

# Machine Learning

## (機器學習)

### Lecture 10: Support Vector Machine (1)

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

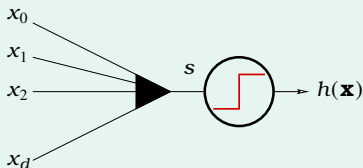
## Lecture 10: Support Vector Machine (1)

- Large-Margin Separating Hyperplane
- Standard Large-Margin Problem
- Support Vector Machine
- Motivation of Dual SVM
- Lagrange Dual SVM
- Solving Dual SVM
- Messages behind Dual SVM

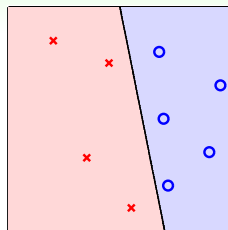
# Linear Classification Revisited

## PLA/pocket

$$h(\mathbf{x}) = \text{sign}(\mathbf{s})$$



plausible err = 0/1  
(small flipping noise)  
minimize **specially**

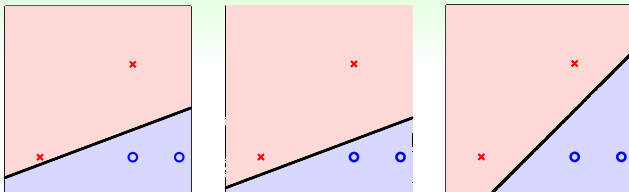


(linear separable)

linear (hyperplane) classifiers:

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})$$

## Which Line Is Best?

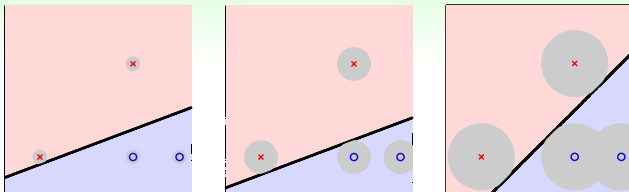


- PLA? depending on randomness
- VC bound? whichever you like!

$$E_{\text{out}}(\mathbf{w}) \leq \underbrace{E_{\text{in}}(\mathbf{w})}_0 + \underbrace{\Omega(\mathcal{H})}_{d_{\text{VC}}=d+1}$$

You? **rightmost one, possibly :-)**

# Why Rightmost Hyperplane?



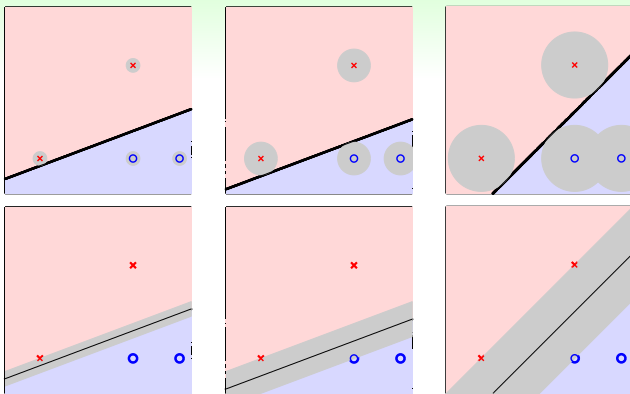
## informal argument

if (Gaussian-like) noise on future  $\mathbf{x} \approx \mathbf{x}_n$ :

$\mathbf{x}_n$ further from hyperplane	distance to closest $\mathbf{x}_n$
$\iff$ tolerate more noise	$\iff$ amount of noise tolerance
$\iff$ more robust to overfitting	$\iff$ robustness of hyperplane

rightmost one: **more robust**  
because of **larger distance to closest  $\mathbf{x}_n$**

