

# Outline

- Basic concepts: SVM and kernels
- SVM primal/dual problems



# Outline

- Basic concepts: SVM and kernels
- SVM primal/dual problems



# Data Classification

- Given training data in different classes (labels **known**)  
Predict test data (labels **unknown**)
- Training and testing



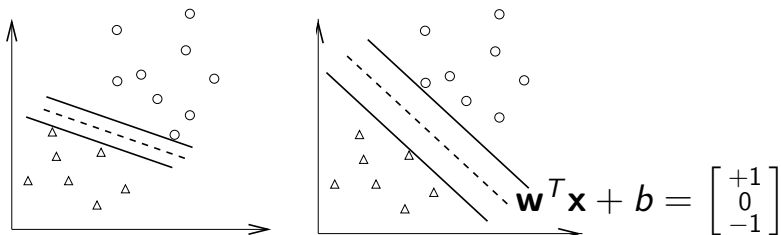
# Support Vector Classification

- **Training** vectors :  $\mathbf{x}_i, i = 1, \dots, l$
- Feature vectors. For example,  
A patient = [height, weight, ...]<sup>T</sup>
- Consider a simple case with **two classes**:  
Define an **indicator** vector  $\mathbf{y}$

$$y_i = \begin{cases} 1 & \text{if } \mathbf{x}_i \text{ in class 1} \\ -1 & \text{if } \mathbf{x}_i \text{ in class 2} \end{cases}$$

- A hyperplane which separates all data





- A separating hyperplane:  $\mathbf{w}^T \mathbf{x} + b = 0$

$$\begin{aligned} (\mathbf{w}^T \mathbf{x}_i) + b &\geq 1 && \text{if } y_i = 1 \\ (\mathbf{w}^T \mathbf{x}_i) + b &\leq -1 && \text{if } y_i = -1 \end{aligned}$$

- Decision function  $f(\mathbf{x}) = \text{sgn}(\mathbf{w}^T \mathbf{x} + b)$ ,  $\mathbf{x}$ : test data  
 Many possible choices of  $\mathbf{w}$  and  $b$



# Maximal Margin

- Distance between  $\mathbf{w}^T \mathbf{x} + b = 1$  and  $-1$ :

$$2/\|\mathbf{w}\| = 2/\sqrt{\mathbf{w}^T \mathbf{w}}$$

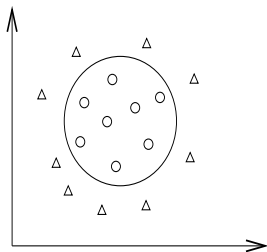
- A **quadratic programming** problem (Boser et al., 1992)

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{subject to} \quad & y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \\ & i = 1, \dots, l. \end{aligned}$$



# Data May Not Be Linearly Separable

- An example:



- Allow training errors
- Higher dimensional (maybe infinite) feature space

$$\phi(\mathbf{x}) = [\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots]^T.$$



- Standard SVM (Boser et al., 1992; Cortes and Vapnik, 1995)

$$\min_{\mathbf{w}, b, \xi} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^l \xi_i$$

subject to  $y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1 - \xi_i,$   
 $\xi_i \geq 0, \quad i = 1, \dots, l.$

- Example:  $\mathbf{x} \in R^3, \phi(\mathbf{x}) \in R^{10}$

$$\phi(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_3, x_1^2, x_2^2, x_3^2, \sqrt{2}x_1x_2, \sqrt{2}x_1x_3, \sqrt{2}x_2x_3]^T$$





# Finding the Decision Function

- $\mathbf{w}$ : maybe **infinite** variables
- The **dual** problem: **finite** number of variables

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^T Q \alpha - \mathbf{e}^T \alpha \\ \text{subject to} \quad & 0 \leq \alpha_i \leq C, i = 1, \dots, l \\ & \mathbf{y}^T \alpha = 0, \end{aligned}$$

where  $Q_{ij} = y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$  and  $\mathbf{e} = [1, \dots, 1]^T$

- At optimum

$$\mathbf{w} = \sum_{i=1}^l \alpha_i y_i \phi(\mathbf{x}_i)$$

- A **finite** problem: #variables = #training data



# Kernel Tricks

- $Q_{ij} = y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$  needs a **closed** form
- Example:  $\mathbf{x}_i \in R^3, \phi(\mathbf{x}_i) \in R^{10}$

$$\phi(\mathbf{x}_i) = [1, \sqrt{2}(x_i)_1, \sqrt{2}(x_i)_2, \sqrt{2}(x_i)_3, (x_i)_1^2, (x_i)_2^2, (x_i)_3^2, \sqrt{2}(x_i)_1(x_i)_2, \sqrt{2}(x_i)_1(x_i)_3, \sqrt{2}(x_i)_2(x_i)_3]^T$$

Then  $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2$ .

