

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

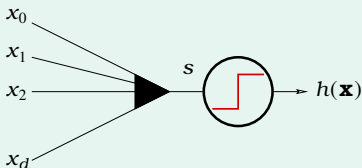
Lecture 10: Support Vector Machine (1)

- 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

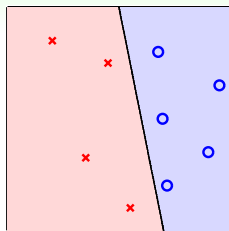
Linear Classification Revisited

PLA/pocket

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



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



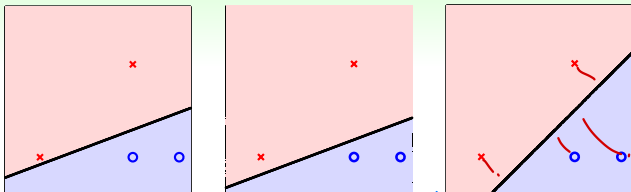
(linear separable)

data是线性可分

linear (hyperplane) classifiers:

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

Which Line Is Best?



↑ 答案不是唯一的 那條線好。

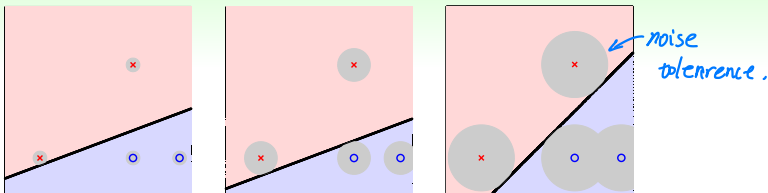
- PLA? depending on randomness → 不一定
- VC bound? whichever you like!

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

parameter 都一樣。

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

Why Rightmost Hyperplane?



informal argument

if (Gaussian-like) noise on future $\mathbf{x} \approx \mathbf{x}_n$ *add noise* *原始数据*

\mathbf{x}_n further from hyperplane

distance to closest \mathbf{x}_n

\iff tolerate more noise

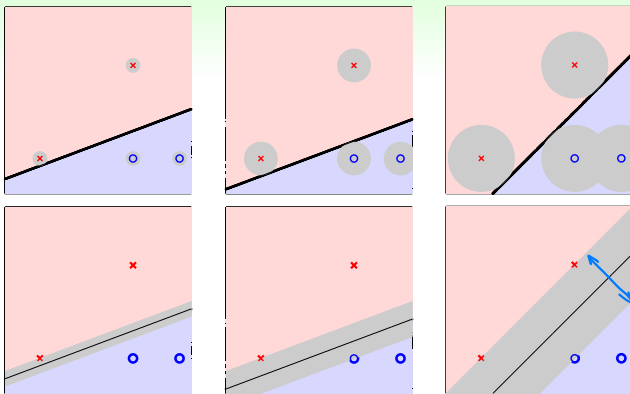
\iff amount of noise tolerance

\iff more robust to overfitting

\iff robustness of hyperplane

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

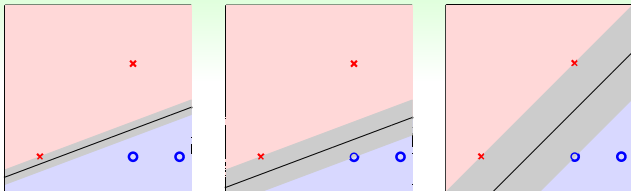
Fat Hyperplane



- **robust** separating hyperplane: **fat**
—far from both sides of examples
- **robustness** \equiv **fatness**: distance to closest x_n

goal: find **fattest** separating hyperplane

Large-Margin Separating Hyperplane

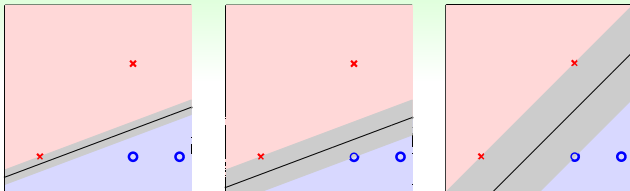


$\max_{\mathbf{w}}$ **fatness**(\mathbf{w}) *找最胖的线* *是线性可分的*
 subject to \mathbf{w} classifies every (\mathbf{x}_n, y_n) correctly
fatness(\mathbf{w}) = $\min_{n=1, \dots, N} \text{distance}(\mathbf{x}_n, \mathbf{w})$ *找最小距离*

