

# Machine Learning (機器學習)

Lecture 11: Support Vector Machine (2)

Hsuan-Tien Lin (林軒田)

htlin@csie.ntu.edu.tw

Department of Computer Science  
& Information Engineering

National Taiwan University  
(國立台灣大學資訊工程系)



# Roadmap

- ① When Can Machines Learn?
- ② Why Can Machines Learn?
- ③ How Can Machines Learn?
- ④ How Can Machines Learn Better?
- ⑤ Embedding Numerous Features: Kernel Models

## Lecture 11: Support Vector Machine (2)

- Kernel Trick
- Polynomial Kernel
- Gaussian Kernel
- Comparison of Kernels
- Motivation and Primal Problem
- Dual Problem
- Messages behind Soft-Margin SVM
- Soft-Margin SVM as Regularized Model

# Dual SVM Revisited

goal: SVM **without dependence on  $\tilde{d}$**

half-way done:

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^T Q_D \alpha - \mathbf{1}^T \alpha \\ \text{subject to} \quad & \mathbf{y}^T \alpha = 0; \\ & \alpha_n \geq 0, \text{ for } n = 1, 2, \dots, N \end{aligned}$$

- $q_{n,m} = y_n y_m \mathbf{z}_n^T \mathbf{z}_m$ : inner product in  $\mathbb{R}^{\tilde{d}}$
- need:  $\mathbf{z}_n^T \mathbf{z}_m = \Phi(\mathbf{x}_n)^T \Phi(\mathbf{x}_m)$  calculated faster than  $O(\tilde{d})$

$\uparrow$   $\tilde{d} \gg k$ , this step become bottleneck.

can we do so?

# Fast Inner Product for $\Phi_2$

2nd order polynomial transform

$$\Phi_2(\mathbf{x}) = (1, x_1, x_2, \dots, x_d, x_1^2, x_1x_2, \dots, x_1x_d, x_2x_1, x_2^2, \dots, x_2x_d, \dots, x_d^2)$$

—include both  $x_1x_2$  &  $x_2x_1$  for '**simplicity**' :-)

$C_Q^{Q+d-1}$

先轉換再內積

$$\begin{aligned}
 \Phi_2(\mathbf{x})^T \Phi_2(\mathbf{x}') &= 1 + \sum_{i=1}^d x_i x'_i + \sum_{i=1}^d \sum_{j=1}^d x_i x_j x'_i x'_j \\
 &= 1 + \sum_{i=1}^d x_i x'_i + \sum_{i=1}^d x_i x'_i \sum_{j=1}^d x_j x'_j \\
 &= 1 + \underline{\mathbf{x}}^T \underline{\mathbf{x}'} + (\underline{\mathbf{x}}^T \underline{\mathbf{x}'}) (\underline{\mathbf{x}}^T \underline{\mathbf{x}'})
 \end{aligned}$$

$x_1 x_2 \dots x_d x_d$   
 $x_i x_j x'_i x'_j$   
 $x_j x'_j$   
 $\uparrow$  先內積再平方

for  $\Phi_2$ , transform + inner product can be carefully done in  $O(d)$  instead of  $O(d^2)$

# Kernel: Transform + Inner Product

transform  $\Phi \iff \text{kernel function: } K_\Phi(\mathbf{x}, \mathbf{x}') \equiv \Phi(\mathbf{x})^T \Phi(\mathbf{x}')$

$$\Phi_2 \iff K_{\Phi_2}(\mathbf{x}, \mathbf{x}') = 1 + (\mathbf{x}^T \mathbf{x}') + (\mathbf{x}^T \mathbf{x}')^2 \leftarrow \text{類似於軟件}$$

- quadratic coefficient  $q_{n,m} = y_n y_m \mathbf{z}_n^T \mathbf{z}_m = y_n y_m K(\mathbf{x}_n, \mathbf{x}_m)$

- optimal bias  $b$ ? from SV  $(\mathbf{x}_s, y_s)$ ,

給定一個SV算 $b$

$$b = \mathbf{y}_s - \mathbf{w}^T \mathbf{z}_s = \mathbf{y}_s - \left( \sum_{n=1}^N \alpha_n y_n \mathbf{z}_n \right)^T \mathbf{z}_s = \mathbf{y}_s - \sum_{n=1}^N \alpha_n y_n K(\mathbf{x}_n, \mathbf{x}_s)$$

展開 $\mathbf{w}$

- optimal hypothesis  $g_{\text{SVM}}$ : for test input  $\mathbf{x}$ ,

$$g_{\text{SVM}}(\mathbf{x}) = \text{sign}(\mathbf{w}^T \Phi(\mathbf{x}) + b) = \text{sign}\left(\sum_{n=1}^N \alpha_n y_n K(\mathbf{x}_n, \mathbf{x}) + b\right)$$

展開 $\mathbf{w}$

結論 → 已全部都代換成 $K$ 了!

kernel trick: plug in efficient kernel function  
to avoid dependence on  $\tilde{d}$

