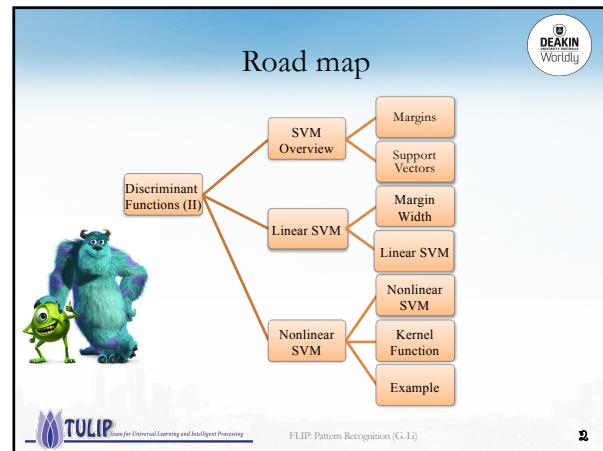


Lecture Notes on  
Pattern Recognition

Session 06(B): Discriminant Functions (II)

Gang Li  
School of Information Technology  
Deakin University, VIC 3125, Australia

DEAKIN Worldly



SVM Overview

- Best Linear Classifiers
- Margins
- Support Vectors
- VC Dimension

FLIP: Pattern Recognition (G. Li)

8

TULIP Team for Universal Learning and Intelligent Processing

Learning through empirical risk minimization

- Estimate  $g(\mathbf{x})$  from a finite set of observations by minimizing some kind of error function, for example, the empirical risk

$$R_{emp}(w, w_0) = \frac{1}{n} \sum_{k=1}^n [z_k - g(x_k, w, w_0)]^2$$

class labels:  $z_k = \begin{cases} +1 & \text{if } \mathbf{x}_k \in \omega_1 \\ -1 & \text{if } \mathbf{x}_k \in \omega_2 \end{cases}$

FLIP: Pattern Recognition (G. Li)

4

TULIP Team for Universal Learning and Intelligent Processing

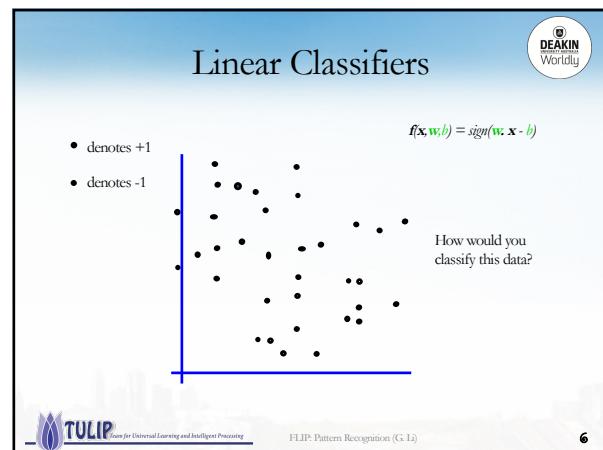
Learning through empirical risk minimization

- Conventional empirical risk minimization over the training data **does not** imply good generalization to novel test data.
  - There could be a number of different functions which all approximate the training data set well.
  - Difficult to determine a function which **best** captures the true underlying structure of the data distribution
    - i.e., has good generalization capabilities

FLIP: Pattern Recognition (G. Li)

5

TULIP Team for Universal Learning and Intelligent Processing



## Linear Classifiers

$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$

- denotes +1
- denotes -1

How would you classify this data?

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li)

## Linear Classifiers

$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$

- denotes +1
- denotes -1

How would you classify this data?

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li)

## Linear Classifiers

$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$

- denotes +1
- denotes -1

How would you classify this data?

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li)

## Linear Classifiers

$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$

- denotes +1
- denotes -1

Any of these would be fine..  
but which is best?