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

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

Large-Margin Separating Hyperplane



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

- 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

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

‘shorten’ \mathbf{x} and \mathbf{w}

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

bias b = w_0 (把 w_0 放进 \mathbf{w}) ~~$x_0 = 1$~~

$$\begin{bmatrix} | \\ \mathbf{w} \\ | \end{bmatrix} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} ; \quad \begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix} = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

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

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

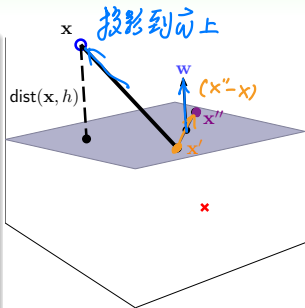
② $\mathbf{w} \perp$ hyperplane:

平面的法向量.

$$\begin{pmatrix} \mathbf{w}^T & (\mathbf{x}'' - \mathbf{x}') \end{pmatrix} = 0$$

vector on hyperplane

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



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

Handwritten notes: $-\frac{\mathbf{w}^T}{\|\mathbf{w}\|} \mathbf{x}' = \frac{+b}{\|\mathbf{w}\|}$ (with an arrow pointing to the b term in the equation)

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

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

correctness

- distance to **separating** hyperplane:

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

乘上 y_n 後, abs 可以拿掉

$$\frac{1}{\|\mathbf{w}\|} \left(\sum \alpha_n y_n x_n^2 + b \right) \frac{1}{\|\mathbf{z}\|}$$

max
 b, \mathbf{w}

margin(b, \mathbf{w})

subject to

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

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

Margin of **Special** Separating Hyperplane

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

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

強制scale

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

$$\begin{aligned}
 & \max_{b, w} \quad \frac{1}{\|w\|} \\
 & \text{subject to} \quad \text{every } y_n(w^T x_n + b) > 0
 \end{aligned}$$

$$\min_{n=1, \dots, N} y_n(w^T x_n + b) = 1$$

← 消掉了, =1 較嚴格
新 condition.

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)

反證法：證不可能所有 x_n 都 > 1 ，會有 $= 1$

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!)
 scale.

放鬆條件
也沒有損失

$\|\mathbf{w}\| \downarrow, \frac{1}{\|\mathbf{w}\|} \uparrow$, 這是解更好

final change: $\max \Rightarrow \min$, remove $\sqrt{\quad}$, add $\frac{1}{2}$

$\min_{b, \mathbf{w}}$

$$\frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$\frac{1}{\|\mathbf{w}\|} \Rightarrow \|\mathbf{w}\| \Rightarrow \mathbf{w}^T \mathbf{w}$$

$$\downarrow$$

$$\frac{1}{2} \mathbf{w}^T \mathbf{w}$$

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

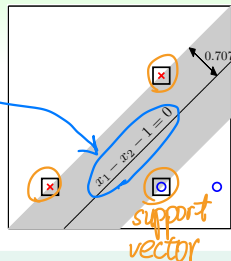
(b, \mathbf{w}) 不是
optimal. $\rightarrow \leftarrow$

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** ← 把這些丟掉也不影响
- 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 $(b, w) = ?$

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

subject to $y_n(w^T x_n + b) \geq 1$,
for $n = 1, 2, \dots, N$

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

u 的一次函数 \leftarrow 二次项系数

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

subject to $a_m^T u \geq c_m$,
for $m = 1, 2, \dots, M$

\leftarrow 二次项系数

一次式限制

objective function:

欲求的 variable.

$$u = \begin{bmatrix} b \\ w \end{bmatrix}$$

$$Q = \begin{bmatrix} 0 & 0_d^T \\ 0_d & I_d \end{bmatrix}$$

$$p = \begin{bmatrix} b \\ 0_{d+1} \end{bmatrix}$$

没有一次项

constraints:

$$a_n^T = y_n \begin{bmatrix} 1 & x_n^T \end{bmatrix}; \quad c_n = 1; \quad M = N$$

SVM with general QP solver:
easy if you've read the manual \leftarrow wtf

SVM with QP Solver

Linear Hard-Margin SVM Algorithm

$$① \quad Q = \begin{bmatrix} 0 & \mathbf{0}_d^T \\ \mathbf{0}_d & \mathbf{I}_d \end{bmatrix}; \quad \mathbf{p} = \mathbf{0}_{d+1}; \quad \mathbf{a}_n^T = y_n [1 \quad \mathbf{x}_n^T]; \quad c_n = 1$$

$$② \quad \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix} \leftarrow \text{QP}(Q, \mathbf{p}, A, \mathbf{c})$$

③ return b & \mathbf{w} as g_{SVM}

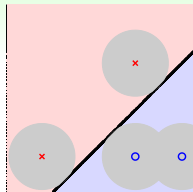
- * hard-margin: nothing violate 'fat boundary' ← 保持所有点都要分开
- * linear: \mathbf{x}_n ← 找分隔线

want **non-linear**?