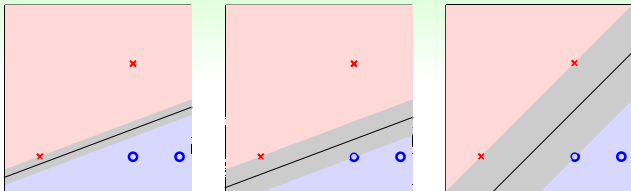
# Fat Hyperplane



- **robust** separating hyperplane: **fat**  
—far from both sides of examples
- **robustness**  $\equiv$  **fatness**: distance to closest  $\mathbf{x}_n$

goal: find **fattest** separating hyperplane

# Large-Margin Separating Hyperplane

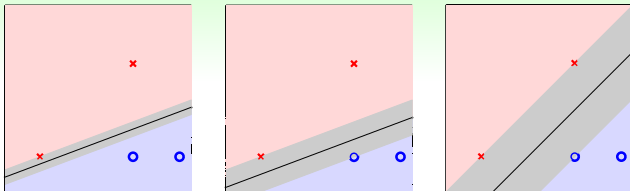


$$\begin{aligned}
 & \max_{\mathbf{w}} \quad \text{fatness}(\mathbf{w}) \\
 & \text{subject to} \quad \mathbf{w} \text{ classifies every } (\mathbf{x}_n, y_n) \text{ correctly} \\
 & \quad \text{fatness}(\mathbf{w}) = \min_{n=1, \dots, N} \text{distance}(\mathbf{x}_n, \mathbf{w})
 \end{aligned}$$

- fatness: formally called **margin**
- **correctness**:  $y_n = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$

goal: find **largest-margin**  
**separating** hyperplane

# Large-Margin Separating Hyperplane



$$\begin{aligned}
 & \max_{\mathbf{w}} \quad \text{margin}(\mathbf{w}) \\
 & \text{subject to} \quad \text{every } y_n \mathbf{w}^T \mathbf{x}_n > 0 \\
 & \quad \text{margin}(\mathbf{w}) = \min_{n=1, \dots, N} \text{distance}(\mathbf{x}_n, \mathbf{w})
 \end{aligned}$$

- fatness: formally called **margin**
- **correctness**:  $y_n = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$

goal: find **largest-margin**  
**separating** hyperplane



**Questions?**

# Distance to Hyperplane: Preliminary

$$\begin{aligned}
 & \max_{\mathbf{w}} \quad \text{margin}(\mathbf{w}) \\
 & \text{subject to} \quad \text{every } y_n \mathbf{w}^T \mathbf{x}_n > 0 \\
 & \quad \text{margin}(\mathbf{w}) = \min_{n=1, \dots, N} \text{distance}(\mathbf{x}_n, \mathbf{w})
 \end{aligned}$$

‘shorten’  $\mathbf{x}$  and  $\mathbf{w}$

distance needs  $w_0$  and  $(w_1, \dots, w_d)$  differently (to be derived)

$$\begin{array}{c} \textcolor{red}{b} \\ \textcolor{blue}{\mathbf{w}} \end{array} = \begin{array}{c} \textcolor{red}{w_0} \\ \textcolor{blue}{\begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}} \end{array} ; \quad \begin{array}{c} \textcolor{red}{x_0} \\ \textcolor{blue}{\mathbf{x}} \end{array} = \begin{array}{c} \textcolor{red}{1} \\ \textcolor{blue}{\begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}} \end{array}$$

for this part:  $h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + \textcolor{red}{b})$

# Distance to Hyperplane

want: distance( $\mathbf{x}$ ,  $b$ ,  $\mathbf{w}$ ), with hyperplane  $\mathbf{w}^T \mathbf{x}' + b = 0$

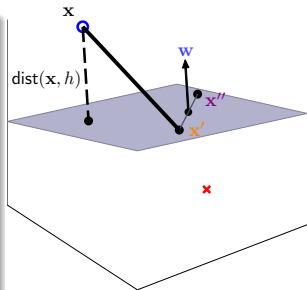
consider  $\mathbf{x}'$ ,  $\mathbf{x}''$  on hyperplane

①  $\mathbf{w}^T \mathbf{x}' = -b$ ,  $\mathbf{w}^T \mathbf{x}'' = -b$

②  $\mathbf{w} \perp$  hyperplane:

$$\begin{pmatrix} \mathbf{w}^T & \underbrace{(\mathbf{x}'' - \mathbf{x}')}_{\text{vector on hyperplane}} \end{pmatrix} = 0$$

