* phi-function maps them to a space where the patterns are linearly separable.
- ▶ Unfortunately, the curse of dimensionality makes it hard to capitalize this in practice. A complete QDF involves  $(d+1)(d+2)/2$  terms; for modest values of  $d$ , say  $d=50$ , this requires many terms



# Representer Theorem

**Theorem 4.2 (Representer Theorem)** Denote by  $\Omega : [0, \infty) \rightarrow \mathbb{R}$  a strictly monotonic increasing function, by  $\mathcal{X}$  a set, and by  $c : (\mathcal{X} \times \mathbb{R}^2)^m \rightarrow \mathbb{R} \cup \{\infty\}$  an arbitrary loss function. Then each minimizer  $f \in \mathcal{H}$  of the regularized risk

$$c((x_1, y_1, f(x_1)), \dots, (x_m, y_m, f(x_m))) + \Omega(\|f\|_{\mathcal{H}}) \quad (4.4)$$

admits a representation of the form

$$f(x) = \sum_{i=1}^m \alpha_i k(x_i, x). \quad (4.5)$$



# Kernelized Ridge Regression

$$\mathbf{w}^* = \operatorname{argmin} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \mathbf{w})^2 + \lambda \sum_{j=1}^p w_j^2$$

$$\mathbf{w}^* = (X X^T + \lambda I)^{-1} X \mathbf{y}$$

► Woodbury matrix identity

$$(\mathbf{P}^{-1} + \mathbf{B}^T \mathbf{R}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{R}^{-1} = \mathbf{P} \mathbf{B}^T (\mathbf{B} \mathbf{P} \mathbf{B}^T + \mathbf{R})^{-1} \quad (158)$$

$$\mathbf{w}^* = X (X^T X + \lambda I)^{-1} \mathbf{y}$$

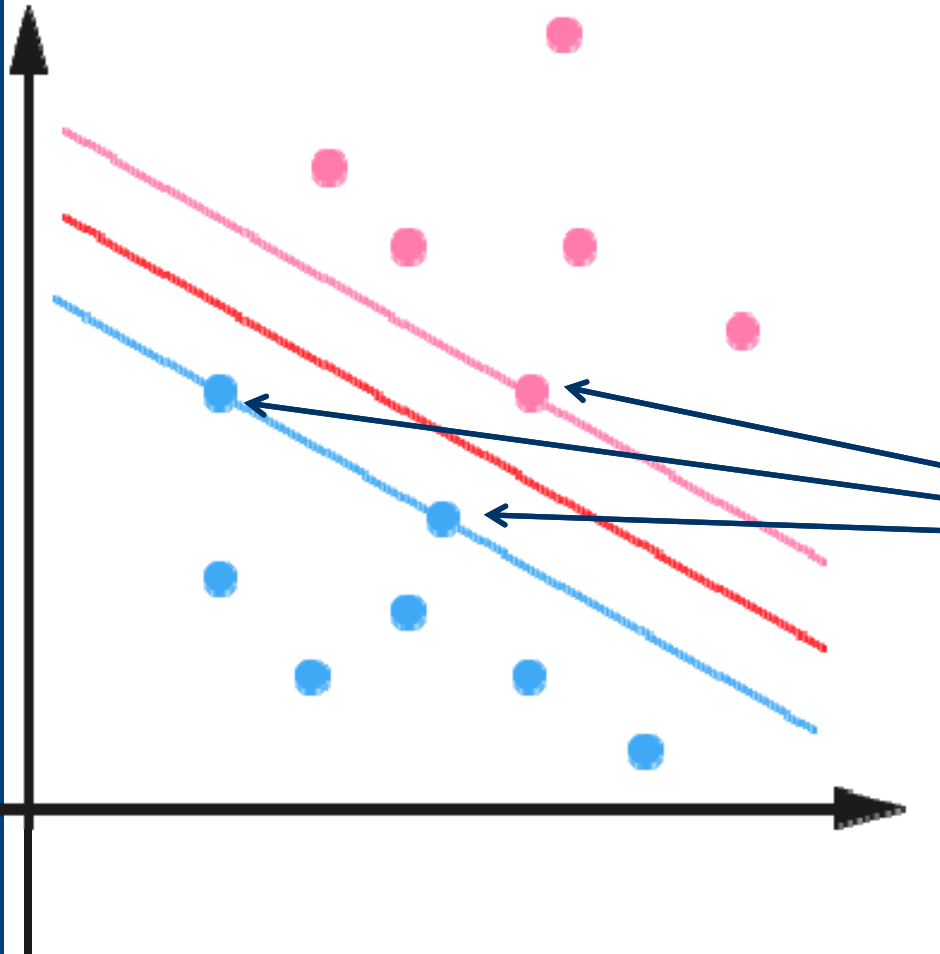
$$\boldsymbol{\alpha} \triangleq (X^T X + \lambda I)^{-1} \mathbf{y}$$

$$\mathbf{w}^* = X \boldsymbol{\alpha} = \sum_{i=1}^n \alpha_i \mathbf{x}_i$$

$$\hat{f}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = \sum_{i=1}^n \alpha_i \mathbf{x}_i^T \mathbf{x} = \sum_{i=1}^n \alpha_i \kappa(\mathbf{x}_i, \mathbf{x}) \quad \kappa(\mathbf{x}_i, \mathbf{x}) = \mathbf{x}_i^T \mathbf{x}$$



# Support Vector Machine



Hyper plane of maximum margin is *supported* by those points (vectors) on the margin. Those are called **Support Vectors**.

Non-support vectors can move freely without affecting the position of the hyperplane as long as they don't exceed the margin.



# Support Vector Machine

- ▶ The final classifier is

$$\text{sgn}(\mathbf{w}^T \mathbf{x} + b) = \text{sgn}\left(\sum_{i=1}^n \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle + b\right)$$

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

- ▶ Note: for non-support vectors, the corresponding  $\alpha_i$  is zero.



# Kernels

- ▶ Let  $\kappa(\mathbf{x}, \mathbf{x}') \geq 0$  be some measure of similarity between objects  $\mathbf{x}, \mathbf{x}' \in \chi$ , where  $\chi$  is some abstract space; we will call  $\kappa$  a **kernel function**.
  - Typically the function is symmetric, and non-negative
- ▶ Examples
  - Linear kernels  $\kappa(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$
  - Polynomial kernels  $\kappa(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + 1)^d$
  - RBF kernels  $\kappa(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right)$



# The advantages of kernel methods

- ▶ Non-linear classifiers
  - The kernel  $\rightarrow$  Nonlinearity of the learned function.
- ▶ The samples can not be represented as feature vectors
  - But we can get the similarity of two samples
  - String kernels
  - Graph kernels