## Kernel SVM with QP

## Kernel Hard-Margin SVM Algorithm

- ①  $q_{n,m} = y_n y_m K(\mathbf{x}_n, \mathbf{x}_m); \mathbf{p} = -\mathbf{1}_N; (\mathbf{A}, \mathbf{c})$  for equ./bound constraints
- ②  $\alpha \leftarrow \text{QP}(\mathbf{Q}_D, \mathbf{p}, \mathbf{A}, \mathbf{c})$
- ③  $b \leftarrow \left( y_s - \sum_{\substack{\text{SV indices } n \\ \text{SV}}} \alpha_n y_n K(\mathbf{x}_n, \mathbf{x}_s) \right)$  with SV  $(\mathbf{x}_s, y_s)$
- ④ return SVs and their  $\alpha_n$  as well as  $b$  such that for new  $\mathbf{x}$ ,  

$$g_{\text{SVM}}(\mathbf{x}) = \text{sign} \left( \sum_{\substack{\text{SV indices } n \\ \text{SV}}} \alpha_n y_n K(\mathbf{x}_n, \mathbf{x}) + b \right)$$

- ①: time complexity  $O(N^2) \cdot (\text{kernel evaluation})$
- ②: QP with  $N$  variables and  $N + 1$  constraints
- ③ & ④: time complexity  $O(\#\text{SV}) \cdot (\text{kernel evaluation})$

都跟  $\tilde{d}$   
无关。

## kernel SVM:

use computational shortcut to avoid  $\tilde{d}$  & predict with SV only

# Questions?

## General Poly-2 Kernel

不同的2次 kernel.

$$\Phi_2(\mathbf{x}) = (1, x_1, \dots, x_d, x_1^2, \dots, x_d^2) \Leftrightarrow K_{\Phi_2}(\mathbf{x}, \mathbf{x}') = 1 + \mathbf{x}^T \mathbf{x}' + (\mathbf{x}^T \mathbf{x}')^2$$

$$\Phi_2(\mathbf{x}) = (1, \sqrt{2}x_1, \dots, \sqrt{2}x_d, x_1^2, \dots, x_d^2) \Leftrightarrow K_2(\mathbf{x}, \mathbf{x}') = 1 + 2\mathbf{x}^T \mathbf{x}' + (\mathbf{x}^T \mathbf{x}')^2$$

加常数.

$$\Phi_2(\mathbf{x}) = (1, \sqrt{2\gamma}x_1, \dots, \sqrt{2\gamma}x_d, \gamma x_1^2, \dots, \gamma x_d^2)$$

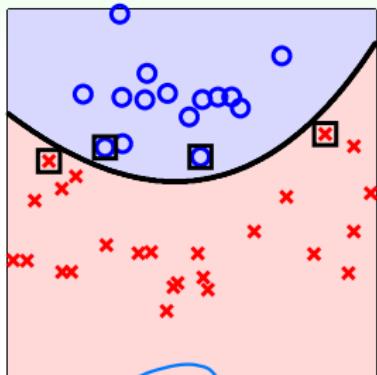
$$\Leftrightarrow K_2(\mathbf{x}, \mathbf{x}') = 1 + (2\gamma)\mathbf{x}^T \mathbf{x}' + (\gamma^2)(\mathbf{x}^T \mathbf{x}')^2$$

$$K_2(\mathbf{x}, \mathbf{x}') = (1 + \gamma \mathbf{x}^T \mathbf{x}')^2 \text{ with } \gamma > 0$$

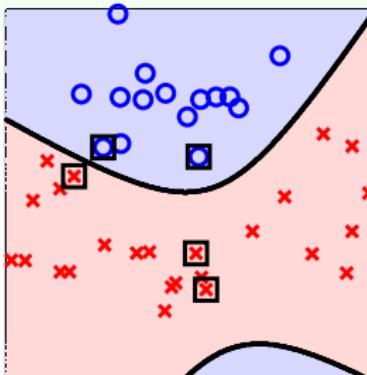
- $K_2$ : somewhat 'easier' to calculate than  $K_{\Phi_2}$
- $\Phi_2$  and  $\Phi_2$ : equivalent **power**,  
different inner product  $\Rightarrow$  different **geometry**

$K_2$  commonly used

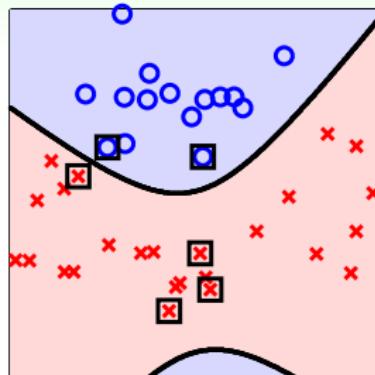
# Poly-2 Kernels in Action



$$(1 + \textcolor{green}{0.001} \mathbf{x}^T \mathbf{x}')^2$$



$$1 + \mathbf{x}^T \mathbf{x}' + (\mathbf{x}^T \mathbf{x}')^2$$



$$(1 + \textcolor{green}{1000} \mathbf{x}^T \mathbf{x}')^2$$

- $g_{\text{SVM}}$  different, SVs different ← 这一不同  
—‘hard’ to say which is better before learning
- change of kernel ⇔ change of margin definition  
↑ kernel很重要.

need selecting  $K$ , just like selecting  $\Phi$

# General Polynomial Kernel

$$K_2(\mathbf{x}, \mathbf{x}') = (\zeta + \gamma \mathbf{x}^T \mathbf{x}')^2 \text{ with } \gamma > 0, \zeta \geq 0$$

