

Machine Learning (機器學習)

Lecture 06: Beyond Basic Linear Models

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)



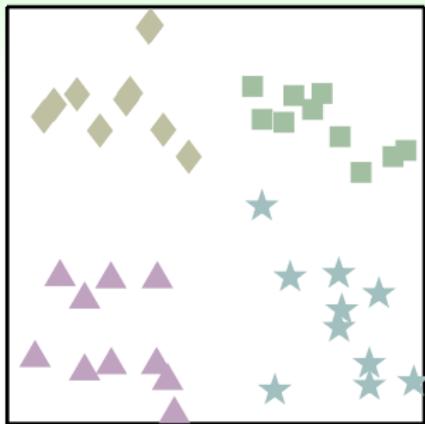
Roadmap

- ① When Can Machines Learn?
- ② Why Can Machines Learn?
- ③ **How** Can Machines Learn?

Lecture 06: Beyond Basic Linear Models

- Multiclass via Logistic Regression
- Multiclass via Binary Classification
- Quadratic Hypotheses
- Nonlinear Transform
- Price of Nonlinear Transform
- Structured Hypothesis Sets

Multiclass Classification



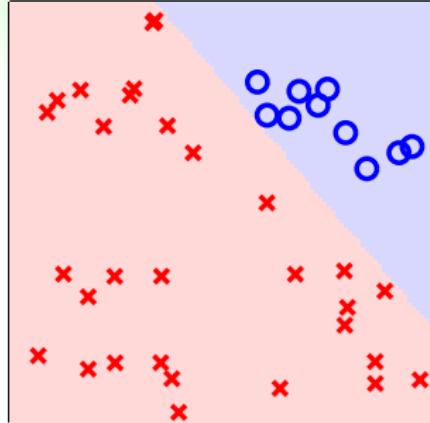
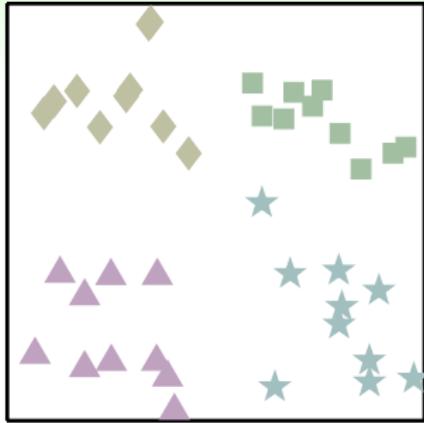
多選擇題

- $\mathcal{Y} = \{\square, \diamond, \triangle, \star\}$
(4-class classification)
- **many applications** in practice, especially for 'recognition'

next: use **tools for $\{\times, \circ\}$ classification** to
 $\{\square, \diamond, \triangle, \star\}$ classification

Binary classification

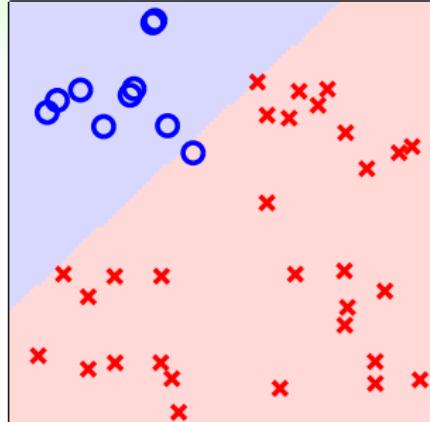
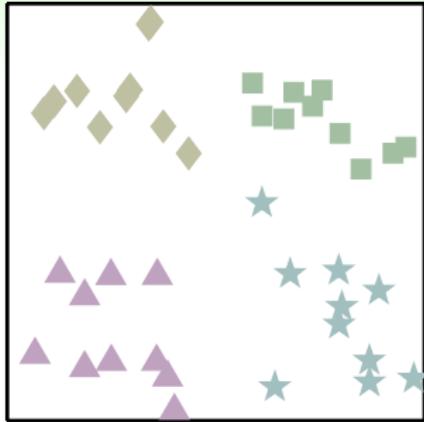
One Class at a Time



□ or not? {□ = ○, ◊ = ✕, △ = ✕, ⋆ = ✕}

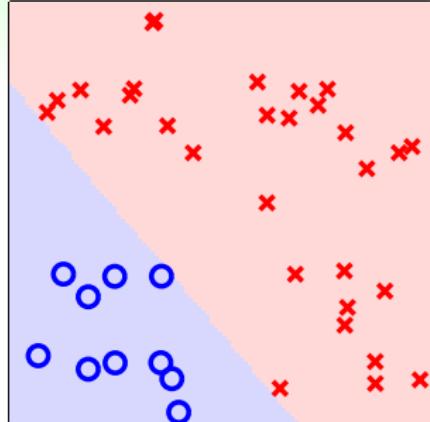
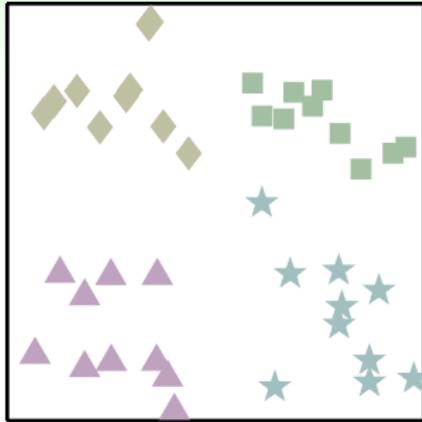
↑ 把口跟其他人分开

