**Slide 1**

Lecture Notes on

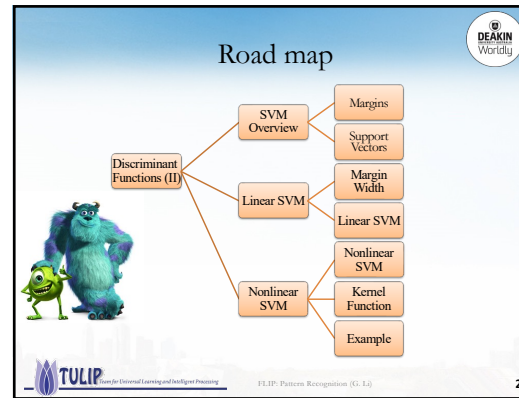## Pattern Recognition

**Session 07(B): Discriminant Functions (II)**

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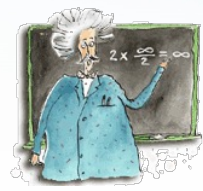Deakin University, VIC 3125, Australia

---

**Slide 2**

## Road map

- Discriminant Functions (II)
  - SVM Overview
    - Margins
    - Support Vectors
  - Linear SVM
    - Margin Width
    - Linear SVM
  - Nonlinear SVM
    - Nonlinear SVM
    - Kernel Function
    - Example

TULIP *Team for Universal Learning and Intelligent Processing*  FLIP: Pattern Recognition (G. Li)  2

---

**Slide 3**

## SVM Overview

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

TULIP *Team for Universal Learning and Intelligent Processing*  FLIP: Pattern Recognition (G. Li)  3

---

**Slide 4**

## Learning through empirical risk minimization

- Estimate $g(x)$ from a finite set of observations by minimizing some kind of error function, for example, the *empirical risk (sample risk)*:

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

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

TULIP *Team for Universal Learning and Intelligent Processing*  FLIP: Pattern Recognition (G. Li)  4

---

**Slide 5**

## 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

TULIP *Team for Universal Learning and Intelligent Processing*  FLIP: Pattern Recognition (G. Li)  5

---

**Slide 6**

## Linear Classifiers

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

- denotes +1
- denotes -1

How would you classify this data?

TULIP *Team for Universal Learning and Intelligent Processing*  FLIP: Pattern Recognition (G. Li)  6

7



8



9



10



11



12

## Slide 13

### Linear Classifiers

- denotes +1
- denotes -1

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

Support Vectors are those data points 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)
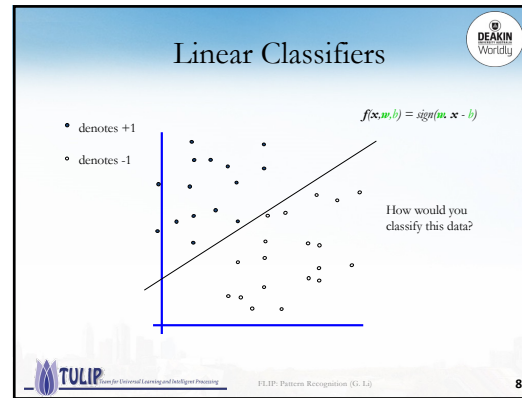
Linear SVM

FLIP: Pattern Recognition (G. Li)

13

---

## Slide 14

### Linear Classifiers

- denotes +1
- denotes -1

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

Support Vectors are those datapoints that the margin pushes up against

1. Intuitively this feels safest.
2. If we have made a small error in the location of the boundary (it's been jolted 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 data points.
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 very well.

FLIP: Pattern Recognition (G. Li)

14

---

## Slide 15

### 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

FLIP: Pattern Recognition (G. Li)

15

---

## Slide 16

### 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.

FLIP: Pattern Recognition (G. Li)

16

---

## Slide 17

### 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-δ):

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

n: training set size

FLIP: Pattern Recognition (G. Li)

17

---

## Slide 18

### 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 hyper-plane is the one giving the **largest margin** of separation between the classes.