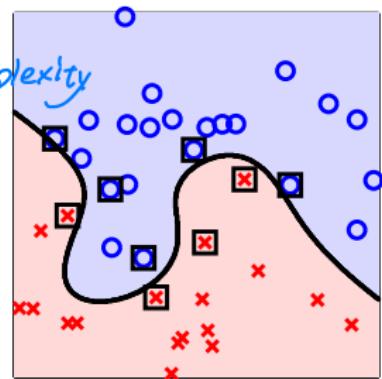
$$K_3(\mathbf{x}, \mathbf{x}') = (\zeta + \gamma \mathbf{x}^T \mathbf{x}')^3 \text{ with } \gamma > 0, \zeta \geq 0$$

⋮

$$K_Q(\mathbf{x}, \mathbf{x}') = (\zeta + \gamma \mathbf{x}^T \mathbf{x}')^Q \text{ with } \gamma > 0, \zeta \geq 0$$

- embeds  $\Phi_Q$  specially with parameters  
 $(\gamma, \zeta) \leftarrow$  可調參數 幫助你控制 complexity
- allows computing large-margin polynomial classification without dependence on  $\tilde{d}$

SVM + Polynomial Kernel: Polynomial SVM



10-th order polynomial  
 $Q=10$  with margin 0.1

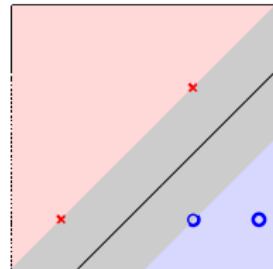
## Special Case: Linear Kernel

$$K_1(\mathbf{x}, \mathbf{x}') = (0 + 1 \cdot \mathbf{x}^T \mathbf{x}')^1 \leftarrow \text{linear Kernel}$$

⋮

$$K_Q(\mathbf{x}, \mathbf{x}') = (\zeta + \gamma \mathbf{x}^T \mathbf{x}')^Q \text{ with } \gamma > 0, \zeta \geq 0$$

- $K_1$ : just **usual inner product**, called **linear kernel**.
- ‘even easier’: can be solved (often in primal form) **efficiently**



try linear first  
linear first remember? :-)

# Questions?

# Kernel of Infinite Dimensional Transform

infinite dimensional  $\Phi(\mathbf{x})$ ? Yes, if  $K(\mathbf{x}, \mathbf{x}')$  efficiently computable!

$$\begin{aligned}
 \text{when } \mathbf{x} = (x), K(x, x') &= \exp(-(x - x')^2) \quad \text{高斯函數.} \\
 &= \exp(-x^2) \exp(-x'^2) \exp(2xx') \\
 \stackrel{\text{Taylor}}{=} &\exp(-x^2) \exp(-x'^2) \left( \sum_{i=0}^{\infty} \frac{(2xx')^i}{i!} \right) \\
 &= \sum_{i=0}^{\infty} \left( \exp(-x^2) \exp(-x'^2) \sqrt{\frac{2^i}{i!}} \sqrt{\frac{2^i}{i!}} (x)^i (x')^i \right) \\
 &= \Phi(x)^T \Phi(x')
 \end{aligned}$$

with infinite dimensional  $\Phi(x) = \exp(-x^2) \cdot \left( 1, \sqrt{\frac{2}{1!}}x, \sqrt{\frac{2^2}{2!}}x^2, \dots \right)$

$\downarrow \mathbf{x}$  處向量  $\infty$  dimension

more generally, Gaussian kernel

$$K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2) \text{ with } \gamma > 0$$

# Hypothesis of Gaussian SVM

$$\text{Gaussian kernel } K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$$

$$\begin{aligned} g_{\text{SVM}}(\mathbf{x}) &= \text{sign} \left( \sum_{\text{SV}} \alpha_n y_n \underbrace{K(\mathbf{x}_n, \mathbf{x})}_{\text{高斯}} + b \right) \\ &= \text{sign} \left( \sum_{\text{SV}} \alpha_n y_n \exp \left( -\gamma \|\mathbf{x} - \mathbf{x}_n\|^2 \right) + b \right) \end{aligned}$$

- linear combination of Gaussians centered at SVs  $\mathbf{x}_n$
- also called Radial Basis Function (RBF) kernel

(高斯基底)

Gaussian SVM:

find  $\alpha_n$  to combine Gaussians centered at  $\mathbf{x}_n$   
& achieve large margin in infinite-dim. space

# Support Vector Mechanism

	<b>large-margin hyperplanes</b>
	+ higher-order transforms with <b>kernel trick</b>
# boundary	<b>not many</b> <b>sophisticated</b>

不做 Z, 改做 K

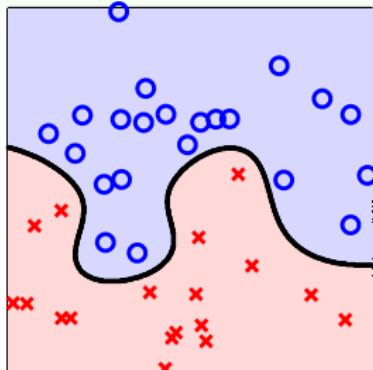
- transformed vector  $\mathbf{z} = \Phi(\mathbf{x}) \Rightarrow$  efficient kernel  $K(\mathbf{x}, \mathbf{x}')$
- store optimal  $\mathbf{w} \Rightarrow$  store **a few SVs** and  $\alpha_n$

不存, 存 SVs,  $\alpha_n$ .

new possibility by Gaussian SVM:  
infinite-dimensional linear classification, with  
 generalization ‘guarded by’ large-margin :-)