One Class at a Time



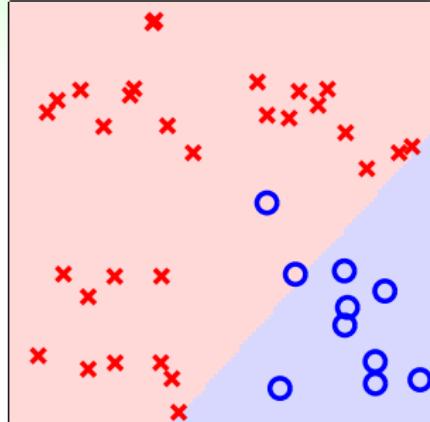
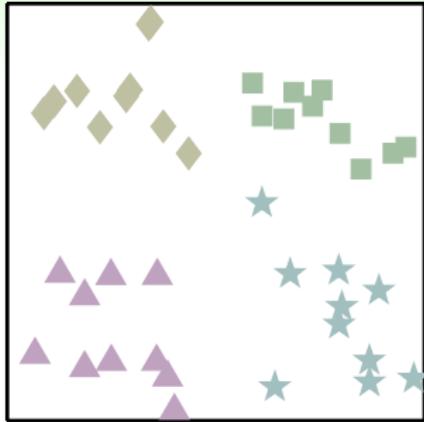
◊ or not? {□ = x, ◊ = o, △ = x, ★ = x}

One Class at a Time



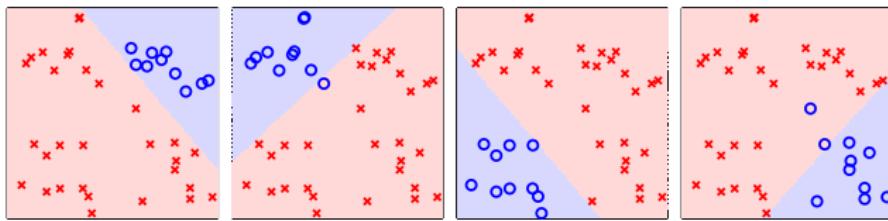
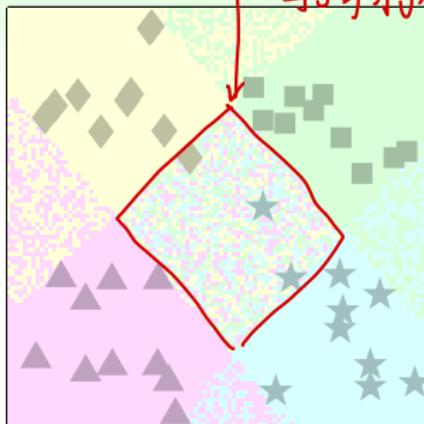
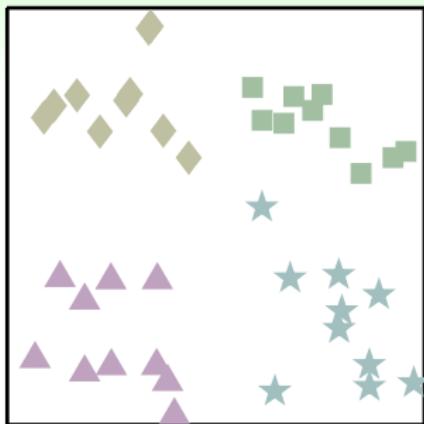
\triangle or not? $\{\square = \times, \diamond = \times, \triangle = \circ, \star = \times\}$

One Class at a Time



* or not? $\{\square = \times, \diamond = \times, \triangle = \times, \star = \circ\}$

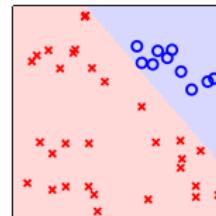
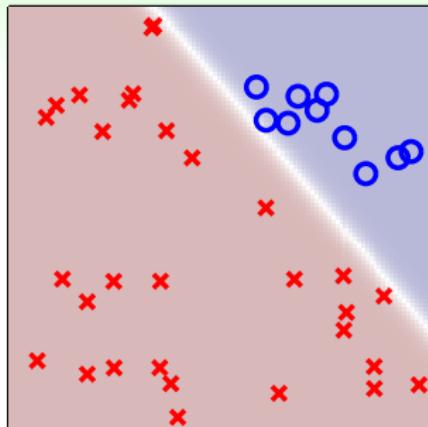
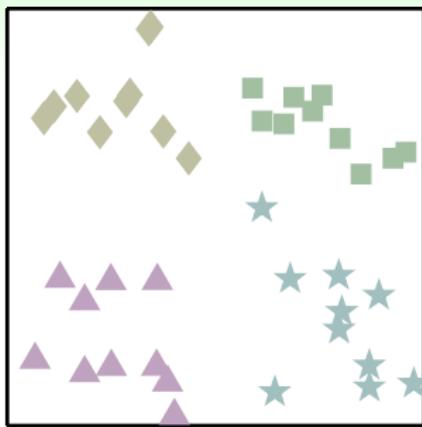
Multiclass Prediction: Combine Binary Classifiers