③ distance = project  $(\mathbf{x} - \mathbf{x}')$  to  $\perp$  hyperplane



$$\text{distance}(\mathbf{x}, b, \mathbf{w}) = \left| \frac{\mathbf{w}^T}{\|\mathbf{w}\|} (\mathbf{x} - \mathbf{x}') \right| \stackrel{\text{①}}{=} \frac{1}{\|\mathbf{w}\|} |\mathbf{w}^T \mathbf{x} + b|$$

# Distance to **Separating** Hyperplane

$$\text{distance}(\mathbf{x}, b, \mathbf{w}) = \frac{1}{\|\mathbf{w}\|} |\mathbf{w}^T \mathbf{x} + b|$$

- separating** hyperplane: for every  $n$

$$y_n(\mathbf{w}^T \mathbf{x}_n + b) > 0$$

- distance to **separating** hyperplane:

$$\text{distance}(\mathbf{x}_n, b, \mathbf{w}) = \frac{1}{\|\mathbf{w}\|} y_n(\mathbf{w}^T \mathbf{x}_n + b)$$

$$\begin{array}{ll} \max_{b, \mathbf{w}} & \text{margin}(b, \mathbf{w}) \\ \text{subject to} & \text{every } y_n(\mathbf{w}^T \mathbf{x}_n + b) > 0 \\ & \text{margin}(b, \mathbf{w}) = \min_{n=1, \dots, N} \frac{1}{\|\mathbf{w}\|} y_n(\mathbf{w}^T \mathbf{x}_n + b) \end{array}$$

# Margin of **Special** Separating Hyperplane

$$\begin{aligned}
 & \max_{b, \mathbf{w}} \quad \text{margin}(\mathbf{b}, \mathbf{w}) \\
 & \text{subject to} \quad \text{every } y_n(\mathbf{w}^T \mathbf{x}_n + b) > 0 \\
 & \quad \text{margin}(\mathbf{b}, \mathbf{w}) = \min_{n=1, \dots, N} \frac{1}{\|\mathbf{w}\|} y_n(\mathbf{w}^T \mathbf{x}_n + b)
 \end{aligned}$$

- $\mathbf{w}^T \mathbf{x} + b = 0$  same as  $3\mathbf{w}^T \mathbf{x} + 3b = 0$ : scaling does not matter
- **special** scaling: only consider separating  $(b, \mathbf{w})$  such that

$$\min_{n=1, \dots, N} y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1 \implies \text{margin}(\mathbf{b}, \mathbf{w}) = \frac{1}{\|\mathbf{w}\|}$$

$$\begin{aligned}
 & \max_{b, \mathbf{w}} \quad \frac{1}{\|\mathbf{w}\|} \\
 & \text{subject to} \quad \text{every } y_n(\mathbf{w}^T \mathbf{x}_n + b) > 0 \\
 & \quad \min_{n=1, \dots, N} y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1
 \end{aligned}$$

# Standard Large-Margin Hyperplane Problem

$$\max_{b, \mathbf{w}} \frac{1}{\|\mathbf{w}\|} \quad \text{subject to} \quad \min_{n=1, \dots, N} y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$$

necessary constraints:  $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1$  for all  $n$

original constraint:  $\min_{n=1, \dots, N} y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$   
 want: optimal  $(b, \mathbf{w})$  here (inside)

if optimal  $(b, \mathbf{w})$  outside, e.g.  $y_n(\mathbf{w}^T \mathbf{x}_n + b) > 1.126$  for all  $n$   
 —can scale  $(b, \mathbf{w})$  to “more optimal”  $(\frac{b}{1.126}, \frac{\mathbf{w}}{1.126})$  (contradiction!)

final change:  $\max \implies \min$ , remove  $\sqrt{\phantom{x}}$ , add  $\frac{1}{2}$

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

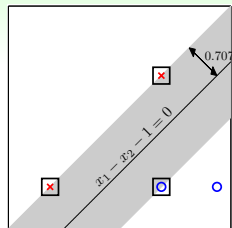
subject to  $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1$  for all  $n$

# Questions?

# Support Vector Machine (SVM)

optimal solution:  $(w_1 = 1, w_2 = -1, b = -1)$

$$\text{margin}(b, \mathbf{w}) = \frac{1}{\|\mathbf{w}\|} = \frac{1}{\sqrt{2}}$$



- examples on boundary: 'locates' fattest hyperplane  
other examples: **not needed**
- call boundary example **support vector** (candidate)