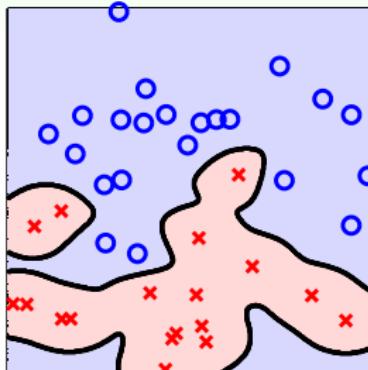
# Gaussian SVM in Action

$\gamma \Rightarrow \text{Gaussian} \text{ 放大了}$



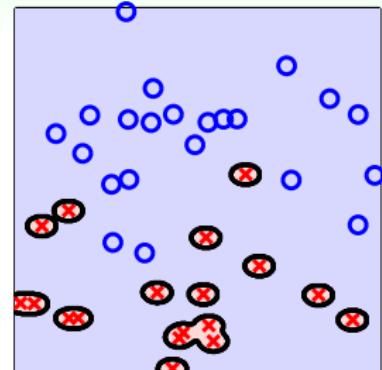
$$\exp(-1\|\mathbf{x} - \mathbf{x}'\|^2)$$

$\gamma = 1$



$$\exp(-10\|\mathbf{x} - \mathbf{x}'\|^2)$$

$\gamma = 10$



$$\exp(-100\|\mathbf{x} - \mathbf{x}'\|^2)$$

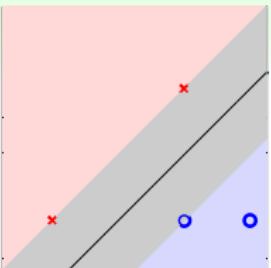
$\gamma = 100$

- large  $\gamma \Rightarrow$  sharp Gaussians  $\Rightarrow$  'overfit'?
- **warning: SVM can still overfit :-)**

Gaussian SVM: need careful selection of  $\gamma$

# Questions?

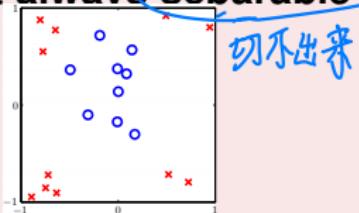
# Linear Kernel: Cons and Pros



$$K(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$$

## Cons

- restricted  
—not always separable?

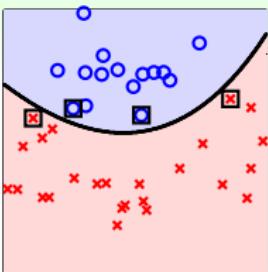


## Pros

- safe—**linear first, remember? :-)**
- fast—with **special QP solver** in primal
- very explainable—**w and SVs** say something

linear kernel: an important **basic** tool

# Polynomial Kernel: Cons and Pros



$$K(\mathbf{x}, \mathbf{x}') = (\zeta + \gamma \mathbf{x}^T \mathbf{x}')^Q$$

## Cons

- numerical difficulty for large  $Q$  次方項太多
  - three parameters ( $\gamma, \zeta, Q$ ) —more difficult to select
- 大 很多參數要選.

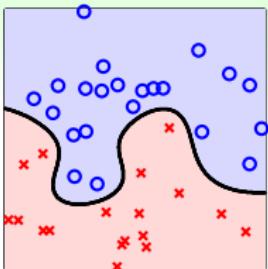
## Pros

- less restricted than linear
- strong physical control —'knows' degree  $Q$  人為設定  $Q$ .

eg.  $Q=2, 3$

polynomial kernel: perhaps small- $Q$  only  
—sometimes efficiently done by **linear on  $\Phi_Q(\mathbf{x})$**

# Gaussian Kernel: Cons and Pros

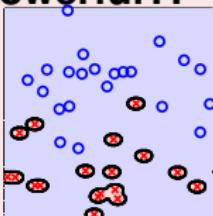


$$K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$$

无限多次方

## Cons

- mysterious—no w 解釋性低
- slower than linear
- too powerful?!



## Pros

- more powerful than linear/poly.
- bounded—less numerical difficulty than poly.
- one parameter only—easier to select than poly.

Gaussian kernel: **one of most popular** but shall **be used with care**

# Other Valid Kernels

- kernel represents **special** similarity:  $\Phi(\mathbf{x})^T \Phi(\mathbf{x}')$  的相似性
- any similarity  $\Rightarrow$  valid kernel? **not really** 不見得
- necessary & sufficient conditions for valid kernel:  
**Mercer's condition** 充份且必要 成為 Kernel 的條件
  - symmetric
  - let  $k_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$ , the matrix  $K$