- Kernel:  $K(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^T \phi(\mathbf{y})$ ; common kernels:

$$e^{-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2}, \text{ (Radial Basis Function)}$$

$$(\mathbf{x}_i^T \mathbf{x}_j / a + b)^d \text{ (Polynomial kernel)}$$



Can be inner product in **infinite** dimensional space

Assume  $x \in R^1$  and  $\gamma > 0$ .

$$\begin{aligned}
 e^{-\gamma \|x_i - x_j\|^2} &= e^{-\gamma(x_i - x_j)^2} = e^{-\gamma x_i^2 + 2\gamma x_i x_j - \gamma x_j^2} \\
 &= e^{-\gamma x_i^2 - \gamma x_j^2} \left( 1 + \frac{2\gamma x_i x_j}{1!} + \frac{(2\gamma x_i x_j)^2}{2!} + \frac{(2\gamma x_i x_j)^3}{3!} + \dots \right) \\
 &= e^{-\gamma x_i^2 - \gamma x_j^2} \left( 1 \cdot 1 + \sqrt{\frac{2\gamma}{1!}} x_i \cdot \sqrt{\frac{2\gamma}{1!}} x_j + \sqrt{\frac{(2\gamma)^2}{2!}} x_i^2 \cdot \sqrt{\frac{(2\gamma)^2}{2!}} x_j^2 \right. \\
 &\quad \left. + \sqrt{\frac{(2\gamma)^3}{3!}} x_i^3 \cdot \sqrt{\frac{(2\gamma)^3}{3!}} x_j^3 + \dots \right) = \phi(x_i)^T \phi(x_j),
 \end{aligned}$$

where

$$\phi(x) = e^{-\gamma x^2} \left[ 1, \sqrt{\frac{2\gamma}{1!}} x, \sqrt{\frac{(2\gamma)^2}{2!}} x^2, \sqrt{\frac{(2\gamma)^3}{3!}} x^3, \dots \right]^T.$$



# Decision function

- At optimum

$$\mathbf{w} = \sum_{i=1}^l \alpha_i y_i \phi(\mathbf{x}_i)$$

- Decision function

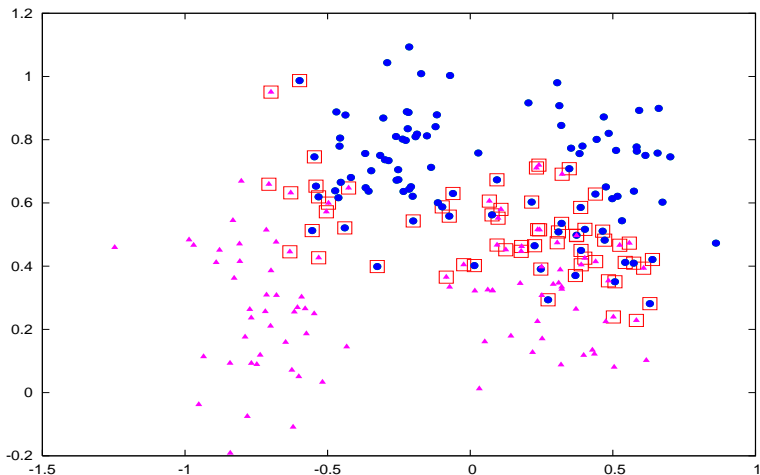
$$\begin{aligned} & \mathbf{w}^T \phi(\mathbf{x}) + b \\ &= \sum_{i=1}^l \alpha_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x}) + b \\ &= \sum_{i=1}^l \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b \end{aligned}$$