四分類器？

but ties? :-)

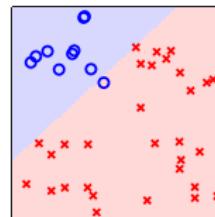
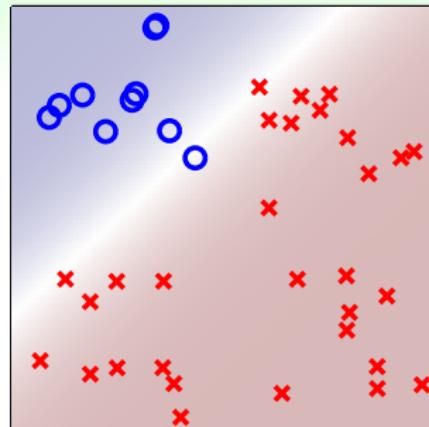
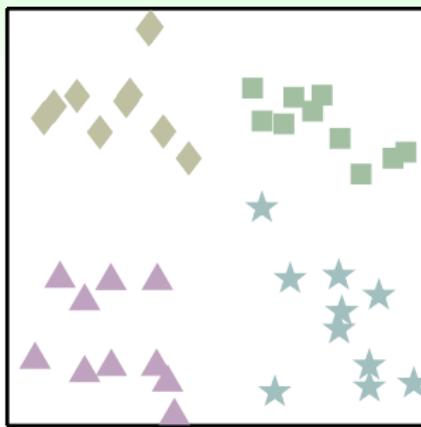
One Class at a Time **Softly**



logistic regression
未估計

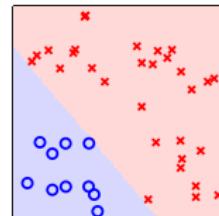
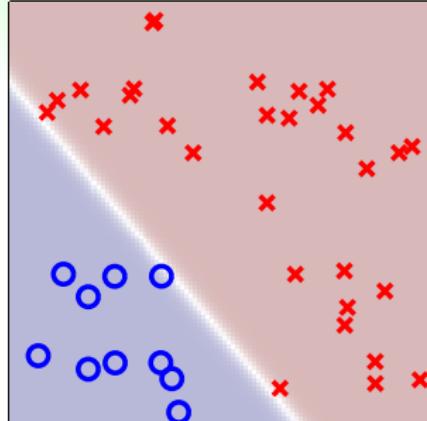
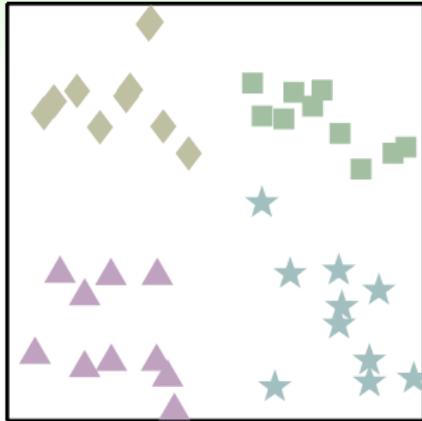
$$P(\square|x)? \{ \square = \circ, \diamond = \times, \triangle = \times, \star = \times \}$$

One Class at a Time **Softly**



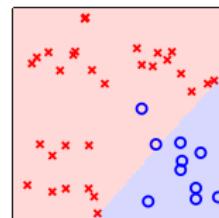
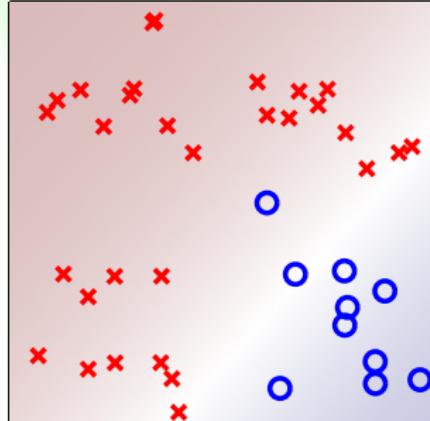
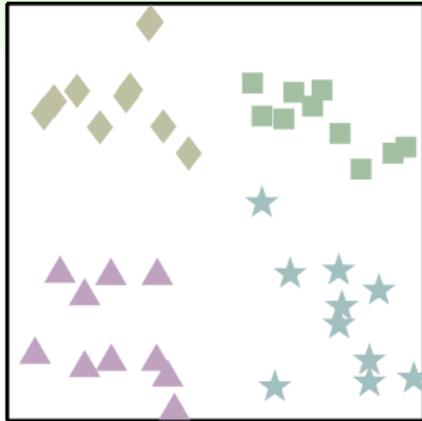
$P(\diamond|x)$? $\{\square = \times, \diamond = \circ, \triangle = \times, \star = \times\}$