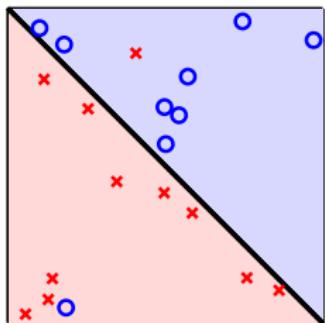
$$\begin{aligned}
 &= \begin{bmatrix} \Phi(\mathbf{x}_1)^T \Phi(\mathbf{x}_1) & \Phi(\mathbf{x}_1)^T \Phi(\mathbf{x}_2) & \dots & \Phi(\mathbf{x}_1)^T \Phi(\mathbf{x}_N) \\ \Phi(\mathbf{x}_2)^T \Phi(\mathbf{x}_1) & \Phi(\mathbf{x}_2)^T \Phi(\mathbf{x}_2) & \dots & \Phi(\mathbf{x}_2)^T \Phi(\mathbf{x}_N) \\ \vdots & \ddots & \ddots & \vdots \\ \Phi(\mathbf{x}_N)^T \Phi(\mathbf{x}_1) & \Phi(\mathbf{x}_N)^T \Phi(\mathbf{x}_2) & \dots & \Phi(\mathbf{x}_N)^T \Phi(\mathbf{x}_N) \end{bmatrix} \\
 &= [\mathbf{z}_1 \ \mathbf{z}_2 \ \dots \ \mathbf{z}_N]^T [\mathbf{z}_1 \ \mathbf{z}_2 \ \dots \ \mathbf{z}_N] \quad \text{平方}
 \end{aligned}$$

define your own kernel: possible, **but hard**

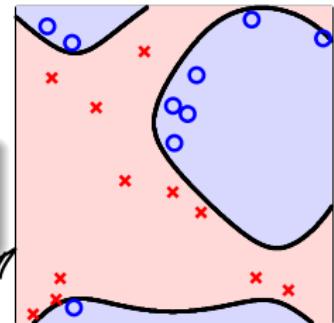
# Questions?

# Cons of Hard-Margin SVM

recall: SVM can still overfit :-(

 $\Phi_1$ 

- part of reasons:  $\Phi$
- other part: **separable**

 $\Phi_4$ 

overfit

if always insisting on **separable** ( $\Rightarrow$  shatter),  
have power to **overfit to noise**

# Give Up on Some Examples

want: **give up** on some noisy examples

## min.-error perceptron

$$\min_{b, w} \sum_{n=1}^N \left[ y_n \neq \text{sign}(w^T z_n + b) \right]$$

*pocket*

做錯的愈少愈好.

## hard-margin SVM

$$\begin{aligned} \min_{b, w} & \frac{1}{2} w^T w \\ \text{s.t. } & y_n(w^T z_n + b) \geq 1 \text{ for all } n \end{aligned}$$

combination:

$$\begin{aligned} \min_{b, w} & \frac{1}{2} w^T w + C \sum_{n=1}^N \left[ y_n \neq \text{sign}(w^T z_n + b) \right] \\ \text{s.t. } & y_n(w^T z_n + b) \geq 1 \text{ for correct } n \\ & y_n(w^T z_n + b) \geq -\infty \text{ for incorrect } n \end{aligned}$$

*權* *做錯的* *數*.

**C**: trade-off of large margin & noise tolerance

# Soft-Margin SVM (1/2)

$$\begin{array}{ll} \min_{b, w} & \frac{1}{2} w^T w + C \cdot \sum_{n=1}^N \llbracket y_n \neq \text{sign}(w^T z_n + b) \rrbracket \\ \text{s.t.} & y_n(w^T z_n + b) \geq 1 - \infty \cdot \llbracket y_n \neq \text{sign}(w^T z_n + b) \rrbracket \end{array}$$

如果做錯的話，不計算這不等式

*boolean.*

- $\llbracket \cdot \rrbracket$ : non-linear, not QP anymore :-(
  - what about dual? kernel?

*無法分析 error 的大小*
- cannot distinguish **small error** (slightly away from fat boundary)  
or **large error** (a...w...a...y... from fat boundary)

- record '**margin violation**' by  $\xi_n$  — **linear constraints**

- penalize with **margin violation** instead of **error count**  
 — **quadratic objective**

soft-margin SVM:  $\min_{b, w, \xi} \frac{1}{2} w^T w + C \cdot \sum_{n=1}^N \xi_n$

s.t.  $y_n(w^T z_n + b) \geq 1 - \xi_n$  and  $\xi_n \geq 0$  for all  $n$

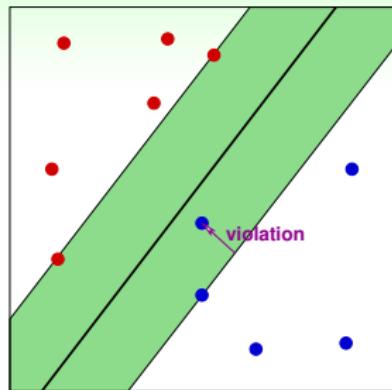
## Soft-Margin SVM (2/2)

- record 'margin violation' by  $\xi_n$
- penalize with margin violation

$$\min_{b, \mathbf{w}, \xi} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N \xi_n$$

控制 penalty

s.t.  $y_n(\mathbf{w}^T \mathbf{z}_n + b) \geq 1 - \xi_n$  and  $\xi_n \geq 0$  for all  $n$



- parameter  $C$ : trade-off of large margin & margin violation
  - large  $C$ : want less margin violation
  - small  $C$ : want large margin
- QP of  $\tilde{d} + 1 + N$  variables,  $\underline{2N}$  constraints

next: remove dependence on  $\tilde{d}$  by  
soft-margin SVM primal  $\Rightarrow$  dual?