**support vector** machine (SVM):  
learn **fattest hyperplanes**  
(with help of **support vectors** )



# Solving General SVM

$$\begin{aligned} \min_{b, \mathbf{w}} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{subject to} \quad & y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \text{ for all } n \end{aligned}$$

- **not easy manually, of course :-)**
  - gradient descent? **not easy with constraints**
  - luckily:
    - (convex) quadratic objective function of  $(b, \mathbf{w})$
    - linear constraints of  $(b, \mathbf{w})$
- **quadratic programming**

**quadratic programming (QP):**  
'easy' optimization problem

# Quadratic Programming

optimal  $(\mathbf{b}, \mathbf{w}) = ?$

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

subject to  $y_n(\mathbf{w}^T \mathbf{x}_n + \mathbf{b}) \geq 1,$   
for  $n = 1, 2, \dots, N$

optimal  $\mathbf{u} \leftarrow \text{QP}(\mathbf{Q}, \mathbf{p}, \mathbf{A}, \mathbf{c})$

$$\min_{\mathbf{u}} \quad \frac{1}{2} \mathbf{u}^T \mathbf{Q} \mathbf{u} + \mathbf{p}^T \mathbf{u}$$

subject to  $\mathbf{a}_m^T \mathbf{u} \geq \mathbf{c}_m,$   
for  $m = 1, 2, \dots, M$

objective function:  $\mathbf{u} = \begin{bmatrix} \mathbf{b} \\ \mathbf{w} \end{bmatrix}; \mathbf{Q} = \begin{bmatrix} 0 & \mathbf{0}_d^T \\ \mathbf{0}_d & \mathbf{I}_d \end{bmatrix}; \mathbf{p} = \mathbf{0}_{d+1}$

constraints:  $\mathbf{a}_n^T = y_n \begin{bmatrix} 1 & \mathbf{x}_n^T \end{bmatrix}; \mathbf{c}_n = 1; M = N$

SVM with general QP solver:  
easy **if you've read the manual :-)**

## SVM with QP Solver

## Linear Hard-Margin SVM Algorithm

- 1  $Q = \begin{bmatrix} 0 & \mathbf{0}_d^T \\ \mathbf{0}_d & I_d \end{bmatrix}; \mathbf{p} = \mathbf{0}_{d+1}; \mathbf{a}_n^T = y_n [1 \quad \mathbf{x}_n^T]; c_n = 1$
- 2  $\begin{bmatrix} b \\ \mathbf{w} \end{bmatrix} \leftarrow \text{QP}(Q, \mathbf{p}, A, \mathbf{c})$
- 3 return  $b$  &  $\mathbf{w}$  as  $g_{\text{SVM}}$

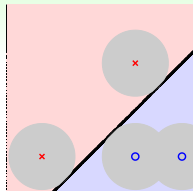
- **hard-margin**: nothing violate 'fat boundary'
- **linear**:  $\mathbf{x}_n$

want **non-linear**?

$\mathbf{z}_n = \Phi(\mathbf{x}_n)$ —**remember? :-)**

# Why Large-Margin Hyperplane?

$$\begin{aligned} \min_{b, \mathbf{w}} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{subject to} \quad & y_n(\mathbf{w}^T \mathbf{z}_n + b) \geq 1 \text{ for all } n \end{aligned}$$



	minimize	constraint
regularization	$E_{\text{in}}$	$\mathbf{w}^T \mathbf{w} \leq C$
SVM	$\mathbf{w}^T \mathbf{w}$	$E_{\text{in}} = 0$ [and more]

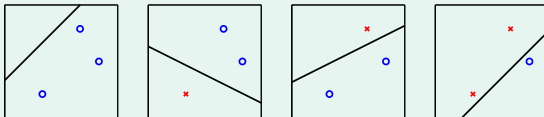
SVM (large-margin hyperplane):  
**'weight-decay regularization' within  $E_{\text{in}} = 0$**

