

# 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  $O(\tilde{d})$

**can we do so?**

Fast Inner Product for  $\Phi_2$ 

## 2nd order polynomial transform

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

—include both  $x_1 x_2$  &  $x_2 x_1$  for 'simplicity' :-)

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

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

# Kernel: Transform + Inner Product

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

$$\phi_2 \iff K_{\phi_2}(\mathbf{x}, \mathbf{x}') = 1 + (\mathbf{x}^T \mathbf{x}') + (\mathbf{x}^T \mathbf{x}')^2$$

- 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**  $(\mathbf{x}_s, y_s)$ ,

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

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

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

**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  $b \leftarrow \left( y_s - \sum_{\text{SV indices } n} \alpha_n y_n K(\mathbf{x}_n, \mathbf{x}_s) \right)$  with SV  $(\mathbf{x}_s, y_s)$
- ④ return SVs and their  $\alpha_n$  as well as  $b$  such that for new  $\mathbf{x}$ ,
 
$$g_{\text{SVM}}(\mathbf{x}) = \text{sign} \left( \sum_{\text{SV indices } n} \alpha_n y_n K(\mathbf{x}_n, \mathbf{x}) + b \right)$$

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

kernel SVM:

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

**Questions?**

# General Poly-2 Kernel

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

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

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

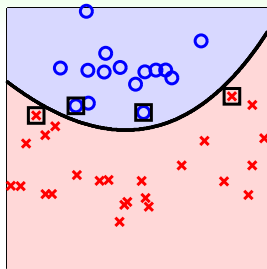
$$K_2(\mathbf{x}, \mathbf{x}') = (1 + \gamma \mathbf{x}^T \mathbf{x}')^2 \text{ with } \gamma > 0$$

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

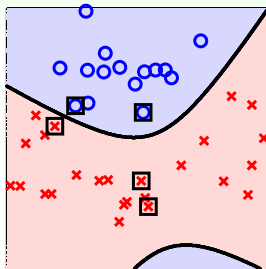
$K_2$  commonly used



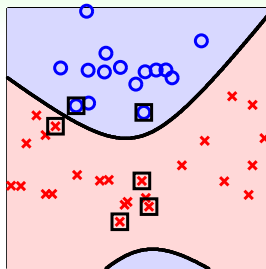
## Poly-2 Kernels in Action



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



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



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

- $g_{\text{SVM}}$  **different**, SVs **different**  
—‘hard’ to say which is better before learning
- change of **kernel**  $\Leftrightarrow$  change of **margin definition**

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

# General Polynomial Kernel

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

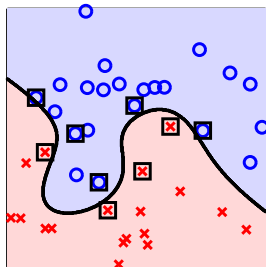
$$K_3(\mathbf{x}, \mathbf{x}') = (\zeta + \gamma \mathbf{x}^T \mathbf{x}')^3 \text{ with } \gamma > 0, \zeta \geq 0$$

$$\vdots$$

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

- embeds  $\Phi_Q$  specially with parameters  $(\gamma, \zeta)$
- allows computing large-margin **polynomial** classification **without dependence on  $\tilde{d}$**

SVM + **Polynomial** Kernel: **Polynomial** SVM

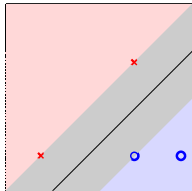


10-th order polynomial  
with margin 0.1

# Special Case: Linear Kernel

$$\begin{aligned}K_1(\mathbf{x}, \mathbf{x}') &= (0 + 1 \cdot \mathbf{x}^T \mathbf{x}')^1 \\&\vdots \\K_Q(\mathbf{x}, \mathbf{x}') &= (\zeta + \gamma \mathbf{x}^T \mathbf{x}')^Q \text{ with } \gamma > 0, \zeta \geq 0\end{aligned}$$

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



**linear first, remember? :-)**

**Questions?**

# Kernel of Infinite Dimensional Transform

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

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

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

more generally, **Gaussian kernel**

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

# Hypothesis of Gaussian SVM

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

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

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

Gaussian SVM:

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

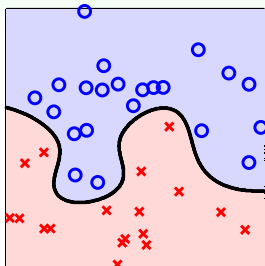
# Support Vector Mechanism

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

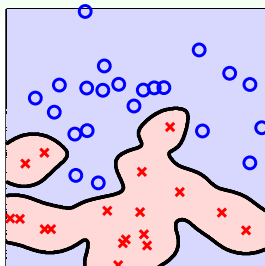
- transformed vector  $\mathbf{z} = \Phi(\mathbf{x}) \implies$  efficient kernel  $K(\mathbf{x}, \mathbf{x}')$
- store optimal  $\mathbf{w} \implies$  store a few SVs and  $\alpha_n$

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

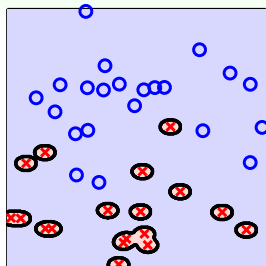
# Gaussian SVM in Action



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



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



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

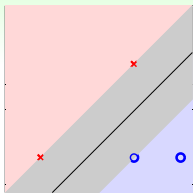
- large  $\gamma \implies$  sharp Gaussians  $\implies$  'overfit'?
- **warning: SVM can still overfit :-)**

Gaussian SVM: need careful selection of  $\gamma$



**Questions?**

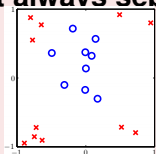
# Linear Kernel: Cons and Pros



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

## Cons

- restricted  
— **not always separable?!**

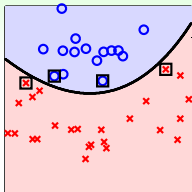


## Pros

- 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

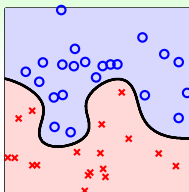
- **numerical difficulty** for large  $Q$ 
  - $|\zeta + \gamma \mathbf{x}^T \mathbf{x}'| < 1: K \rightarrow 0$
  - $|\zeta + \gamma \mathbf{x}^T \mathbf{x}'| > 1: K \rightarrow \text{big}$
- three parameters ( $\gamma, \zeta, Q$ )  
—**more difficult to select**

## Pros

- **less restricted** than linear
- strong physical control  
—‘knows’ **degree  $Q$**

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

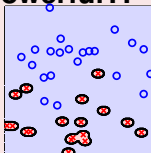
# Gaussian Kernel: Cons and Pros



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

## Cons

- **mysterious**—no  $\mathbf{w}$
- **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:  $\Phi(\mathbf{x})^T \Phi(\mathbf{x}')$
- any similarity  $\implies$  valid kernel? **not really**
- necessary **& sufficient** conditions for valid kernel:  
**Mercer's condition**
  - symmetric
  - let  $k_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$ , the matrix **K**

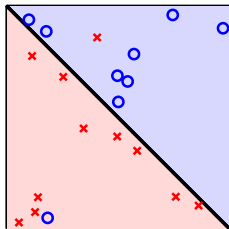
$$\begin{aligned}
 &= \begin{bmatrix} \Phi(\mathbf{x}_1)^T \Phi(\mathbf{x}_1) & \Phi(\mathbf{x}_1)^T \Phi(\mathbf{x}_2) & \dots & \Phi(\mathbf{x}_1)^T \Phi(\mathbf{x}_N) \\ \Phi(\mathbf{x}_2)^T \Phi(\mathbf{x}_1) & \Phi(\mathbf{x}_2)^T \Phi(\mathbf{x}_2) & \dots & \Phi(\mathbf{x}_2)^T \Phi(\mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ \Phi(\mathbf{x}_N)^T \Phi(\mathbf{x}_1) & \Phi(\mathbf{x}_N)^T \Phi(\mathbf{x}_2) & \dots & \Phi(\mathbf{x}_N)^T \Phi(\mathbf{x}_N) \end{bmatrix} \\
 &= \begin{bmatrix} \mathbf{z}_1 & \mathbf{z}_2 & \dots & \mathbf{z}_N \end{bmatrix}^T \begin{bmatrix} \mathbf{z}_1 & \mathbf{z}_2 & \dots & \mathbf{z}_N \end{bmatrix} \\
 &= \mathbf{Z}\mathbf{Z}^T \text{ must always be positive semi-definite}
 \end{aligned}$$