# Questions?

Lagrange Dual *拉格朗日对偶.*

primal:  $\min_{b, \mathbf{w}, \xi} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \cdot \sum_{n=1}^N \xi_n$

s.t.  $y_n(\mathbf{w}^T \mathbf{z}_n + b) \geq 1 - \xi_n$  and  $\xi_n \geq 0$  for all  $n$

Lagrange function with Lagrange multipliers  $\underline{\alpha_n}$  and  $\underline{\beta_n}$

$$\begin{aligned} \mathcal{L}(b, \mathbf{w}, \xi, \underline{\alpha}, \underline{\beta}) = & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \cdot \sum_{n=1}^N \xi_n \\ & + \sum_{n=1}^N \underline{\alpha_n} \cdot (1 - \xi_n - y_n(\mathbf{w}^T \mathbf{z}_n + b)) + \sum_{n=1}^N \underline{\beta_n} \cdot (-\xi_n) \end{aligned}$$

*Lagrange multiplier.*

want: Lagrange dual

$$\max_{\underline{\alpha_n} \geq 0, \underline{\beta_n} \geq 0} \left( \min_{b, \mathbf{w}, \xi} \mathcal{L}(b, \mathbf{w}, \xi, \underline{\alpha}, \underline{\beta}) \right)$$

Simplify  $\xi_n$  and  $\beta_n$ 

$$\max_{\alpha_n \geq 0, \beta_n \geq 0} \left( \min_{b, w, \xi} \frac{1}{2} w^T w + \cancel{\alpha} \cdot \sum_{n=1}^N \xi_n \right) \mathcal{U}_n.$$

✓  $\sum_{n=1}^N (\mathcal{U}_n - \alpha_n) \cdot -\xi_n$

$$+ \sum_{n=1}^N \alpha_n \cdot (1 - \xi_n - y_n(w^T z_n + b)) + \sum_{n=1}^N \beta_n \cdot (-\xi_n)$$

✓  $\uparrow$

- $\frac{\partial \mathcal{L}}{\partial \xi_n} = 0 = \cancel{\alpha} - \alpha_n - \beta_n \quad \alpha_n + \beta_n = C$
- no loss of optimality if solving with implicit constraint  $\beta_n = C - \alpha_n$  and explicit constraint  $0 \leq \alpha_n \leq C$ :  $\beta_n$  removed

$\xi$  can also be removed :-), like how we removed  $b$

$$\max_{0 \leq \alpha_n \leq C, \beta_n = C - \alpha_n} \left( \min_{b, w, \xi} \frac{1}{2} w^T w + \sum_{n=1}^N \alpha_n (1 - y_n(w^T z_n + b)) \right)$$

$\cancel{\beta_n} \rightarrow \xi_n$

$$+ \sum_{n=1}^N (C - \alpha_n - \beta_n) \cdot \xi_n$$

$\leftarrow$  把  $\beta_n$  去掉

# Other Simplifications

$$\max_{0 \leq \alpha_n \leq C, \beta_n = C - \alpha_n} \sum_{n=1}^N \alpha_n (1 - y_n (\mathbf{w}^T \mathbf{z}_n + b))$$

$$\min_{b, \mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

familiar? :-)

*PP hard-margin - f(z).*

- inner problem same as hard-margin SVM
- $\frac{\partial \mathcal{L}}{\partial b} = 0$ : no loss of optimality if solving with constraint  $\sum_{n=1}^N \alpha_n y_n = 0$
- $\frac{\partial \mathcal{L}}{\partial w_i} = 0$ : no loss of optimality if solving with constraint

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{z}_n$$

standard dual can be derived  
using the same steps as Lecture 10

## Standard Soft-Margin SVM Dual

$$\min_{\alpha} \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y_n y_m \mathbf{z}_n^T \mathbf{z}_m - \sum_{n=1}^N \alpha_n$$

subject to  $\sum_{n=1}^N y_n \alpha_n = 0;$  ← 1个 constraint

$0 \leq \alpha_n \leq C$  for  $n = 1, 2, \dots, N;$  ← 2N个 constraint

implicitly  $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{z}_n;$

$\beta_n = C - \alpha_n$  for  $n = 1, 2, \dots, N$

—only difference to hard-margin: upper bound on  $\alpha_n$

another (convex) QP,  
with **N variables** &  **$2N + 1$  constraints**

# Questions?

# Kernel Soft-Margin SVM

## Kernel Soft-Margin SVM Algorithm

- 1  $q_{n,m} = y_n y_m K(\mathbf{x}_n, \mathbf{x}_m); \mathbf{p} = -\mathbf{1}_N; (\mathbf{A}, \mathbf{c})$  for equ./lower-bound/upper-bound constraints
  - 2  $\alpha \leftarrow QP(Q_D, \mathbf{p}, \mathbf{A}, \mathbf{c})$  C
  - ~~3~~  $b \leftarrow ?$
  - 4 return SVs and their  $\alpha_n$  as well as  $b$  such that for new  $\mathbf{x}$ ,
- $$g_{\text{SVM}}(\mathbf{x}) = \text{sign} \left( \sum_{\text{SV indices } n} \alpha_n y_n K(\mathbf{x}_n, \mathbf{x}) + b \right)$$

- almost the same as hard-margin
- more flexible than hard-margin  
—primal/dual always solvable

How to get  $b$

remaining question: step 3?