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

Why Large-Margin Hyperplane?

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



- 係2面

對調

	minimize	constraint
regularization	E_{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 $E_{\text{in}} = 0$

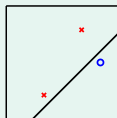
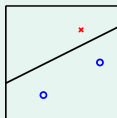
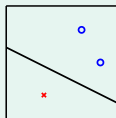
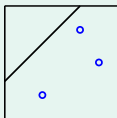
Large-Margin Restricts Dichotomies

consider 'large-margin algorithm' \mathcal{A}_ρ :

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

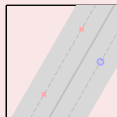
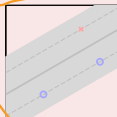
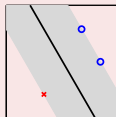
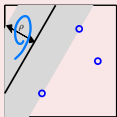
\mathcal{A}_0 : like PLA \implies shatter 'general' 3 inputs

3 input 能
shatter



$\mathcal{A}_{1.26}$: more strict than SVM \implies cannot shatter any 3 inputs

但多了限制
 ρ 宽度



找不到 $\rho=1.26$
的 "linear"

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

VC Dimension of Large-Margin Algorithm

fewer dichotomies \implies smaller 'VC dim.'

considers $d_{VC}(\mathcal{A}_\rho)$ [data-dependent, need more than VC]
 instead of $d_{VC}(\mathcal{H})$ [data-independent, covered by VC]

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

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

比原本小

↑
 是限制的寬度

Benefits of Large-Margin Hyperplanes

	large-margin hyperplanes	hyperplanes	hyperplanes + feature transform ϕ
#	even fewer	not many	many
boundary	simple	simple	sophisticated

- **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 ϕ ✓
#	not many
boundary	sophisticated } 22 min

Questions?

Non-Linear Support Vector Machine Revisited

Non-Linear Hard-Margin SVM

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

控制模型複雜度

- 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) W dimension 很大?
- 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

只跟 N 有关
 $N = \# \text{ of data}$

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

对偶问题

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

Key Tool: Lagrange MultipliersRegularization by
Constrained-Minimizing E_{in}

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

W 长度不能太大

Regularization by
Minimizing E_{aug}

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

乘上係数的 constant

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

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

- regularization: view λ as given parameter instead of C , and solve 'easily'
- dual SVM: view λ 's as unknown given the constraints, and **solve them as variables instead** *把 λ 當變數*

how many λ 's as variables? N —one per constraint*SVM 有 N constraints.*

Starting Point: Constrained to 'Unconstrained'

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

Lagrange Function

with Lagrange multipliers ~~α_n~~

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

係數.

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}) : $\max_{\text{all } \alpha_n \geq 0} \left(\square + \sum_n \alpha_n (\text{some positive}) \right) \rightarrow \infty$
 壞 (b, \mathbf{w}) 得到 ∞
- 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

使得式子看起來沒 constraint
 optimal value.

Questions?

Strong Duality of Quadratic Programming

min, max 做交换

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

- ‘=’: **strong duality**, true for **QP** if

- **convex** primal
- **feasible** primal (true if Φ -separable) *弱*
- **linear** constraints, *SVM 本来就是线性*

— called constraint qualification \Rightarrow *能直接解 RHS*

exists **primal-dual** optimal
solution (b, w, α) for both sides.

Solving Lagrange Dual: Simplifications (1/2)

RHS

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

沒有條件了

- inner problem 'unconstrained', at optimal:

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

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

but wait, **b can be removed**

只考慮微分=0的 α . ↓ 變成單變數.

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

Solving Lagrange Dual: Simplifications (2/2)

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

- inner problem 'unconstrained', at optimal:

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

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

極值發生的位置

but wait!