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

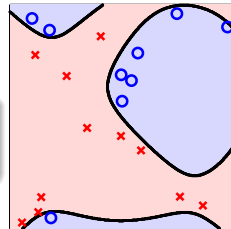
**Questions?**

# Cons of Hard-Margin SVM

recall: SVM can still overfit :-)

 $\Phi_1$ 

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

 $\Phi_4$ 

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

# Give Up on Some Examples

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

## min.-error perceptron

$$\min_{b, \mathbf{w}} \sum_{n=1}^N \mathbb{I}[y_n \neq \text{sign}(\mathbf{w}^T \mathbf{z}_n + b)]$$

## hard-margin SVM

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

$$\text{s.t.} \quad y_n(\mathbf{w}^T \mathbf{z}_n + b) \geq 1 \text{ for all } n$$

combination:

$$\min_{b, \mathbf{w}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \cdot \sum_{n=1}^N \mathbb{I}[y_n \neq \text{sign}(\mathbf{w}^T \mathbf{z}_n + b)]$$

$$\text{s.t.} \quad y_n(\mathbf{w}^T \mathbf{z}_n + b) \geq 1 \text{ for } \text{correct } n$$

$$y_n(\mathbf{w}^T \mathbf{z}_n + b) \geq -\infty \text{ for } \text{incorrect } n$$

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



## Soft-Margin SVM (1/2)

$$\begin{aligned}
 \min_{b, \mathbf{w}} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \cdot \sum_{n=1}^N \mathbb{I}[y_n \neq \text{sign}(\mathbf{w}^T \mathbf{z}_n + b)] \\
 \text{s.t.} \quad & y_n(\mathbf{w}^T \mathbf{z}_n + b) \geq 1 - \infty \cdot \mathbb{I}[y_n \neq \text{sign}(\mathbf{w}^T \mathbf{z}_n + b)]
 \end{aligned}$$

- $\mathbb{I}[\cdot]$ : non-linear, **not QP anymore** :-(  
—what about dual? kernel?
- 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**

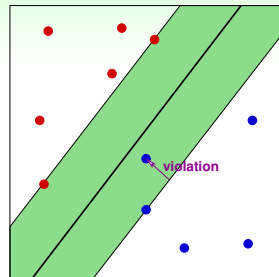
$$\begin{aligned}
 \text{soft-margin SVM: } \min_{b, \mathbf{w}, \xi} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \cdot \sum_{n=1}^N \xi_n \\
 \text{s.t.} \quad & y_n(\mathbf{w}^T \mathbf{z}_n + b) \geq 1 - \xi_n \text{ and } \xi_n \geq 0 \text{ for all } n
 \end{aligned}$$

## Soft-Margin SVM (2/2)

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

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

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



- 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,  $2N$  constraints

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

**Questions?**

# Lagrange Dual

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

Lagrange function with Lagrange multipliers  $\alpha_n$  and  $\beta_n$

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

want: Lagrange dual

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

Simplify  $\xi_n$  and  $\beta_n$ 

$$\max_{\alpha_n \geq 0, \beta_n \geq 0} \left( \min_{b, \mathbf{w}, \xi} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \cdot \sum_{n=1}^N \xi_n + \sum_{n=1}^N \alpha_n \cdot (1 - \xi_n - y_n(\mathbf{w}^T \mathbf{z}_n + b)) + \sum_{n=1}^N \beta_n \cdot (-\xi_n) \right)$$

- $\frac{\partial \mathcal{L}}{\partial \xi_n} = 0 = C - \alpha_n - \beta_n$
- 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  $b$

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

## Other Simplifications

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

familiar? :-)

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

$$\begin{aligned}
 \min_{\alpha} \quad & \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y_n y_m \mathbf{z}_n^T \mathbf{z}_m - \sum_{n=1}^N \alpha_n \\
 \text{subject to} \quad & \sum_{n=1}^N y_n \alpha_n = 0; \\
 & 0 \leq \alpha_n \leq C, \text{ for } n = 1, 2, \dots, N; \\
 \text{implicitly} \quad & \mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{z}_n; \\
 & \beta_n = C - \alpha_n, \text{ for } n = 1, 2, \dots, N
 \end{aligned}$$

—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  $\alpha \leftarrow \text{QP}(\mathbf{Q}_D, \mathbf{p}, \mathbf{A}, \mathbf{c})$
- 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

remaining question: step ③?