One Class at a Time **Softly**



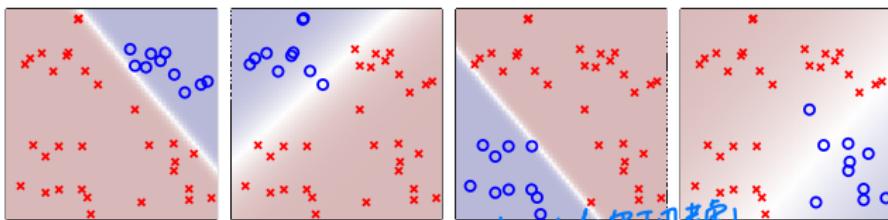
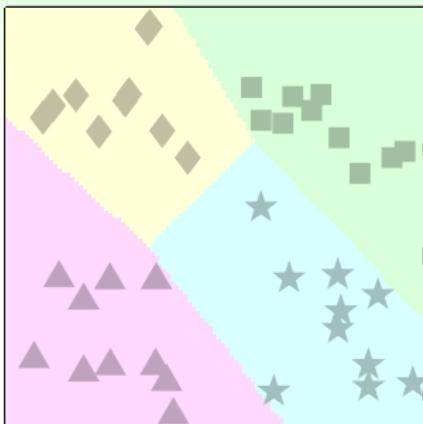
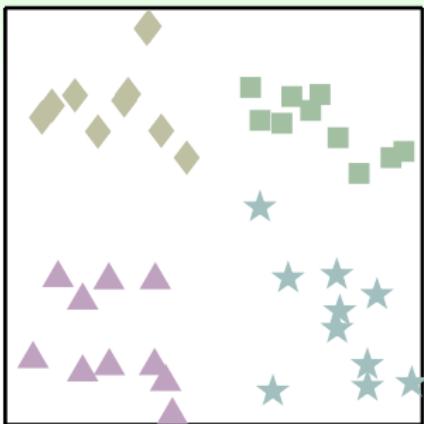
$$P(\Delta | \mathbf{x})? \{ \square = \times, \diamond = \times, \triangle = \circ, \star = \times \}$$

One Class at a Time **Softly**



$P(\star|\mathbf{x})?$ $\{\square = \times, \diamond = \times, \triangle = \times, \star = \circ\}$

Multiclass Prediction: Combine Soft Classifiers



$$g(\mathbf{x}) = \operatorname{argmax}_{k \in \mathcal{Y}} \theta\left(\mathbf{w}^T_{[k]} \mathbf{x} \right)$$

Sigmoid (巴布爾函數)

選可能性最大者

One-Versus-All (OVA) Decomposition

① for $k \in \mathcal{Y}$

obtain $\mathbf{w}_{[k]}$ by running logistic regression on

$$\mathcal{D}_{[k]} = \{(\mathbf{x}_n, y'_n = 2[\underline{y_n = k}] - 1)\}_{n=1}^N$$

N training data.
是不是這 class

② return $g(\mathbf{x}) = \operatorname{argmax}_{k \in \mathcal{Y}} (\mathbf{w}_{[k]}^T \mathbf{x})$ 选最有可能的 class

- pros: efficient,
can be coupled with any logistic regression-like approaches
- cons: often unbalanced $\mathcal{D}_{[k]}$ when K large ← 如果 training data 大部份是 X , 則會 train 不好.
- extension: multinomial ('coupled') logistic regression

多類別

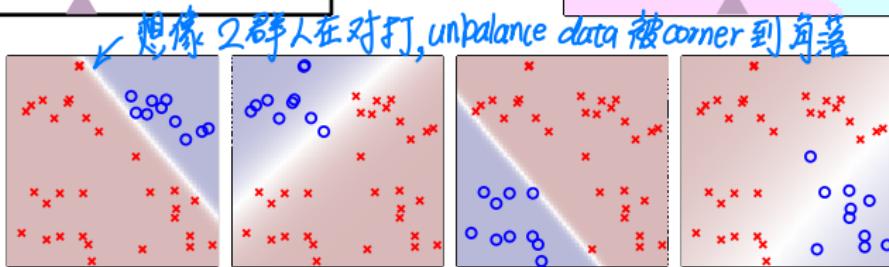
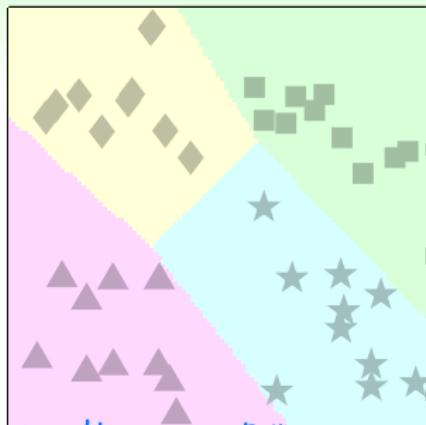
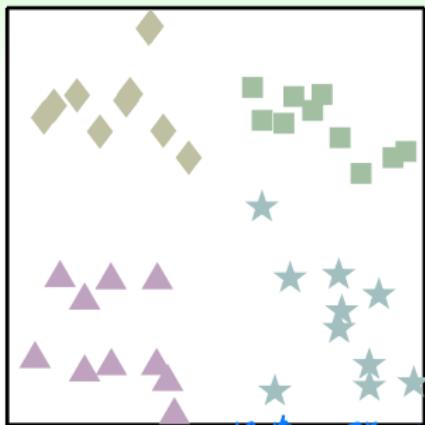
OVA: a simple multiclass meta-algorithm

to keep in your toolbox

非常常用

Questions?