- Only  $\phi(\mathbf{x}_i)$  of  $\alpha_i > 0$  used  $\Rightarrow$  **support vectors**



# Support Vectors: More Important Data

Only  $\phi(\mathbf{x}_i)$  of  $\alpha_i > 0$  used  $\Rightarrow$  support vectors



# Outline

- Basic concepts: SVM and kernels
- SVM primal/dual problems



# Deriving the Dual

- Consider the problem without  $\xi_i$

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{subject to} \quad & y_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1, \quad i = 1, \dots, l. \end{aligned}$$

- Its dual

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^T Q \alpha - \mathbf{e}^T \alpha \\ \text{subject to} \quad & 0 \leq \alpha_i, \quad i = 1, \dots, l, \\ & \mathbf{y}^T \alpha = 0. \end{aligned}$$



# Lagrangian Dual

$$\max_{\alpha \geq 0} (\min_{\mathbf{w}, b} L(\mathbf{w}, b, \alpha)),$$

where

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^l \alpha_i (y_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b) - 1)$$

Strong duality

$$\min \text{ Primal} = \max_{\alpha \geq 0} (\min_{\mathbf{w}, b} L(\mathbf{w}, b, \alpha))$$





- Simplify the dual. **When  $\alpha$  is fixed,**

$$\min_{\mathbf{w}, b} L(\mathbf{w}, b, \alpha) =$$

$$\begin{cases} -\infty & \text{if } \sum_{i=1}^l \alpha_i y_i \neq 0 \\ \min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^l \alpha_i [y_i (\mathbf{w}^T \phi(\mathbf{x}_i) - 1)] & \text{if } \sum_{i=1}^l \alpha_i y_i = 0 \end{cases}$$

- If  $\sum_{i=1}^l \alpha_i y_i \neq 0$ ,  
decrease

$$-b \sum_{i=1}^l \alpha_i y_i$$

in  $L(\mathbf{w}, b, \alpha)$  to  $-\infty$



- If  $\sum_{i=1}^l \alpha_i y_i = 0$ , optimum of the **strictly convex**  $\frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^l \alpha_i [y_i (\mathbf{w}^T \phi(\mathbf{x}_i) - 1)]$  happens when

$$\frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}, b, \alpha) = 0.$$

- Thus,

$$\mathbf{w} = \sum_{i=1}^l \alpha_i y_i \phi(\mathbf{x}_i).$$



- Note that

$$\begin{aligned} \mathbf{w}^T \mathbf{w} &= \left( \sum_{i=1}^l \alpha_i y_i \phi(\mathbf{x}_i) \right)^T \left( \sum_{j=1}^l \alpha_j y_j \phi(\mathbf{x}_j) \right) \\ &= \sum_{i,j} \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \end{aligned}$$

- The dual is

$$\max_{\alpha \geq 0} \begin{cases} \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) & \text{if } \sum_{i=1}^l \alpha_i y_i = 0, \\ -\infty & \text{if } \sum_{i=1}^l \alpha_i y_i \neq 0. \end{cases}$$



- Lagrangian dual:  $\max_{\alpha \geq 0} (\min_{\mathbf{w}, b} L(\mathbf{w}, b, \alpha))$
  - $-\infty$  definitely **not** maximum of the dual
- Dual optimal solution not happen when

$$\sum_{i=1}^l \alpha_i y_i \neq 0$$

- Dual simplified to

$$\begin{aligned} \max_{\alpha \in R^l} \quad & \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \\ \text{subject to} \quad & \mathbf{y}^T \boldsymbol{\alpha} = 0, \\ & \alpha_i \geq 0, i = 1, \dots, l. \end{aligned}$$



- Our problems may be **infinite** dimensional
  - Can still use Lagrangian duality
- See a rigorous discussion in Lin (2001)



# References I

- B. E. Boser, I. Guyon, and V. Vapnik. A training algorithm for optimal margin classifiers. In *Proceedings of the Fifth Annual Workshop on Computational Learning Theory*, pages 144–152. ACM Press, 1992.
- C. Cortes and V. Vapnik. Support-vector network. *Machine Learning*, 20:273–297, 1995.
- C.-J. Lin. Formulations of support vector machines: a note from an optimization point of view. *Neural Computation*, 13(2):307–317, 2001.