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li)

## Linear Classifiers

$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$

- denotes +1
- denotes -1

Define the margin of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li)

## Linear Classifiers

$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$

- denotes +1
- denotes -1

The maximum margin linear classifier is the linear classifier with the, um, maximum margin.  
This is the simplest kind of SVM (Called an LSVM)

Linear SVM

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li)

## Linear Classifiers

$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$

- denotes +1
- denotes -1

Support Vectors are those datapoints that the margin pushes up against

The **maximum margin linear classifier** is the linear classifier with the, um, maximum margin. This is the simplest kind of SVM (Called an LSVM)

Linear SVM

18

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

## Linear Classifiers

$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$

- denotes +1
- denotes -1

Support Vectors are those datapoints that the margin pushes up against

1. Intuitively this feels safest.

2. If we've made a small error in the location of the boundary (it's been tilted in its perpendicular direction) this gives us least chance of causing a misclassification.

3. LOOCV is easy since the model is immune to removal of any non-support-vector datapoints.

4. There's some theory (using VC dimension) that is related to (but not the same as) the proposition that this is a good thing.

5. Empirically it works very well.

19

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

## Statistical Learning: Capacity and VC dimension

- To guarantee good generalization performance, the capacity of the learned functions must be controlled.
- Functions with high capacity are more complicated (i.e., have many degrees of freedom).

low capacity

high capacity

width

length

salmon

sea bass

15

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

## Statistical Learning: Capacity and VC dimension

- In statistical learning, the **Vapnik-Chervonenkis (VC) dimension** is one of the most popular measures of capacity.
  - The **VC dimension** can predict a **probabilistic upper bound** on the test error (generalization error) of a classification model.

16

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

## Statistical Learning: Capacity and VC dimension

- A function that
  1. minimizes the empirical risk and
  2. has low VC dimension

will generalize well regardless of the dimensionality of the input space with probability (1- $\delta$ ):

$$\text{err}_{\text{true}} \leq \text{err}_{\text{train}} + \sqrt{\frac{VC(\log(2n/VC) + 1) - \log(\delta/4)}{n}}$$

n: training set size

17

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

## VC dimension and margin of separation

- Vapnik has shown that **maximizing** the **margin** of separation between classes is equivalent to **minimizing** the **VC dimension**.
  - The optimal hyperplane is the one giving the **largest margin** of separation between the classes.

y<sub>2</sub>

y<sub>1</sub>

optimal hyperplane

R<sub>1</sub>

R<sub>2</sub>

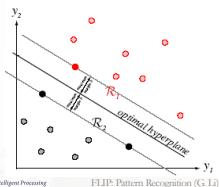
18

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

## Margin of separation and support vectors

- The margin (i.e., empty area around the decision boundary) is defined by the distance to the nearest training patterns which we refer to as **support vectors**.
  - Intuitively speaking, these are the most difficult patterns to classify.



19

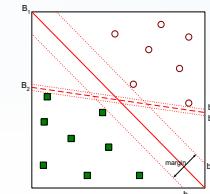
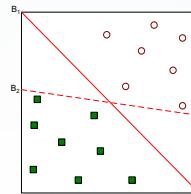


TULIP Team for Universal Learning and Intelligent Processing



## Margin of separation and support vectors

different solutions      corresponding margins



FLIP- Pattern Recognition (G. Li)

20

## SVM Overview

- SVMs perform **structural risk minimization** to achieve good generalization performance.
- The optimization criterion is the **margin** of separation between classes.
- Training is equivalent to solving a **quadratic programming** problem with **linear constraints**.
- Primarily **two-class** classifiers but can be extended to **multiple** classes.



TULIP Team for Universal Learning and Intelligent Processing



FLIP- Pattern Recognition (G. Li)

21

## Linear SVM



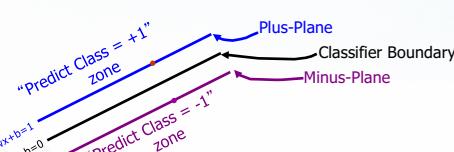
- Margin Width
- Linear SVM

FLIP- Pattern Recognition (G. Li)

22

## Specifying a line and margin

- How do we represent this mathematically?
  - ...in  $m$  input dimensions?