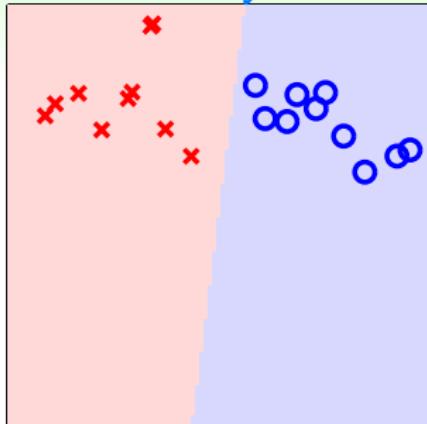
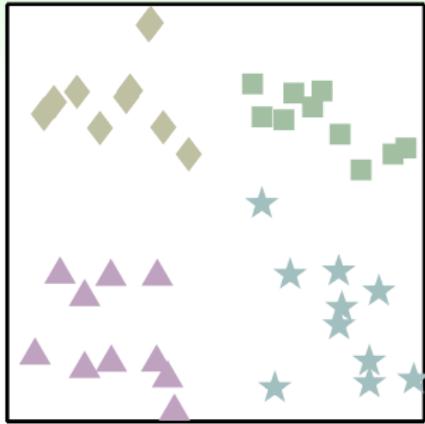
Source of Unbalance: One versus All



idea: make binary classification problems
more **balanced** by one versus **one**

One versus One at a Time

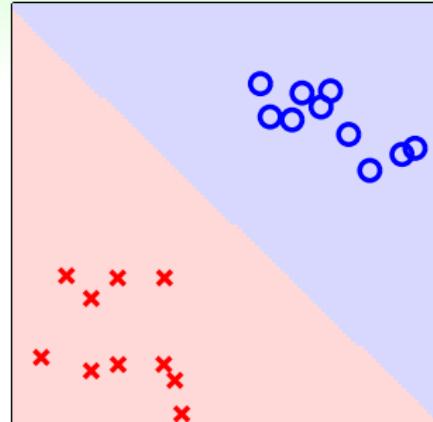
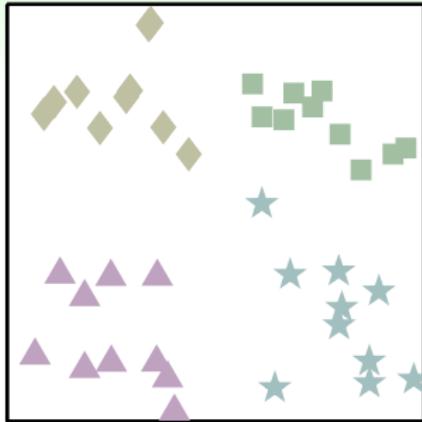
所有类选 2 个 class
↓ 随机对打



□ or ◇? {□ = ○, ◇ = ✕, △ = nil, ★ = nil}

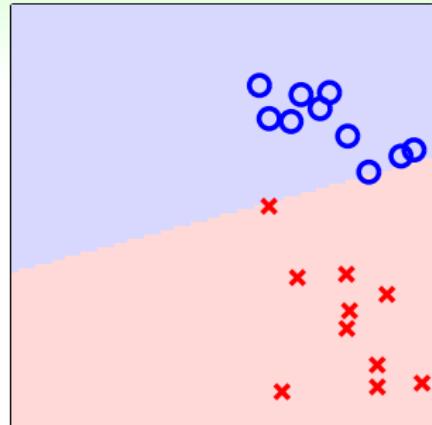
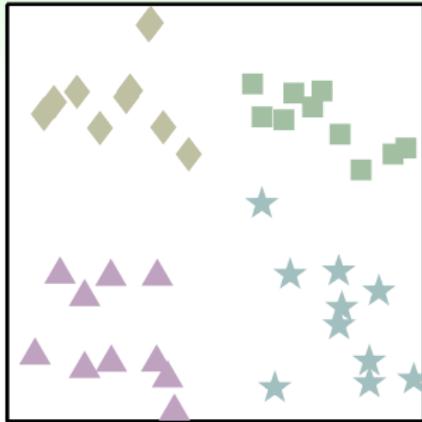
只用 2 个 classes

One versus One at a Time



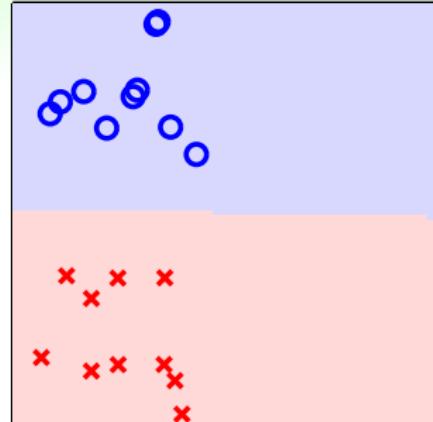
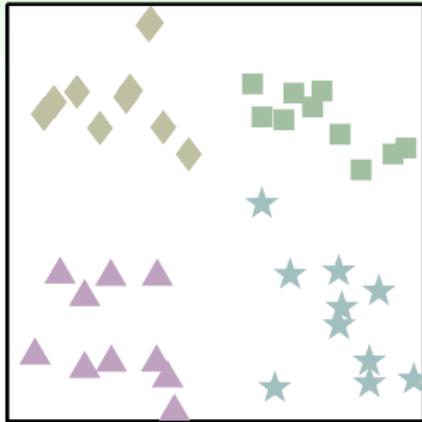
\square or \triangle ? $\{\square = \circ, \diamond = \text{nil}, \triangle = \times, \star = \text{nil}\}$