只找極值處

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

3 conditions

KKT Optimality Conditions

人名

幾乎只有 α $\frac{1}{\|w\|} = \text{distance.}$

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

大約 ≤ 0

if **primal-dual** optimal (b, \mathbf{w}, α) , $-\frac{1}{2} \|\mathbf{w}\|^2 + (-\|\Sigma\|^2 + 2\sum \alpha)^{-\frac{1}{2}} = \frac{1}{\|w\|}$

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

$$\frac{\alpha_n}{y_n^T \mathbf{z}_n + b} = 0$$

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

$$y_n \mathbf{w}^T \mathbf{z}_n + y_n b = 1$$

← 2个至少有1个等于零

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

will use **KKT** to 'solve' (b, \mathbf{w}) from optimal α

Questions?

Dual Formulation of Support Vector Machine

希望求 min, 取负号

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

standard hard-margin SVM dual

平方项展开

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

subject to

$$\sum_{n=1}^N y_n \alpha_n = 0; \quad \leftarrow 1 \text{ constraint.}$$

$\alpha_n \geq 0, \text{ for } n = 1, 2, \dots, N \quad \leftarrow N \text{ conditions}$

$N+1$

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

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

Dual SVM with QP Solver

解2次

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

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

$$\alpha_n \geq 0,$$

for $n = 1, 2, \dots, N$

表成QP

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 $i = 1, 2, \dots$

- $q_{n,m} = y_n y_m \mathbf{z}_n^T \mathbf{z}_m$ 2次項係數.
- $\mathbf{p} = -\mathbf{1}_N$ 1次項係數.
- $\mathbf{a}_{\geq} = \mathbf{y}, \mathbf{a}_{\leq} = -\mathbf{y};$
 $\mathbf{a}_n^T = n\text{-th unit direction}.$
- $c_{\geq} = 0, c_{\leq} = 0; c_n = 0$

特殊的bound.

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

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

- $q_{n,m} = y_n y_m \mathbf{z}_n^T \mathbf{z}_m$, often non-zero. 不稀疏
- if $N = 30,000$, dense \mathbf{Q}_D (N by N symmetric) takes $> 3\text{G RAM}$
- need special solver for \mathbf{Q}_D \uparrow 3万等, 存一半
 • not storing whole \mathbf{Q}_D \leftarrow 特殊加速
 • utilizing special constraints properly

to scale up to large N

usually better to use special solver in practice

Optimal (b, \mathbf{w}) 有以後解得(b, \mathbf{w})

KKT conditions

if primal-dual optimal (b, \mathbf{w}, α), 中間產物

- primal feasible: $y_n(\mathbf{w}^T \mathbf{z}_n + b) \geq 1$
- 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):

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

- optimal $\alpha \Rightarrow$ optimal \mathbf{w} ? easy above! 其中有一人=零
- optimal $\alpha \Rightarrow$ optimal b ? a range from primal feasible & equality from comp. slackness if one $\alpha_n > 0 \Rightarrow b = y_n - \mathbf{w}^T \mathbf{z}_n$

該點剛好是 support vector

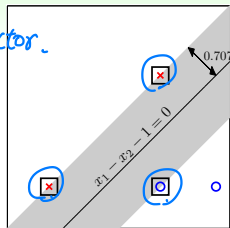
comp. slackness:

 $\alpha_n > 0 \Rightarrow$ on fat boundary (SV!) 求得 b .

Questions?

Support Vectors Revisited

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



- only **SV** needed to compute \mathbf{w} : $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{z}_n = \sum_{\text{SV}} \alpha_n y_n \mathbf{z}_n$ α 的 linear combine. 只要 SV 的 α.
- only **SV** needed to compute b : $b = y_n - \mathbf{w}^T \mathbf{z}_n$ with any **SV** (\mathbf{z}_n, y_n) 任一方 SV.

SVM: learn **fattest hyperplane** \odot
by identifying support vectors \Rightarrow 找 SV \odot 用 SV 找 optimal
with **dual** optimal solution

Summary: Two Forms of Hard-Margin SVM

Primal Hard-Margin SVM

原始的

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

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

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

Dual Hard-Margin SVM

$$\min_{\alpha} \quad \frac{1}{2} \alpha^T Q_D \alpha - \mathbf{1}^T \alpha$$

$$\text{s.t.} \quad \mathbf{y}^T \alpha = 0;$$

$$\alpha_n \geq 0 \text{ for } n = 1, \dots, N$$

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

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

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

≠ 用 (b, \mathbf{w}) 做分类的 prediction.

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!

如何避开 \tilde{d} 的计算

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