Solving for  $b$ 

## hard-margin SVM

complementary slackness:

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

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

## soft-margin SVM

complementary slackness:

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

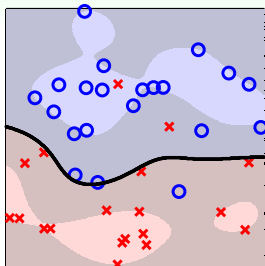
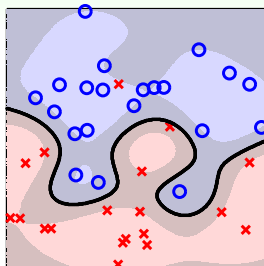
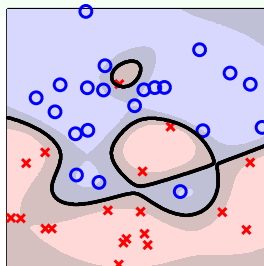
- SV ( $\alpha_s > 0$ )  
 $\Rightarrow b = y_s - y_s \xi_s - \mathbf{w}^T \mathbf{z}_s$
- free ( $\alpha_s < C$ )  
 $\Rightarrow \xi_s = 0$

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

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

—range of  $b$  otherwise

# Soft-Margin Gaussian SVM in Action

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

- large  $C \implies$  less noise tolerance  $\implies$  'overfit'?
- **warning: SVM can still overfit :-)**

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

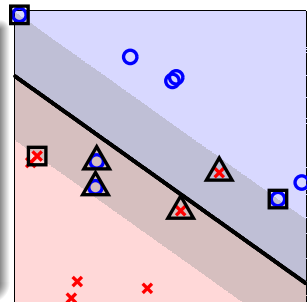
Physical Meaning of  $\alpha_n$ 

complementary slackness:

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

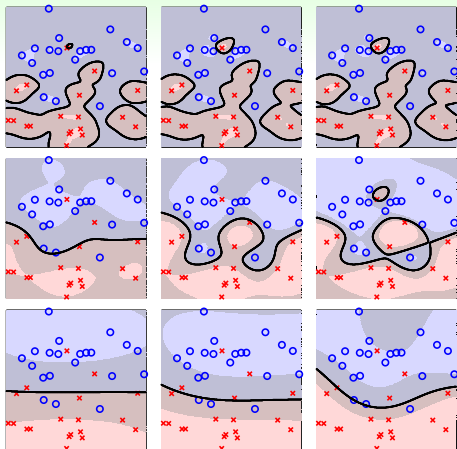
$$(C - \alpha_n)\xi_n = 0$$

- non SV ( $0 = \alpha_n$ ):  $\xi_n = 0$ ,  
'away from'/on **fat boundary**
- $\square$  free SV ( $0 < \alpha_n < C$ ):  $\xi_n = 0$ ,  
on **fat boundary**, locates  $b$
- $\triangle$  bounded SV ( $\alpha_n = C$ ):  
 $\xi_n$  = violation amount,  
'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

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

**Questions?**

# Wrap-Up

## Hard-Margin Primal

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

## Soft-Margin Primal

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

## Hard-Margin Dual

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

## Soft-Margin Dual

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

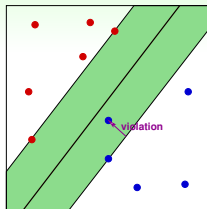
**soft**-margin preferred in practice;  
linear: LIBLINEAR; non-linear: LIBSVM

Slack Variables  $\xi_n$ 

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

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

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



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:

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



# Unconstrained Form

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

familiar? :-)

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

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

# SVM as Regularized Model

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

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

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

larger  $C$  or  $C \iff$  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**