FLIP: Pattern Recognition (G. Li)

18

## Slide 19

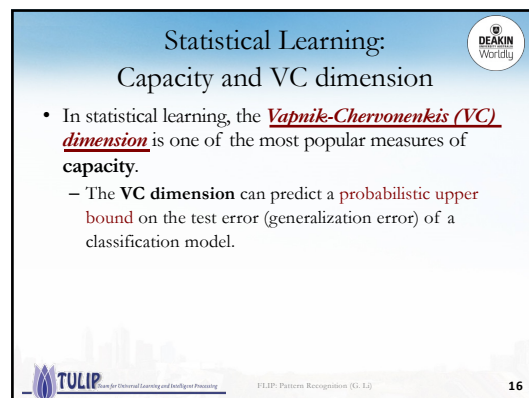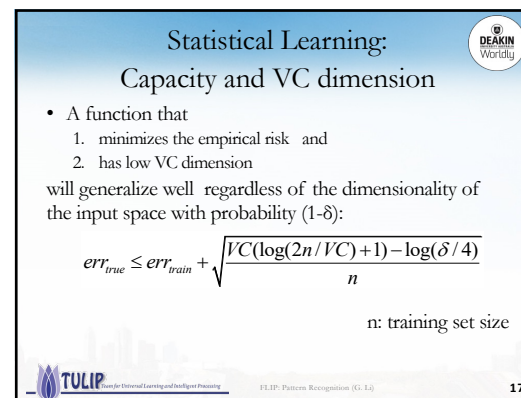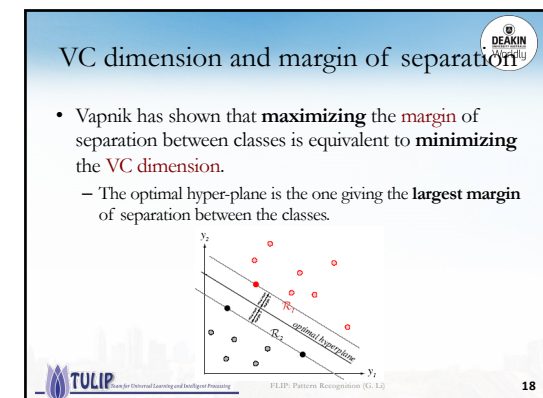### Margin of separation and support vectors

- 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.

FLIP: Pattern Recognition (G. Li)    19

19

## Slide 20

### Margin of separation and support vectors

different solutions          corresponding margins

FLIP: Pattern Recognition (G. Li)    20

20

## Slide 21

### 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.

FLIP: Pattern Recognition (G. Li)    21

21

## Slide 22

### Linear SVM

- Margin Width
- Linear SVM

FLIP: Pattern Recognition (G. Li)    22

22

## Slide 23

### Specifying a line and margin

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

FLIP: Pattern Recognition (G. Li)    23

23

## Slide 24

### Specifying a line and margin

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

Classify as..   +1    if    $\mathbf{w} \cdot \mathbf{x} + w_0 \geq 1$
                -1    if    $\mathbf{w} \cdot \mathbf{x} + w_0 \leq -1$
                ☺     if    $-1 < \mathbf{w} \cdot \mathbf{x} + w_0 < 1$

FLIP: Pattern Recognition (G. Li)    24

24

## Slide 25

### Specifying a line and margin

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

"Predict Class = +1" zone
Plus-Plane
Classifier Boundary
Minus-Plane
wx+w0=1
wx+w0=0
wx+w0=-1
"Predict Class = -1" zone

FLIP: Pattern Recognition (G. Li) — 25

25

## Slide 26

### Specifying a line and margin

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

And so of course the vector **w** is also perpendicular to the Minus Plane

Let **u** and **v** be two vectors on the Plus Plane. What is $w \cdot (u - v)$?

"Predict Class = +1" zone
Classifier Boundary
Minus-Plane
wx+w0=1
wx+w0=0
wx+w0=-1
"Predict Class = -1" zone

FLIP: Pattern Recognition (G. Li) — 26

26

## Slide 27

### Specifying a line and margin

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

Any location in $R^m$: not necessarily a data point

"Predict Class = +1" zone
$x^+$
Plus-Plane
Classifier Boundary
Minus-Plane
wx+w0=1
wx+w0=0
wx+w0=-1
$x^-$
"Predict Class = -1" zone

FLIP: Pattern Recognition (G. Li) — 27

27

## Slide 28

### Specifying a line and margin

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

"Predict Class = +1" zone
$x^+$
Plus-Plane
Classifier Boundary
Minus-Plane
wx+w0=0
wx+w0=-1
$x^-$
"Predict Class = -1" zone

FLIP: Pattern Recognition (G. Li) — 28

28

## Slide 29

### Specifying a line and margin

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

- The line from $x^-$ to $x^+$ is perpendicular to the planes.
- So to get from $x^-$ to $x^+$ travel some distance in direction **w**.

"Predict Class = +1" zone
$x^+$
Plus-Plane
Classifier Boundary
Minus-Plane
wx+w0=0
wx+w0=-1
$x^-$
"Predict Class = -1" zone

FLIP: Pattern Recognition (G. Li) — 29

29

## Slide 30

### Specifying a line and margin

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

"Predict Class = +1" zone
$x^+$
$M$ = Margin Width
$x^-$
"Predict Class = -1" zone
wx+w0=1
wx+w0=0
wx+w0=-1

- How do we compute $M$ in terms of $w$ and $b$?