One versus One at a Time



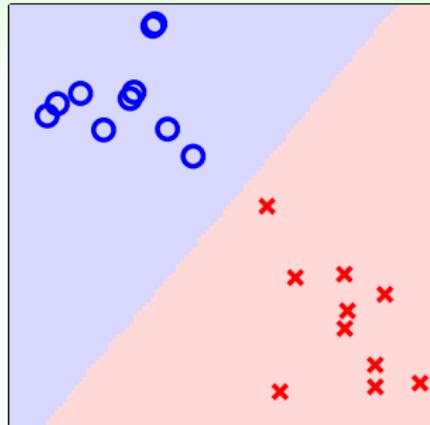
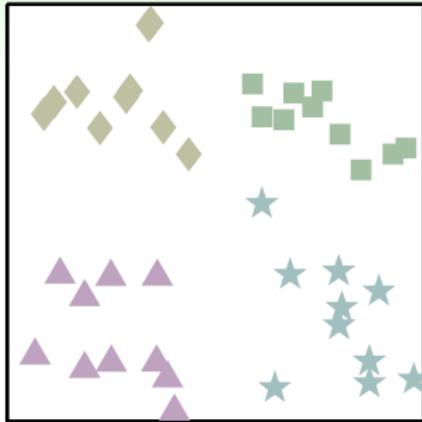
□ or ⋆? {□ = ○, ◊ = nil, △ = nil, ⋆ = ✕}

One versus One at a Time



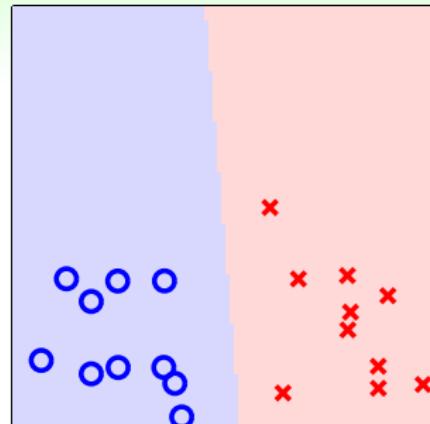
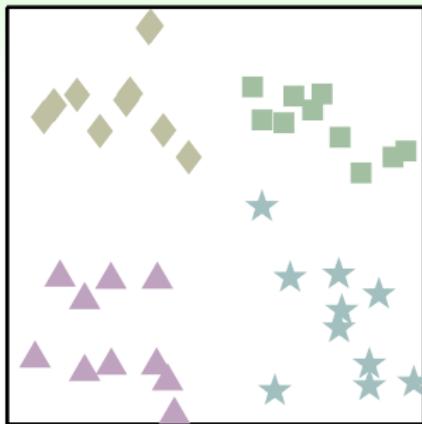
◊ or △? {□ = nil, ◊ = o, △ = x, * = nil}