# Large-Margin Restricts Dichotomies

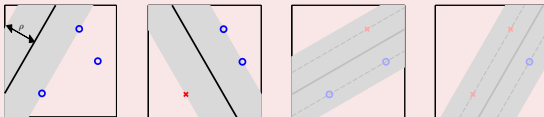
consider 'large-margin algorithm'  $\mathcal{A}_\rho$ :

either **returns  $g$  with  $\text{margin}(g) \geq \rho$  (if exists)**, or 0 otherwise

$\mathcal{A}_0$ : like PLA  $\implies$  shatter 'general' 3 inputs



$\mathcal{A}_{1.126}$ : more strict than SVM  $\implies$  cannot shatter any 3 inputs



fewer dichotomies  $\implies$  smaller 'VC dim.'  $\implies$  **better generalization**

# VC Dimension of Large-Margin Algorithm

fewer dichotomies  $\implies$  smaller ‘VC dim.’

**considers  $d_{\text{VC}}(\mathcal{A}_\rho)$  [data-dependent, need more than VC]**

instead of  $d_{\text{VC}}(\mathcal{H})$  [data-independent, covered by VC]

generally, when  $\mathcal{X}$  in radius- $R$  hyperball:

$$d_{\text{VC}}(\mathcal{A}_\rho) \leq \min \left( \frac{R^2}{\rho^2}, d \right) + 1 \leq \underbrace{d+1}_{d_{\text{VC}}(\text{perceptrons})}$$

# Benefits of Large-Margin Hyperplanes

	large-margin hyperplanes	hyperplanes	hyperplanes + feature transform $\phi$
#	even fewer	<b>not many</b>	many
boundary	simple	simple	<b>sophisticated</b>

- **not many** good, for  $d_{VC}$  and generalization
- **sophisticated** good, for possibly better  $E_{in}$

a new possibility: non-linear SVM

	large-margin hyperplanes + numerous feature transform $\phi$
#	<b>not many</b>
boundary	<b>sophisticated</b>

**Questions?**



# Non-Linear Support Vector Machine Revisited

## Non-Linear Hard-Margin SVM

$$\begin{aligned}
 \min_{b, \mathbf{w}} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\
 \text{s. t.} \quad & y_n (\mathbf{w}^T \underbrace{\mathbf{z}_n}_{\Phi(\mathbf{x}_n)} + b) \geq 1, \\
 & \text{for } n = 1, 2, \dots, N
 \end{aligned}$$

- 1  $Q = \begin{bmatrix} 0 & \mathbf{0}_{\tilde{d}}^T \\ \mathbf{0}_{\tilde{d}} & I_{\tilde{d}} \end{bmatrix}; \mathbf{p} = \mathbf{0}_{\tilde{d}+1};$   
 $\mathbf{a}_n^T = y_n \begin{bmatrix} 1 & \mathbf{z}_n^T \end{bmatrix}; \mathbf{c}_n = 1$
- 2  $\begin{bmatrix} b \\ \mathbf{w} \end{bmatrix} \leftarrow \text{QP}(Q, \mathbf{p}, \mathbf{A}, \mathbf{c})$
- 3 return  $b \in \mathbb{R}$  &  $\mathbf{w} \in \mathbb{R}^{\tilde{d}}$  with  
 $g_{\text{SVM}}(\mathbf{x}) = \text{sign}(\mathbf{w}^T \Phi(\mathbf{x}) + b)$

- demanded: **not many** (large-margin), but **sophisticated** boundary (feature transform)
- QP with  $\tilde{d} + 1$  variables and  $N$  constraints  
 —challenging if  $\tilde{d}$  large, **or infinite?! :-)**

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

Todo: SVM ‘without’  $\tilde{d}$ 

## Original SVM

(convex) QP of

- $\tilde{d} + 1$  variables
- $N$  constraints

## ‘Equivalent’ SVM

(convex) QP of

- $N$  variables
- $N + 1$  constraints

## Warning: Heavy Math!!!!!!

- introduce some necessary math without rigor to help **understand SVM deeper**
- ‘**claim**’ **some results** if details unnecessary  
—like how we ‘claimed’ Hoeffding

‘Equivalent’ SVM: based on some  
**dual problem** of Original SVM

## Key Tool: Lagrange Multipliers

Regularization by  
Constrained-Minimizing  $E_{\text{in}}$

$$\min_{\mathbf{w}} E_{\text{in}}(\mathbf{w}) \text{ s.t. } \mathbf{w}^T \mathbf{w} \leq C$$