FLIP: Pattern Recognition (G. Li) — 30

30

## Slide 31

### Specifying a line and margin

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

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

$M$ = Margin Width

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

FLIP: Pattern Recognition (G. Li)  31

31

## Slide 32

### Specifying a line and margin

- What we know:
  - $\mathbf{w} \cdot \mathbf{x}^+ + w_0 = +1$
  - $\mathbf{w} \cdot \mathbf{x}^- + w_0 = -1$
  - $\mathbf{x}^+ = \mathbf{x}^- + \lambda \mathbf{w}$
  - $|\mathbf{x}^+ - \mathbf{x}^-| = M$
  - $\lambda = \frac{2}{\mathbf{w} \cdot \mathbf{w}}$

$$M = |\mathbf{x}^+ - \mathbf{x}^-| = |\lambda \mathbf{w}|$$
$$= \lambda |\mathbf{w}| = \lambda \sqrt{\mathbf{w} \cdot \mathbf{w}}$$
$$= \frac{2\sqrt{\mathbf{w} \cdot \mathbf{w}}}{\mathbf{w} \cdot \mathbf{w}} = \frac{2}{\sqrt{\mathbf{w} \cdot \mathbf{w}}}$$

$M$ = Margin Width

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

FLIP: Pattern Recognition (G. Li)  32

32

## Slide 33

### 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 data points. *How?*

$M$ = Margin Width = $\frac{2}{\sqrt{\mathbf{w} \cdot \mathbf{w}}}$

FLIP: Pattern Recognition (G. Li)  33

33

## Slide 34

### Linear SVM: separable case

- Linear discriminant

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

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

- Class labels

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

- Normalized version

$$z_k g(\mathbf{x}_k) > 0 \quad or \quad z_k(\mathbf{w}^t \mathbf{x}_k + w_0) > 0, \quad for \ k = 1, 2, ..., n$$

34

## Slide 35

### Linear SVM: separable case

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

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

- To ensure uniqueness, impose:

$$b\|w\| = 1$$

- The above *constraint* becomes:

$$z_k g(\mathbf{x}_k) \ge 1 \quad or \quad z_k(\mathbf{w}^t \mathbf{x}_k + w_0) > 1 \quad where \quad b = \frac{1}{\|w\|}$$

FLIP: Pattern Recognition (G. Li)  35

35

## Slide 36

### Linear SVM: separable case

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

**Problem 1**: Minimize $\frac{1}{2}\|w\|^2$

subject to $z_k(\mathbf{w}^t x_k + w_0) \ge 1, \quad k = 1, 2, ..., n$

- Quadratic programming problem !

36

## Linear SVM: separable case

- Use **_Langrange_** optimization, we seek to minimize:

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

- Easier to solve the "dual" problem
  - **_Kuhn-Tucker construction_**:

  **Problem 2**: Maximize $\sum_{k=1}^{n} \lambda_k - \frac{1}{2} \sum_{k,j}^{n} \lambda_k \lambda_j z_k z_j x_j^t x_k$

  subject to $\sum_{k=1}^{n} z_k \lambda_k = 0, \quad \lambda_k \geq 0, k = 1, 2, \ldots, n$

FLIP: Pattern Recognition (G. Li)   **37**

37

## Linear SVM: separable case

- The solution is given by:

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

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

- dot product

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

- 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**

38

## Linear SVM: non-separable case

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

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

**Problem 3**: Minimize $\frac{1}{2} \|w\|^2 + c \sum_{k=1}^{n} \psi_k$

subject to $z_k(w^t x_k + w_0) \geq 1 - \psi_k, \quad k = 1, 2, \ldots, n$

- 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**

39

## Linear SVM: non-separable case

- Easier to solve the "dual" problem
  - **_Kuhn-Tucker construction_**

  **Problem 4**: Maximize $\sum_{k=1}^{n} \lambda_k - \frac{1}{2} \sum_{k,j}^{n} \lambda_k \lambda_j z_k z_j x_j^t x_k$

  subject to $\sum_{k=1}^{n} z_k \lambda_k = 0$ and $0 \leq \lambda_k \leq c, k = 1, 2, \ldots, n$

  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**

40

## Nonlinear SVM

- Nonlinear SVM
- Kernel Functions
- Example

FLIP: Pattern Recognition (G. Li)   **41**

41

## Nonlinear SVM

- Extending these concepts to the non-linear case involves mapping the data to a high-dimensional space $h$:

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

Linearly Separable in Higher Dimension

- 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**

42

## Slide 43

### Nonlinear SVM

- linear SVM:

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

- non-linear SVM:

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

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

FLIP: Pattern Recognition (G. Li)    43

43

## Slide 44

### Nonlinear SVM

- The ***disadvantage*** of this approach is that the mapping might be very computationally intensive to compute!

$$\mathbf{x}_k \to \Phi(\mathbf{x}_k)$$

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

FLIP: Pattern Recognition (G. Li)    44

44

## Slide 45

### Nonlinear SVM

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

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

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

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

FLIP: Pattern Recognition (G. Li)    45

45

## Slide 46

### Nonlinear SVM

- **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.