TULIP Team for Universal Learning and Intelligent Processing

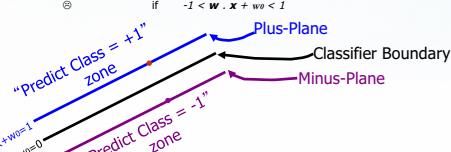
FLIP- Pattern Recognition (G. Li)

23

## Specifying a line and margin

- Plus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + w_0 = +1 \}$
- Minus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + w_0 = -1 \}$

Classify as..	+1	if $\mathbf{w} \cdot \mathbf{x} + w_0 >= 1$
	-1	if $\mathbf{w} \cdot \mathbf{x} + w_0 <= -1$
	0	if $-1 < \mathbf{w} \cdot \mathbf{x} + w_0 < 1$

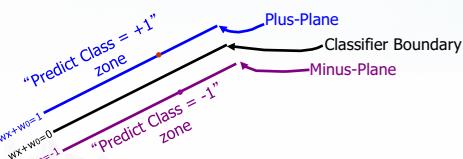


FLIP- Pattern Recognition (G. Li)

24

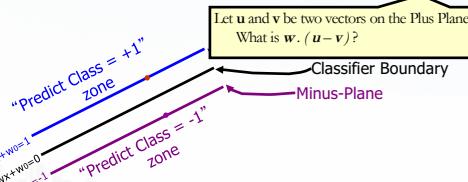
## Specifying a line and margin

- Plus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + w_0 = +1 \}$
- Minus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + w_0 = -1 \}$
- Claim: The vector  $\mathbf{w}$  is perpendicular to the Plus Plane. Why?



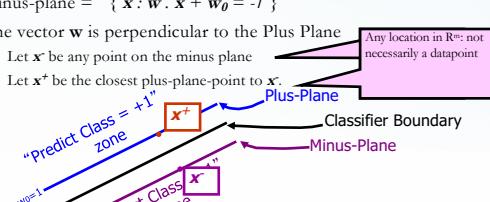
## Specifying a line and margin

- Plus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + w_0 = +1 \}$
- Minus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + w_0 = -1 \}$
- Claim: The vector  $\mathbf{w}$  is perpendicular to the Plus Plane. Why?



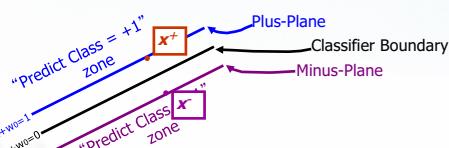
## Specifying a line and margin

- Plus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + w_0 = +1 \}$
- Minus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + w_0 = -1 \}$
- The vector  $\mathbf{w}$  is perpendicular to the Plus Plane
  - Let  $\mathbf{x}^*$  be any point on the minus plane
  - Let  $\mathbf{x}^+$  be the closest plus-plane-point to  $\mathbf{x}^*$ .



## Specifying a line and margin

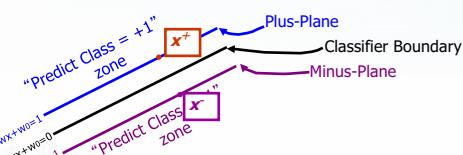
- Plus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + w_0 = +1 \}$
- Minus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + w_0 = -1 \}$
- The vector  $\mathbf{w}$  is perpendicular to the Plus Plane
- Claim:  $\mathbf{x}^+ = \mathbf{x}^* + \lambda \mathbf{w}$  for some value of  $\lambda$ . Why?