Regularization by  
Minimizing  $E_{\text{aug}}$

$$\min_{\mathbf{w}} E_{\text{aug}}(\mathbf{w}) = E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$$

- $C$  equivalent to some  $\lambda \geq 0$  by checking **optimality condition**

$$\nabla E_{\text{in}}(\mathbf{w}) + \frac{2\lambda}{N} \mathbf{w} = \mathbf{0}$$

- regularization: view  $\lambda$  as **given parameter instead of  $C$** , and solve 'easily'
- dual SVM: view  $\lambda$ 's as unknown given the constraints, and **solve them as variables instead**

how many  $\lambda$ 's as variables?  
 $N$ —one per constraint

## Starting Point: Constrained to 'Unconstrained'

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

## Lagrange Function

with Lagrange multipliers  ~~$\alpha_n$~~ ,

$$\mathcal{L}(b, \mathbf{w}, \boldsymbol{\alpha}) = \underbrace{\frac{1}{2} \mathbf{w}^T \mathbf{w}}_{\text{objective}} + \sum_{n=1}^N \alpha_n \underbrace{(1 - y_n(\mathbf{w}^T \mathbf{z}_n + b))}_{\text{constraint}}$$

## Claim

$$\text{SVM} \equiv \min_{b, \mathbf{w}} \left( \max_{\text{all } \alpha_n \geq 0} \mathcal{L}(b, \mathbf{w}, \boldsymbol{\alpha}) \right) = \min_{b, \mathbf{w}} \left( \infty \text{ if violate ; } \frac{1}{2} \mathbf{w}^T \mathbf{w} \text{ if feasible} \right)$$

- any 'violating'  $(b, \mathbf{w})$ :  $\max_{\text{all } \alpha_n \geq 0} \left( \square + \sum_n \alpha_n (\text{some positive}) \right) \rightarrow \infty$
- any 'feasible'  $(b, \mathbf{w})$ :  $\max_{\text{all } \alpha_n \geq 0} \left( \square + \sum_n \alpha_n (\text{all non-positive}) \right) = \square$

constraints now **hidden in** max

**Questions?**

# Strong Duality of Quadratic Programming

$$\underbrace{\min_{b, w} \left( \max_{\text{all } \alpha_n \geq 0} \mathcal{L}(b, w, \alpha) \right)}_{\text{equiv. to original (primal) SVM}} = \underbrace{\max_{\text{all } \alpha_n \geq 0} \left( \min_{b, w} \mathcal{L}(b, w, \alpha) \right)}_{\text{Lagrange dual}}$$

- ‘=’: **strong duality**, true for QP if
  - convex primal
  - feasible primal (true if  $\Phi$ -separable)
  - linear constraints
- called **constraint qualification**

exists **primal-dual** optimal  
solution  $(b, w, \alpha)$  for **both sides**

## Solving Lagrange Dual: Simplifications (1/2)

$$\max_{\text{all } \alpha_n \geq 0} \left( \min_{b, \mathbf{w}} \underbrace{\frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{n=1}^N \alpha_n (1 - y_n (\mathbf{w}^T \mathbf{z}_n + b))}_{\mathcal{L}(b, \mathbf{w}, \alpha)} \right)$$

- inner problem ‘unconstrained’, at optimal:

$$\frac{\partial \mathcal{L}(b, \mathbf{w}, \alpha)}{\partial b} = 0 = - \sum_{n=1}^N \alpha_n y_n$$

- no loss of optimality if solving with constraint  $\sum_{n=1}^N \alpha_n y_n = 0$

but wait,  $b$  can be removed

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

## Solving Lagrange Dual: Simplifications (2/2)

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

- inner problem 'unconstrained', at optimal:

$$\frac{\partial \mathcal{L}(b, \mathbf{w}, \alpha)}{\partial w_i} = 0 = w_i - \sum_{n=1}^N \alpha_n y_n z_{n,i}$$

- no loss of optimality if solving with constraint  $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{z}_n$

but wait!