One versus One at a Time



◊ or ⋆? {□ = nil, ◊ = ○, △ = nil, ⋆ = ✕}

One versus One at a Time

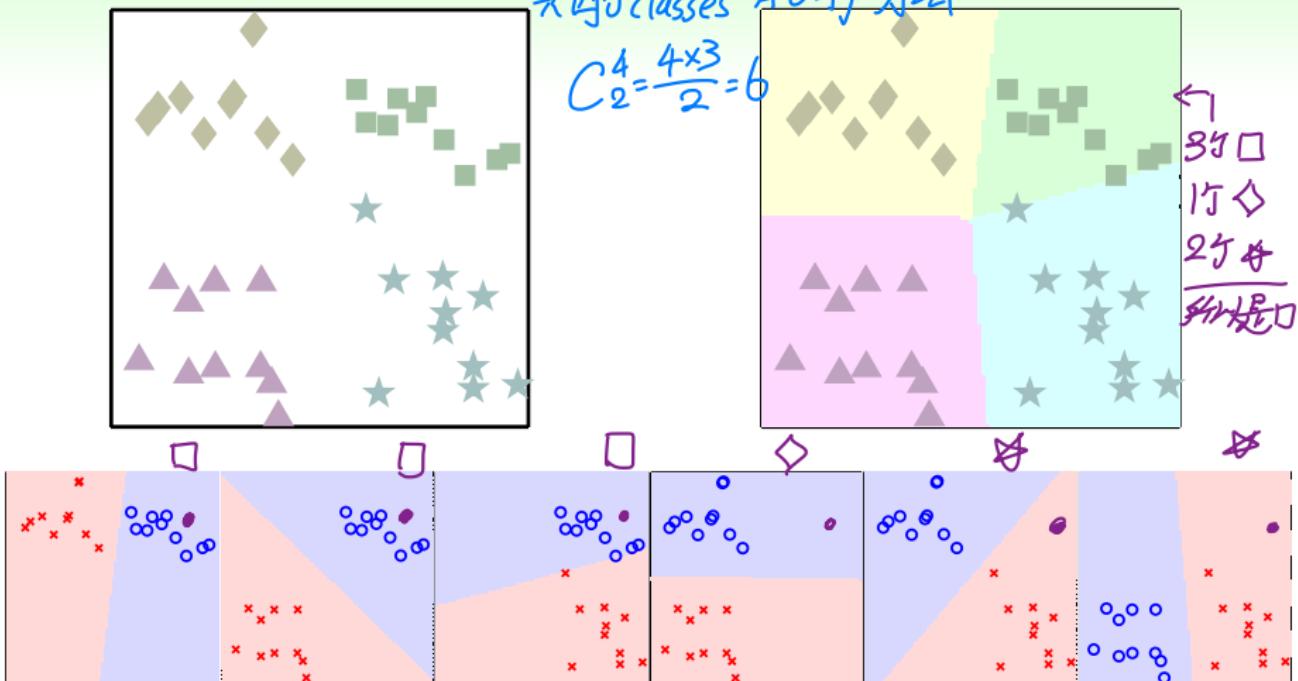


\triangle or \star ? $\{\square = \text{nil}, \diamond = \text{nil}, \triangle = \circ, \star = \times\}$

Multiclass Prediction: Combine **Pairwise** Classifiers

共四 classes 有 6 种方法

$$C_2^4 = \frac{4 \times 3}{2} = 6$$



$g(\mathbf{x}) = \text{tournament champion} \left\{ \mathbf{w}_{[k,\ell]}^T \mathbf{x} \right\}$
 (voting of classifiers)

One-versus-one (OVO) Decomposition

- for $(k, \ell) \in \mathcal{Y} \times \mathcal{Y}$
obtain $\mathbf{w}_{[k, \ell]}$ by running linear binary classification on

$$\mathcal{D}_{[k, \ell]} = \{(\mathbf{x}_n, y'_n = 2[\![y_n = k]\!] - 1) : y_n = k \text{ or } y_n = \ell\}$$

- return $g(\mathbf{x}) = \underbrace{\text{tournament champion}}_{\text{投票最高者}} \left\{ \mathbf{w}_{[k, \ell]}^T \mathbf{x} \right\}$ *25 classes*

- pros: efficient ('smaller' training problems), stable,
can be coupled with any binary classification approaches
- cons: use $O(K^2)$ $\mathbf{w}_{[k, \ell]}$ 雖然要 train 多个 classifier, 但每个 classifier
—**more space, slower prediction, more training** *training data 很少*

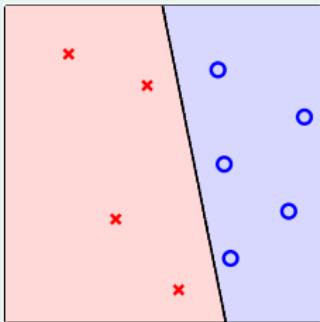
常用在尺小大的情形下.

OVO: another simple multiclass
meta-algorithm to keep in your toolbox

Questions?

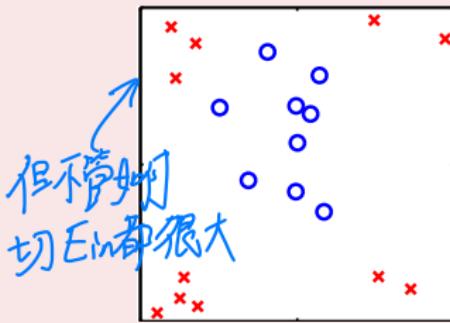
Linear Hypotheses

up to now: linear hypotheses



- visually: 'line-like' boundary
- mathematically: linear scores $s = \mathbf{w}^T \mathbf{x}$

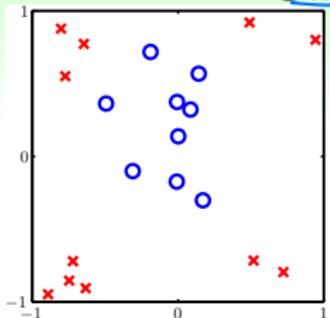
but limited . . .



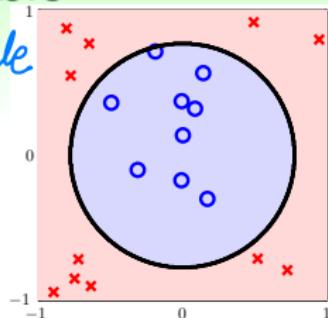
- theoretically: d_{VC} under control :-)
- practically: on some \mathcal{D} , large E_{in} for every line :-(

好處是 VC dimension 很小, E_{in}, E_{out} 不會差太多

how to **break the limit** of linear hypotheses

Circular Separable

不是 linear Separable



- \mathcal{D} not linear separable 圓周可分
- but circular separable by a circle of radius $\sqrt{0.6}$ centered at origin: circle formula

$$h_{\text{SEP}}(\mathbf{x}) = \text{sign} \left(-x_1^2 - x_2^2 + 0.6 \right)$$

↑ 到圓心的距離

re-derive Circular-PLA, Circular-Regression,
blahblah ... all over again? :-)