## Specifying a line and margin

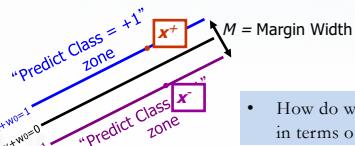
- Plus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + w_0 = +1 \}$
- Minus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + w_0 = -1 \}$
- The vector  $\mathbf{w}$  is perpendicular to the Plus Plane
- Claim:  $\mathbf{x}^+ = \mathbf{x}^* + \lambda \mathbf{w}$  for some value of  $\lambda$ . Why?

• The line from  $\mathbf{x}^*$  to  $\mathbf{x}^+$  is perpendicular to the planes.  
 • So to get from  $\mathbf{x}^*$  to  $\mathbf{x}^+$  travel some distance in direction  $\mathbf{w}$ .



## Specifying a line and margin

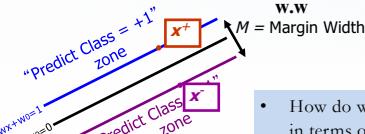
- What we know:
  - $\mathbf{w} \cdot \mathbf{x}^* + w_0 = +1$
  - $\mathbf{w} \cdot \mathbf{x}^+ + w_0 = -1$
  - $\mathbf{x}^+ = \mathbf{x}^* + \lambda \mathbf{w}$
  - $|\mathbf{x}^+ - \mathbf{x}^*| = M$



## Specifying a line and margin

- What we know:
  - $\mathbf{w} \cdot \mathbf{x}^+ + w_0 = +1$
  - $\mathbf{w} \cdot \mathbf{x}^- + w_0 = -1$
  - $\mathbf{x}^+ = \mathbf{x}^- + \lambda \mathbf{w}$
  - $|\mathbf{x}^+ - \mathbf{x}^-| = M$

$$\begin{aligned} \mathbf{w} \cdot (\mathbf{x}^- + \lambda \mathbf{w}) + w_0 &= 1 \\ \Rightarrow \mathbf{w} \cdot \mathbf{x}^- + w_0 + \lambda \mathbf{w} \cdot \mathbf{w} &= 1 \\ \Rightarrow -1 + \lambda \mathbf{w} \cdot \mathbf{w} &= 1 \\ \Rightarrow \lambda &= \frac{2}{\mathbf{w} \cdot \mathbf{w}} \end{aligned}$$



- How do we compute  $M$  in terms of  $\mathbf{w}$  and  $b$ ?

## Specifying a line and margin

- What we know:

$$\mathbf{w} \cdot \mathbf{x}^+ + w_0 = +1$$

$$\mathbf{w} \cdot \mathbf{x}^- + w_0 = -1$$

$$\mathbf{x}^+ = \mathbf{x}^- + \lambda \mathbf{w}$$

$$|\mathbf{x}^+ - \mathbf{x}^-| = M$$

$$\lambda = \frac{2}{\mathbf{w} \cdot \mathbf{w}}$$

$$M = |\mathbf{x}^+ - \mathbf{x}^-| = \lambda \mathbf{w} |$$

$$= \lambda |\mathbf{w}| = \lambda \sqrt{\mathbf{w} \cdot \mathbf{w}}$$

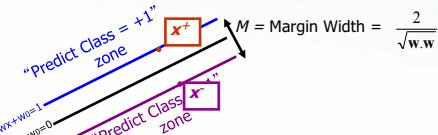
$$= \frac{2\sqrt{\mathbf{w} \cdot \mathbf{w}}}{\mathbf{w} \cdot \mathbf{w}} = \frac{2}{\sqrt{\mathbf{w} \cdot \mathbf{w}}}$$

$$M = \text{Margin Width}$$

- How do we compute  $M$  in terms of  $\mathbf{w}$  and  $b$ ?