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

$$\iff \max_{\text{all } \alpha_n \geq 0, \sum y_n \alpha_n = 0, \mathbf{w} = \sum \alpha_n y_n \mathbf{z}_n} -\frac{1}{2} \left\| \sum_{n=1}^N \alpha_n y_n \mathbf{z}_n \right\|^2 + \sum_{n=1}^N \alpha_n$$



## KKT Optimality Conditions

$$\max_{\text{all } \alpha_n \geq 0, \sum y_n \alpha_n = 0, \mathbf{w} = \sum \alpha_n y_n \mathbf{z}_n} -\frac{1}{2} \left\| \sum_{n=1}^N \alpha_n y_n \mathbf{z}_n \right\|^2 + \sum_{n=1}^N \alpha_n$$

if **primal-dual** optimal  $(\mathbf{b}, \mathbf{w}, \boldsymbol{\alpha})$ ,

- **primal feasible**:  $y_n(\mathbf{w}^T \mathbf{z}_n + \mathbf{b}) \geq 1$
- **dual feasible**:  $\alpha_n \geq 0$
- **dual-inner** optimal:  $\sum y_n \alpha_n = 0$ ;  $\mathbf{w} = \sum \alpha_n y_n \mathbf{z}_n$
- **primal-inner** optimal (at optimal all ‘**Lagrange terms**’ disappear):

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

—called **Karush-Kuhn-Tucker (KKT) conditions**, necessary for optimality [& sufficient here]

will use **KKT** to ‘solve’  $(\mathbf{b}, \mathbf{w})$  from optimal  $\boldsymbol{\alpha}$

**Questions?**

# Dual Formulation of Support Vector Machine

$$\max_{\text{all } \alpha_n \geq 0, \sum y_n \alpha_n = 0, \mathbf{w} = \sum \alpha_n y_n \mathbf{z}_n} -\frac{1}{2} \left\| \sum_{n=1}^N \alpha_n y_n \mathbf{z}_n \right\|^2 + \sum_{n=1}^N \alpha_n$$

standard hard-margin SVM **dual**

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

(convex) QP of  $N$  variables &  $N + 1$  constraints, as promised

how to solve? **yeah, we know QP! :-)**

## Dual SVM with QP Solver

optimal  $\alpha = ?$ 

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

$$- \sum_{n=1}^N \alpha_n$$

subject to

$$\sum_{n=1}^N y_n \alpha_n = 0;$$

$$\alpha_n \geq 0,$$

$$\text{for } n = 1, 2, \dots, N$$

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

$$\min_{\alpha} \quad \frac{1}{2} \alpha^T \mathbf{Q} \alpha + \mathbf{p}^T \alpha$$

subject to  $\mathbf{a}_i^T \alpha \geq c_i,$   
for  $i = 1, 2, \dots$

- $q_{n,m} = y_n y_m \mathbf{z}_n^T \mathbf{z}_m$
- $\mathbf{p} = -\mathbf{1}_N$
- $\mathbf{a}_{\geq} = \mathbf{y}, \mathbf{a}_{\leq} = -\mathbf{y};$   
 $\mathbf{a}_n^T = n\text{-th unit direction}$
- $c_{\geq} = 0, c_{\leq} = 0; c_n = 0$

note: many solvers treat **equality** ( $\mathbf{a}_{\geq}, \mathbf{a}_{\leq}$ ) &  
**bound** ( $\mathbf{a}_n$ ) constraints **specially for numerical stability**

## Dual SVM with Special QP Solver

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

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^T \mathbf{Q}_D \alpha + \mathbf{p}^T \alpha \\ \text{subject to} \quad & \text{special equality and bound constraints} \end{aligned}$$

- $q_{n,m} = y_n y_m \mathbf{z}_n^T \mathbf{z}_m$ , often non-zero
- if  $N = 30,000$ , dense  $\mathbf{Q}_D$  ( $N$  by  $N$  symmetric) takes  $> 3\text{G}$  RAM
- need **special solver** for
  - not storing whole  $\mathbf{Q}_D$
  - utilizing **special constraints** properly

to scale up to large  $N$

usually better to use **special solver** in practice

Optimal ( $\mathbf{b}, \mathbf{w}$ )