# Solving for $b$

## hard-margin SVM

complementary slackness:

$$\alpha_n(1 - y_n(\mathbf{w}^T \mathbf{z}_n + b)) = 0$$

2者不能共存於世

- SV ( $\alpha_s > 0$ )  
 $\Rightarrow b = y_s - \mathbf{w}^T \mathbf{z}_s$

## soft-margin SVM

complementary slackness:

$$\alpha_n(1 - \xi_n - y_n(\mathbf{w}^T \mathbf{z}_n + b)) = 0$$

$$(C - \alpha_n)\xi_n = 0$$

- SV ( $\alpha_s > 0$ )  
 $\Rightarrow b = y_s - y_s \xi_s - \mathbf{w}^T \mathbf{z}_s$
- free ( $\alpha_s < C$ )  
 $\Rightarrow \xi_s = 0$

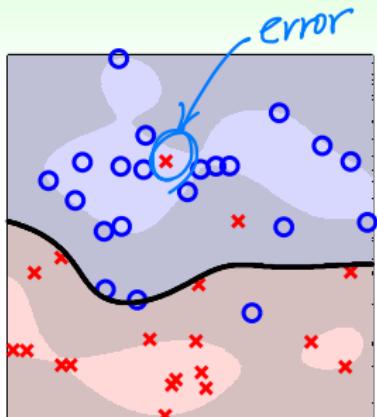
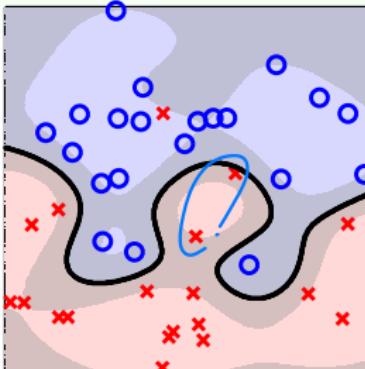
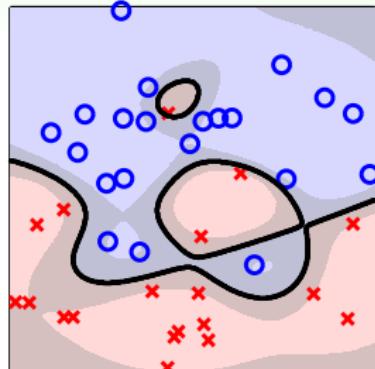
solve unique  $b$  with free SV ( $\mathbf{x}_s, y_s$ ):

任一free SV

$$b = y_s - \sum_{\text{SV indices } n} \alpha_n y_n K(\mathbf{x}_n, \mathbf{x}_s)$$

—range of  $b$  otherwise

# Soft-Margin Gaussian SVM in Action

 $C = 1$  $C = 10$  $C = 100$ 

- large  $C \Rightarrow$  less noise tolerance  $\Rightarrow$  'overfit'?
- warning: SVM can still overfit :-(

过参数

soft-margin Gaussian SVM:  
need careful selection of ( $\gamma$ ,  $C$ )

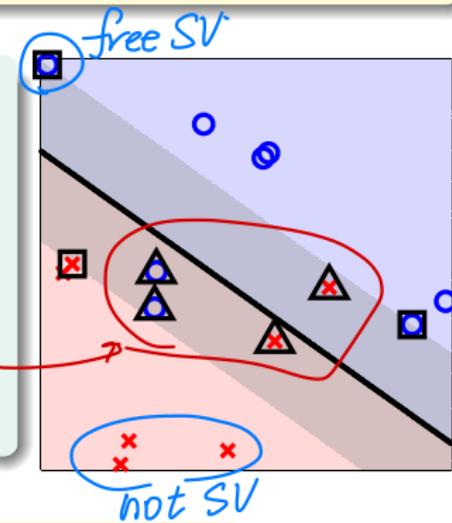
# Physical Meaning of $\alpha_n$

complementary slackness:

$$\left\{ \begin{array}{l} \alpha_n(1 - \xi_n - y_n(\mathbf{w}^T \mathbf{z}_n + b)) = 0 \\ (C - \alpha_n)\xi_n = 0 \end{array} \right.$$

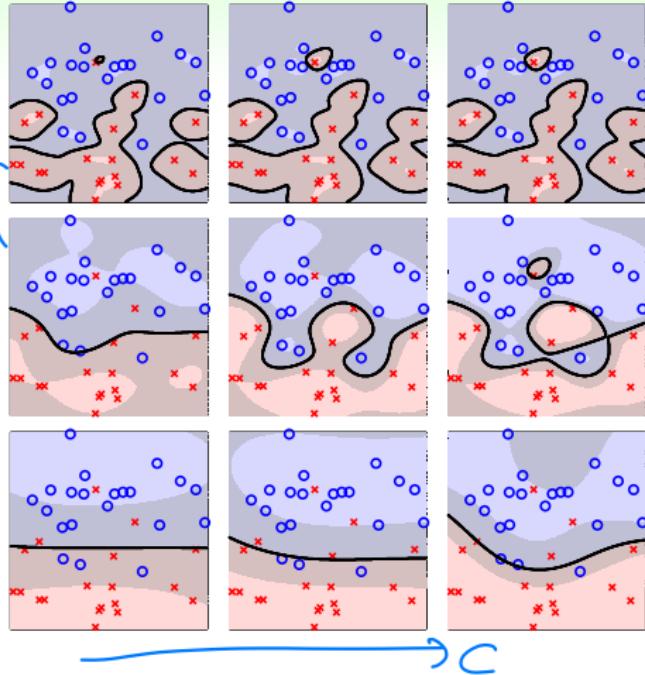
3種类用以分类

- non SV ( $0 = \alpha_n$ ):  $\xi_n = 0$ ,  
‘away from’/on fat boundary
- □ free SV ( $0 < \alpha_n < C$ ):  $\xi_n = 0$ , 找  $b$   
on fat boundary, locates  $b$
- △ bounded SV ( $\alpha_n = C$ ):  
 $\xi_n = \text{violation amount}$ , ← 違反的量。  
‘violate’/on fat boundary