## Specifying a line and margin

- Given a guess of  $\mathbf{w}$  and  $w_0$  we can
  - Compute whether all data points in the correct half-planes
  - Compute the width of the margin
- So now we need to search the space of  $\mathbf{w}$ 's to find the widest margin that matches all the datapoints. **How?**



## Linear SVM: separable case

- Linear discriminant

$$g(\mathbf{x}) = \mathbf{w}' \mathbf{x} + w_0$$

Decide  $\omega_1$  if  $g(\mathbf{x}) > 0$  and  $\omega_2$  if  $g(\mathbf{x}) < 0$

- Class labels

$$z_k = \begin{cases} +1 & \text{if } \mathbf{x}_k \in \omega_1 \\ -1 & \text{if } \mathbf{x}_k \in \omega_2 \end{cases}$$

- Normalized version

$$z_k g(\mathbf{x}_k) > 0 \text{ or } z_k (\mathbf{w}' \mathbf{x}_k + w_0) > 0, \text{ for } k = 1, 2, \dots, n$$

## Linear SVM: separable case

- The distance of a point  $\mathbf{x}_k$  from the separating hyperplane should satisfy the **constraint**:

$$\frac{z_k g(\mathbf{x}_k)}{\|\mathbf{w}\|} \geq b, \quad b > 0$$

- To ensure uniqueness, impose:

$$b \|\mathbf{w}\| = 1$$

- The above **constraint** becomes:

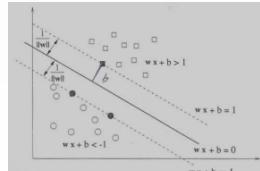
$$z_k g(\mathbf{x}_k) \geq 1 \text{ or } z_k (\mathbf{w}' \mathbf{x}_k + w_0) > 1 \quad \text{where } b = \frac{1}{\|\mathbf{w}\|}$$

## Linear SVM: separable case

- Optimization:

- Maximize the margin  $\frac{2}{\sqrt{\mathbf{w} \cdot \mathbf{w}}}$
- With constraints

Problem 1: Minimize  $\frac{1}{2} \|\mathbf{w}\|^2$   
subject to  $z_k (\mathbf{w}' \mathbf{x}_k + w_0) \geq 1, \quad k = 1, 2, \dots, n$



• Quadratic programming problem!

## Linear SVM: separable case

- Use Langrange optimization, we seek to minimize:

$$L(\mathbf{w}, w_0, \lambda) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{k=1}^n \lambda_k [z_k (\mathbf{w}' \mathbf{x}_k + w_0) - 1], \quad \lambda_k \geq 0$$

- Easier to solve the “dual” problem
  - Kuhn-Tucker construction

**Problem 2:** Maximize  $\sum_{k=1}^n \lambda_k - \frac{1}{2} \sum_{k,j} \lambda_k \lambda_j z_k z_j \mathbf{x}_j' \mathbf{x}_k$   
 subject to  $\sum_{k=1}^n z_k \lambda_k = 0, \quad \lambda_k \geq 0, \quad k = 1, 2, \dots, n$

  FLIP: Pattern Recognition (G. Li) 37

## Linear SVM: separable case

- The solution is given by:

$$\mathbf{w} = \sum_{k=1}^n z_k \lambda_k \mathbf{x}_k, \quad w_0 = z_k - \mathbf{w}' \mathbf{x}_k$$

$$g(\mathbf{x}) = \mathbf{w}' \mathbf{x} + w_0$$

 • dot product  

$$g(\mathbf{x}) = \sum_{k=1}^n z_k \lambda_k (\mathbf{x}'_k \mathbf{x}) + w_0 = \sum_{k=1}^n z_k \lambda_k (\mathbf{x}' \cdot \mathbf{x}_k) + w_0$$

- It can be shown that if  $\mathbf{x}_k$  is not a support vector, then the corresponding  $\lambda_k = 0$ .
  - Only support vectors contribute to the solution!