FLIP: Pattern Recognition (G. Li)    46

46

## Slide 47

### Polynomial Kernel

- $K(x,y) = (x.y)^d$

\* 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.

FLIP: Pattern Recognition (G. Li)    47

47

## Slide 48

### Choice of $\Phi$ is not unique

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

$$(x.y)^2 = (x_1 y_1 + x_2 y_2)^2$$

$$\Phi(x).\Phi(y) = x_1^2 y_1^2 + 2x_1 y_1 x_2 y_2 + x_2^2 y_2^2 = (x_1 y_1 + x_2 y_2)^2$$

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

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

FLIP: Pattern Recognition (G. Li)    48

48

## Slide 49

### Kernel functions

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

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

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

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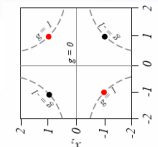(corresponds to a radial basis function (RBF) neural network)

FLIP: Pattern Recognition (G. Li)  49

49

## Slide 50

### 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 $\omega_2$



FLIP: Pattern Recognition (G. Li)  50

50

## Slide 51

### Example

- Consider the following mapping (among others)

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

FLIP: Pattern Recognition (G. Li)  51

51

## Slide 52

### Example

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

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

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

FLIP: Pattern Recognition (G. Li)  52

52

## Slide 53

### Example

- We seek to maximize

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

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

- 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 !

FLIP: Pattern Recognition (G. Li)  53

53

## Slide 54

### Example

- We now compute w

$$w = \sum_{k=1}^{4} z_k \lambda_k \Phi(x_k) = \frac{1}{8}\begin{pmatrix} 1 \\ \sqrt{2} \\ \sqrt{2} \\ \sqrt{2} \\ 1 \\ 1 \end{pmatrix} - \frac{1}{8}\begin{pmatrix} 1 \\ \sqrt{2} \\ -\sqrt{2} \\ -\sqrt{2} \\ 1 \\ 1 \end{pmatrix} + \frac{1}{8}\begin{pmatrix} 1 \\ -\sqrt{2} \\ \sqrt{2} \\ -\sqrt{2} \\ 1 \\ 1 \end{pmatrix} - \frac{1}{8}\begin{pmatrix} 1 \\ -\sqrt{2} \\ -\sqrt{2} \\ \sqrt{2} \\ 1 \\ 1 \end{pmatrix} = \frac{1}{2}\begin{pmatrix} 0 \\ 0 \\ \sqrt{2} \\ 0 \\ 0 \\ 0 \end{pmatrix}$$
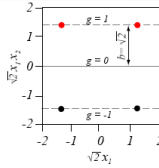
- The solution for $w_0$ can be determined using any support vector, e.g., $x_1$:

$$w^t \Phi(x_1) + w_0 = z_1 \quad \text{or} \quad w0=0$$

FLIP: Pattern Recognition (G. Li)  54

54

## Example

- The margin b is computed as follows

$$\mathbf{w} = \frac{1}{2} \begin{pmatrix} 0 \\ 0 \\ \sqrt{2} \\ 0 \\ 0 \\ 0 \end{pmatrix} \qquad b = \frac{1}{\|w\|} = \sqrt{2}$$
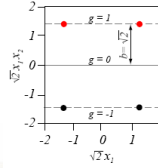
55

## Example

- The discriminant function is as follows

$$g(x) = w^t \Phi(x) + w_0 = x_1 x_2$$

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

$$\mathbf{w} = \frac{1}{2} \begin{pmatrix} 0 \\ 0 \\ \sqrt{2} \\ 0 \\ 0 \\ 0 \end{pmatrix} \qquad \Phi(x) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1 \\ \sqrt{2}x_1 x_2 \\ \sqrt{2}x_2 \\ x_2^2 \\ 1 \end{pmatrix}$$
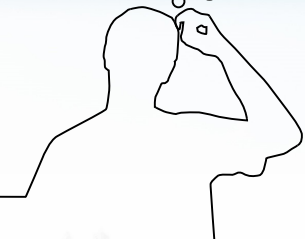
56

## 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.

57

## Questions?

59