$\alpha_n$  can be used for **data analysis**

# Practical Need: Model Selection



- complicated even for  $(C, \gamma)$  of Gaussian SVM
- more combinations if including other kernels or parameters

CROSS validation

how to select? **validation :-)**

# Questions?

## Wrap-Up

## Hard-Margin Primal

原始

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{s.t.} \quad & y_n (\mathbf{w}^T \mathbf{z}_n + b) \geq 1 \end{aligned}$$

## Soft-Margin Primal

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N \xi_n \\ \text{s.t.} \quad & y_n (\mathbf{w}^T \mathbf{z}_n + b) \geq 1 - \xi_n, \xi_n \geq 0 \end{aligned}$$

## Hard-Margin Dual

use  $\alpha$  to present

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^T Q \alpha - \mathbf{1}^T \alpha \\ \text{s.t.} \quad & \mathbf{y}^T \alpha = 0 \\ & 0 \leq \alpha_n \end{aligned}$$

## Soft-Margin Dual

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^T Q \alpha - \mathbf{1}^T \alpha \\ \text{s.t.} \quad & \mathbf{y}^T \alpha = 0 \\ & 0 \leq \alpha_n \leq C \end{aligned}$$

upper bound.

actually not one use hard-margin.

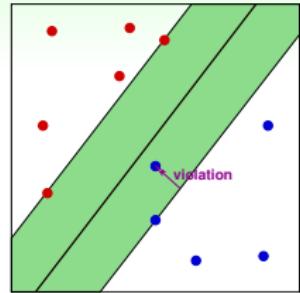
soft-margin preferred in practice;  
 linear: LIBLINEAR; non-linear: LIBSVM

# Slack Variables $\xi_n$

- record ‘margin violation’ by  $\xi_n$
- penalize with margin violation

$$\min_{b, \mathbf{w}, \xi} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \cdot \sum_{n=1}^N \xi_n$$

s.t.  $y_n(\mathbf{w}^T \mathbf{z}_n + b) \geq 1 - \xi_n$  and  $\xi_n \geq 0$  for all  $n$



margin violation amount

on any  $(b, \mathbf{w})$ ,  $\xi_n = \text{margin violation} = \max(1 - y_n(\mathbf{w}^T \mathbf{z}_n + b), 0)$

- {  $(\mathbf{x}_n, y_n)$  violating margin:  $\xi_n = 1 - y_n(\mathbf{w}^T \mathbf{z}_n + b)$  }
- {  $(\mathbf{x}_n, y_n)$  not violating margin:  $\xi_n = 0$  }

‘unconstrained’ form of soft-margin SVM:

*rewrite penalty*

$$\min_{b, \mathbf{w}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N \max(1 - y_n(\mathbf{w}^T \mathbf{z}_n + b), 0)$$

## Unconstrained Form

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N \max(1 - y_n(\mathbf{w}^T \mathbf{z}_n + b), 0)$$

familiar? :-)

$$\min \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum \widehat{\text{err}}$$

*Similar To L2 regularization*

just L2 regularization

$$\min \quad \frac{\lambda}{N} \mathbf{w}^T \mathbf{w} + \frac{1}{N} \sum \text{err}$$

with shorter  $\mathbf{w}$ , another parameter, and special err

why not solve this? :-)

- not QP, no (?) kernel trick
- $\max(\cdot, 0)$  not differentiable, harder to solve *Finge error* *凹函数*

# SVM as Regularized Model

	minimize	constraint
regularization by constraint	$E_{in}$	$\mathbf{w}^T \mathbf{w} \leq C$
hard-margin SVM	$\mathbf{w}^T \mathbf{w}$	$E_{in} = 0$ [and more]
L2 regularization	$\frac{\lambda}{N} \mathbf{w}^T \mathbf{w} + E_{in}$	
soft-margin SVM	$\frac{1}{2} \mathbf{w}^T \mathbf{w} + \widehat{CNE_{in}}$	

large margin  $\iff$  fewer hyperplanes  $\iff$  L2 regularization of short  $\mathbf{w}$

soft margin  $\iff$  special  $\widehat{err}$

larger  $C$  or  $\widehat{C}$   $\iff$  smaller  $\lambda$   $\iff$  less regularization

viewing SVM as regularized model:

allows extending/connecting to other learning models

# Questions?

# Summary

## ① Embedding Numerous Features: Kernel Models

### Lecture 11: Support Vector Machine (2)

- Kernel Trick  
**kernel as shortcut of transform + inner product**
- Polynomial Kernel  
**embeds specially-scaled polynomial transform**
- Gaussian Kernel  
**embeds infinite dimensional transform**
- Comparison of Kernels  
**linear for efficiency or Gaussian for power**
- Motivation and Primal Problem  
**add margin violations  $\xi_n$**
- Dual Problem  
**upper-bound  $\alpha_n$  by  $C$**
- Messages behind Soft-Margin SVM  
**bounded/free SVs for data analysis**
- Soft-Margin SVM as Regularized Model  
**L2-regularization with hinge error measure**