  FLIP: Pattern Recognition (G. Li) 38

## Linear SVM: non-separable case

- Allow miss-classifications (i.e., soft margin classifier) by introducing positive error (slack) variables  $\psi_k$ :

$$z_k (\mathbf{w}' \mathbf{x}_k + w_0) \geq 1 - \psi_k, \quad \psi_k \geq 0, \quad k = 1, 2, \dots, n$$

**Problem 3:** Minimize  $\frac{1}{2} \|\mathbf{w}\|^2 + c \sum_{k=1}^n \psi_k$   
 subject to  $z_k (\mathbf{w}' \mathbf{x}_k + w_0) \geq 1 - \psi_k, \quad k = 1, 2, \dots, n$

- The result is a hyperplane that minimizes the sum of errors  $\psi_k$  while maximizing the margin for the correctly classified data.
- constant  $c$  controls the trade-off between the margin and misclassification errors.
- Aims to prevent outliers from affecting the optimal hyperplane.

  FLIP: Pattern Recognition (G. Li) 39

## Linear SVM: non-separable case

- Easier to solve the “dual” problem
  - Kuhn-Tucker construction

**Problem 4:** Maximize  $\sum_{k=1}^n \lambda_k - \frac{1}{2} \sum_{k,j} \lambda_k \lambda_j z_k z_j \mathbf{x}_j' \mathbf{x}_k$   
 subject to  $\sum_{k=1}^n z_k \lambda_k = 0$  and  $0 \leq \lambda_k \leq c, \quad k = 1, 2, \dots, n$

where the use of error variables  $\psi_k$  constraint the range of the Lagrange coefficients from 0 to  $c$ .

  FLIP: Pattern Recognition (G. Li) 40

## Nonlinear SVM

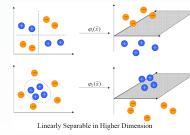
- Nonlinear SVM
- Kernel Functions
- Example



  FLIP: Pattern Recognition (G. Li) 41

## Nonlinear SVM

- Extending these concepts to the non-linear case involves mapping the data to a high-dimensional space  $\mathbb{R}$ :

$$\mathbf{x}_k \rightarrow \Phi(\mathbf{x}_k) = \begin{bmatrix} \varphi_1(\mathbf{x}_k) \\ \varphi_2(\mathbf{x}_k) \\ \vdots \\ \varphi_h(\mathbf{x}_k) \end{bmatrix}$$

- Mapping the data to a sufficiently high dimensional space is likely to cast the data linearly separable in that space.

  FLIP: Pattern Recognition (G. Li) 42

## Nonlinear SVM

**DEAKIN Worldwide**

- linear SVM:  $g(\mathbf{x}) = \sum_{k=1}^n z_k \lambda_k (\mathbf{x} \cdot \mathbf{x}_k) + w_0$

- non-linear SVM:  $g(\mathbf{x}) = \sum_{k=1}^n z_k \lambda_k (\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}_k)) + w_0$

- Decide  $\omega_1$  if  $g(\mathbf{x}) > 0$  and  $\omega_2$  if  $g(\mathbf{x}) < 0$

**TULIP** Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 43

## Nonlinear SVM

**DEAKIN Worldwide**

- The disadvantage of this approach is that the mapping  $\mathbf{x}_k \rightarrow \Phi(\mathbf{x}_k)$  might be very computationally intensive to compute!

- non-linear SVM:  $g(\mathbf{x}) = \sum_{k=1}^n z_k \lambda_k (\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}_k)) + w_0$

**TULIP** Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 44

## Nonlinear SVM

**DEAKIN Worldwide**

- The kernel trick
  - Compute dot products using a kernel function

$$K(\mathbf{x}, \mathbf{x}_k) = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}_k)$$

$$g(\mathbf{x}) = \sum_{k=1}^n z_k \lambda_k K(\mathbf{x}, \mathbf{x}_k) + w_0$$

**TULIP** Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 45

## Nonlinear SVM

**DEAKIN Worldwide**

- The kernel trick
  - Compute dot products using a kernel function
  - Kernel functions which can be expressed as a dot product in some space satisfy the Mercer's condition
    - The Mercer's condition does not tell us how to construct  $\Phi()$  or even what the high dimensional space is.
    - see Burges' paper
  - Advantages of kernel trick
    - no need to know  $\Phi()$
    - computations remain feasible even if the feature space has high dimensionality.