## KKT conditions

if primal-dual optimal ( $\mathbf{b}, \mathbf{w}, \alpha$ ),

- primal feasible:  $y_n(\mathbf{w}^T \mathbf{z}_n + \mathbf{b}) \geq 1$
- dual feasible:  $\alpha_n \geq 0$
- dual-inner optimal:  $\sum y_n \alpha_n = 0$ ;  $\mathbf{w} = \sum \alpha_n y_n \mathbf{z}_n$
- primal-inner optimal (at optimal all 'Lagrange terms' disappear):

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

- optimal  $\alpha \implies$  optimal  $\mathbf{w}$ ? easy above!
- optimal  $\alpha \implies$  optimal  $\mathbf{b}$ ? a range from primal feasible & equality from comp. slackness if one  $\alpha_n > 0 \implies \mathbf{b} = y_n - \mathbf{w}^T \mathbf{z}_n$

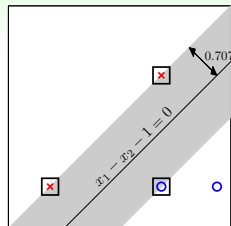
**comp. slackness:**

$$\alpha_n > 0 \implies \text{on fat boundary (SV!)}$$

**Questions?**

# Support Vectors Revisited

- on boundary: 'locates' fattest hyperplane;  
others: **not needed**
- examples with  $\alpha_n > 0$ : on boundary
- call  $\alpha_n > 0$  examples ( $\mathbf{z}_n, y_n$ )  
**support vectors** ~~(candidates)~~
- SV** (positive  $\alpha_n$ )  
 $\subseteq$  SV candidates (on boundary)



- only **SV** needed to compute  $\mathbf{w}$ :  $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{z}_n = \sum_{\text{SV}} \alpha_n y_n \mathbf{z}_n$
- only **SV** needed to compute  $b$ :  $b = y_n - \mathbf{w}^T \mathbf{z}_n$  with any **SV** ( $\mathbf{z}_n, y_n$ )

**SVM**: learn **fattest hyperplane**  
by identifying **support vectors**  
with **dual** optimal solution



# Summary: Two Forms of Hard-Margin SVM

## Primal Hard-Margin SVM

$$\begin{aligned} \min_{b, \mathbf{w}} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{sub. to} \quad & y_n(\mathbf{w}^T \mathbf{z}_n + b) \geq 1, \\ & \text{for } n = 1, 2, \dots, N \end{aligned}$$

- $\tilde{d} + 1$  variables,  
 $N$  constraints  
 —suitable when  $\tilde{d} + 1$  small
- physical meaning: locate  
**specially-scaled**  $(b, \mathbf{w})$

## Dual Hard-Margin SVM

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

- $N$  variables,  
 $N + 1$  simple constraints  
 —suitable when  $N$  small
- physical meaning: locate  
**SVs**  $(\mathbf{z}_n, y_n)$  & their  $\alpha_n$

both eventually result in optimal  $(b, \mathbf{w})$  for fattest hyperplane

$$g_{\text{SVM}}(\mathbf{x}) = \text{sign}(\mathbf{w}^T \Phi(\mathbf{x}) + b)$$

# Are We Done Yet?

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

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

- $N$  variables,  $N + 1$  constraints: **no dependence on  $\tilde{d}$ ?**
- $q_{n,m} = y_n y_m \mathbf{z}_n^T \mathbf{z}_m$ : inner product in  $\mathbb{R}^{\tilde{d}}$   
 —  $O(\tilde{d})$  via naïve computation!

no dependence **only if**  
**avoiding naïve computation (next lecture :-))**

**Questions?**

# Summary

## 1 Embedding Numerous Features: Kernel Models

### Lecture 10: Support Vector Machine (1)

- Large-Margin Separating Hyperplane  
**intuitively more robust against noise**
- Standard Large-Margin Problem  
**minimize 'length of  $w$ ' at special separating scale**
- Support Vector Machine  
**'easy' via quadratic programming**
- Motivation of Dual SVM  
**want to remove dependence on  $\tilde{d}$**
- Lagrange Dual SVM  
**KKT conditions link primal/dual**
- Solving Dual SVM  
**another QP, better solved with special solver**
- Messages behind Dual SVM  
**SVs represent fattest hyperplane**