**TULIP** Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 46

## Polynomial Kernel

**DEAKIN Worldwide**

- $K(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y})^d$

\* It can be shown for the case of polynomial kernels that the data is mapped to a space of dimension  $h = \binom{p+d-1}{d}$  where  $p$  is the original dimensionality.

\* Suppose  $p=256$  and  $d=4$ , then  $h=183,181,376$  !!

\* A dot product in the high dimensional space would require  $O(h)$  computations while the kernel requires only  $O(p)$  computations.

**TULIP** Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 47

## Choice of $\Phi$ is not unique

**DEAKIN Worldwide**

Example: consider  $\mathbf{x} \in \mathbb{R}^2$ ,  $\Phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix} \in \mathbb{R}^3$ , and  $K(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y})^2$

$$(\mathbf{x}, \mathbf{y})^2 = (x_1y_1 + x_2y_2)^2$$

$$\Phi(\mathbf{x}) \cdot \Phi(\mathbf{y}) = x_1^2y_1^2 + 2x_1y_1x_2y_2 + x_2^2y_2^2 = (x_1y_1 + x_2y_2)^2$$

- Note that neither the mapping  $\Phi()$  nor the high dimensional space are unique.

$$\Phi(\mathbf{x}) = \frac{1}{\sqrt{2}} \begin{pmatrix} (x_1^2 - x_2^2) \\ 2x_1x_2 \\ (x_1^2 + x_2^2) \end{pmatrix} \in \mathbb{R}^3 \quad \text{or} \quad \Phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ x_1x_2 \\ x_2^2 \end{pmatrix} \in \mathbb{R}^4$$

**TULIP** Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 48

## Kernel functions

- By using different kernel functions, SVM implement a variety of learning machines, some of which coincide with classical architectures

*polynomial:*  $K(x, x_k) = (x \cdot x_k)^d$

*sigmoidal:*  $K(x, x_k) = \tanh(v_k(x \cdot x_k) + c_k)$   
(corresponds to a two-layer sigmoidal neural network)

*Gaussian:*  $K(x, x_k) = \exp\left(\frac{-\|x - x_k\|^2}{2\sigma_k^2}\right)$   
(corresponds to a radial basis function (RBF) neural network)

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 49

## Example

- Consider the XOR problem which is non-linear separable

(1,1) and (-1, -1) belong to  $\omega_1$   
(1,-1) and (-1, 1) belong to  $\omega_2$

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 50

## Example

- Consider the following mapping (among others)

$$y = \Phi(x) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1 \\ \sqrt{2}x_1x_2 \\ \sqrt{2}x_2 \\ x_2^2 \\ 1 \end{pmatrix} \quad h=6$$

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 51

## Example

- The above transformation maps  $x_k$  to a 6-dimensional space

$$y_1 = \Phi(x_1) = \begin{pmatrix} 1 \\ \sqrt{2} \\ \sqrt{2} \\ 1 \\ 1 \end{pmatrix} \quad y_3 = \Phi(x_3) = \begin{pmatrix} 1 \\ -\sqrt{2} \\ \sqrt{2} \\ 1 \\ 1 \end{pmatrix}$$

$$y_2 = \Phi(x_2) = \begin{pmatrix} 1 \\ -\sqrt{2} \\ -\sqrt{2} \\ 1 \\ 1 \end{pmatrix} \quad y_4 = \Phi(x_4) = \begin{pmatrix} 1 \\ -\sqrt{2} \\ -\sqrt{2} \\ \sqrt{2} \\ 1 \\ 1 \end{pmatrix}$$

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 52

## Example

- We seek to maximize

$$\sum_{k=1}^4 \lambda_k - \frac{1}{2} \sum_{k,j} \lambda_k \lambda_j z_k z_j \Phi(x_j^T) \Phi(x_k)$$

subject to  $\sum_{k=1}^4 z_k \lambda_k = 0$ ,  $\lambda_k \geq 0$ ,  $k = 1, 2, \dots, 4$

- The solution turns out to be:

$$\lambda_1 = \lambda_2 = \lambda_3 = \lambda_4 = \frac{1}{8}$$

- Since all  $\lambda_k \neq 0$ , all  $x_k$  are support vectors !

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 53

## Example

- We now compute  $w$

$$w = \sum_{k=1}^4 z_k \lambda_k \Phi(x_k) = \frac{1}{8} \begin{pmatrix} 1 \\ \sqrt{2} \\ \sqrt{2} \\ 1 \\ 1 \end{pmatrix} - \frac{1}{8} \begin{pmatrix} 1 \\ -\sqrt{2} \\ -\sqrt{2} \\ 1 \\ 1 \end{pmatrix} + \frac{1}{8} \begin{pmatrix} 1 \\ \sqrt{2} \\ \sqrt{2} \\ 1 \\ 1 \end{pmatrix} - \frac{1}{8} \begin{pmatrix} 1 \\ -\sqrt{2} \\ -\sqrt{2} \\ \sqrt{2} \\ 1 \\ 1 \end{pmatrix} = \frac{1}{2} \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

- The solution for  $w_0$  can be determined using any support vector, e.g.,  $x_1$ :

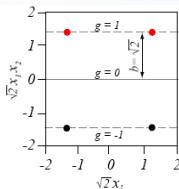
$$w^T \Phi(x_1) + w_0 = z_1 \quad \text{or} \quad w_0 = 0$$

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 54

## Example

- The margin  $b$  is computed as follows

$$\mathbf{w} = \frac{1}{2} \begin{pmatrix} 0 \\ 0 \\ \sqrt{2} \\ 0 \\ 0 \end{pmatrix} \quad b = \frac{1}{\|\mathbf{w}\|} = \sqrt{2}$$



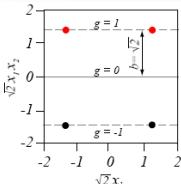
## Example

- The discriminant function is as follows

$$g(\mathbf{x}) = \mathbf{w}^T \Phi(\mathbf{x}) + w_0 = x_1 x_2$$

where we decide  $\omega_1$  if  $g(\mathbf{x}) > 0$  and  $\omega_2$  if  $g(\mathbf{x}) < 0$

$$\mathbf{w} = \frac{1}{2} \begin{pmatrix} 0 \\ 0 \\ \sqrt{2} \\ 0 \\ 0 \end{pmatrix} \quad \Phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1 x_2 \\ \sqrt{2}x_2^2 \\ 1 \end{pmatrix}$$



## Summary of SVM

- SVM is based on exact optimization, not on approximate methods
  - i.e., global optimization method, ***no local optima***
- Avoid ***overfitting*** in high dimensional spaces and ***generalize*** well using a small training set.
- Performance depends on the choice of the ***kernel*** and its parameters.
- Its ***complexity*** depends on the number of support vectors, **not** on the dimensionality of the transformed space.

## Seminar S10

### Topics

- VC Dimension
- PAC Learning

### Requirements

- Prepare a **15 minutes** talk on your chosen topic
- Make **ppt** to assist your talk
- Prepare **at least 3 questions** to ask the audience after your talk
- Get ready to **take questions** from the audience

### Hints

- C.J.C. Burges. A tutorial on support vector machines for pattern recognition. Data Mining and Knowledge Discovery, 2(2):955-974, 1998.  
<http://citeseer.nj.nec.com/burges98tutorial.html>

## Questions?