Circular Separable and Linear Separable

$$\begin{aligned}
 h(\mathbf{x}) &= \text{sign} \left(\underbrace{0.6}_{\tilde{w}_0} \cdot \underbrace{1}_{z_0} + \underbrace{(-1)}_{\tilde{w}_1} \cdot \underbrace{x_1^2}_{z_1} + \underbrace{(-1)}_{\tilde{w}_2} \cdot \underbrace{x_2^2}_{z_2} \right) \\
 &= \text{sign} \left(\tilde{\mathbf{w}}^T \underline{\underline{z}} \right)
 \end{aligned}$$

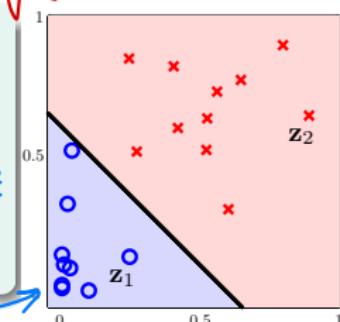
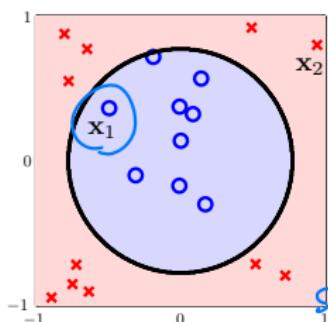
半徑²

data.

圓圓可分 轉到 z-space 變線性可分

- $\{(\mathbf{x}_n, y_n)\}$ circular separable
 $\Rightarrow \{(\mathbf{z}_n, y_n)\}$ linear separable
- $\mathbf{x} \in \mathcal{X} \xrightarrow{\phi} \mathbf{z} \in \mathcal{Z}$:
(nonlinear) feature transform ϕ

轉換到 z-space.



circular separable in $\mathcal{X} \Rightarrow$ linear separable in \mathcal{Z}
vice versa?

Linear Hypotheses in \mathcal{Z} -Space

$$(z_0, z_1, z_2) = \mathbf{z} = \underline{\Phi(\mathbf{x})} = (1, x_1^2, x_2^2)$$

$$h(\mathbf{x}) = \tilde{h}(\mathbf{z}) = \text{sign} \left(\underline{\tilde{\mathbf{w}}^T \Phi(\mathbf{x})} \right) = \text{sign} \left(\tilde{w}_0 + \tilde{w}_1 x_1^2 + \tilde{w}_2 x_2^2 \right)$$

$\tilde{\mathbf{w}} = (\tilde{w}_0, \tilde{w}_1, \tilde{w}_2)$

- $(0.6, -1, -1)$: circle (inside)
- $(-0.6, +1, +1)$: circle (outside)
- $(0.6, -1, -2)$: ellipse. $\text{sign}(0.6 - x_1^2 - 2x_2^2)$
- $(0.6, -1, +2)$: hyperbola
- $(0.6, +1, +2)$: constant $\circ :-)$ $\text{sign}(0.6 + x_1^2 + x_2^2)$ ←永遠是正的

lines in \mathcal{Z} -space

\iff special quadratic curves in \mathcal{X} -space

但圆心一定要过原点 !?

General Quadratic Hypothesis Set

a 'bigger' \mathcal{Z} -space with $\Phi_2(\mathbf{x}) = (1, \underline{x_1}, \underline{x_2}, \underline{x_1^2}, \underline{x_1x_2}, \underline{x_2^2})$

perceptrons in \mathcal{Z} -space \iff quadratic hypotheses in \mathcal{X} -space

$$\mathcal{H}_{\Phi_2} = \left\{ h(\mathbf{x}) : h(\mathbf{x}) = \tilde{h}(\Phi_2(\mathbf{x})) \text{ for some linear } \underline{\tilde{h}} \text{ on } \mathcal{Z} \right\}$$

2R

Linear 的 perception.

- can implement all possible quadratic curve boundaries:
circle, ellipse, rotated ellipse, hyperbola, parabola, ...

ellipse $2(x_1 + x_2 - 3)^2 + (x_1 - x_2 - 4)^2 = 1$

$$\Leftrightarrow \underline{\underline{\mathbf{w}^T}} = [3, 3, -20, -4, 1, 3] \quad \Rightarrow 2(x_1^2 + x_2^2 + 9 + 2x_1x_2 - 6x_2 - 6x_1) + x_1^2 + x_2^2 + (6 - 2x_1x_2 - 8x_1 + 8x_2 = 1)$$

- include lines and constants as degenerate cases 但都能表示

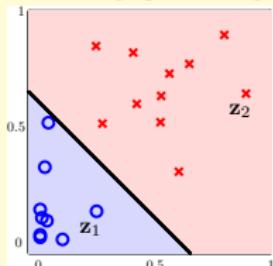
$$\Rightarrow \underline{\underline{3x_1^2}} + \underline{\underline{3x_2^2}} + \underline{\underline{2x_1x_2}} + \underline{\underline{-20x_1}} + \underline{\underline{-4x_2}} + \underline{\underline{33}} = 0$$

next: learn a good quadratic hypothesis g

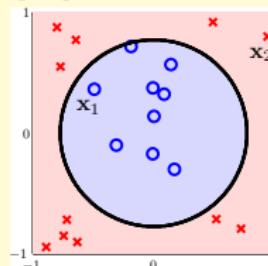
Questions?

Good Quadratic Hypothesis

\mathcal{Z} -space
perceptrons
good perceptron
separating perceptron



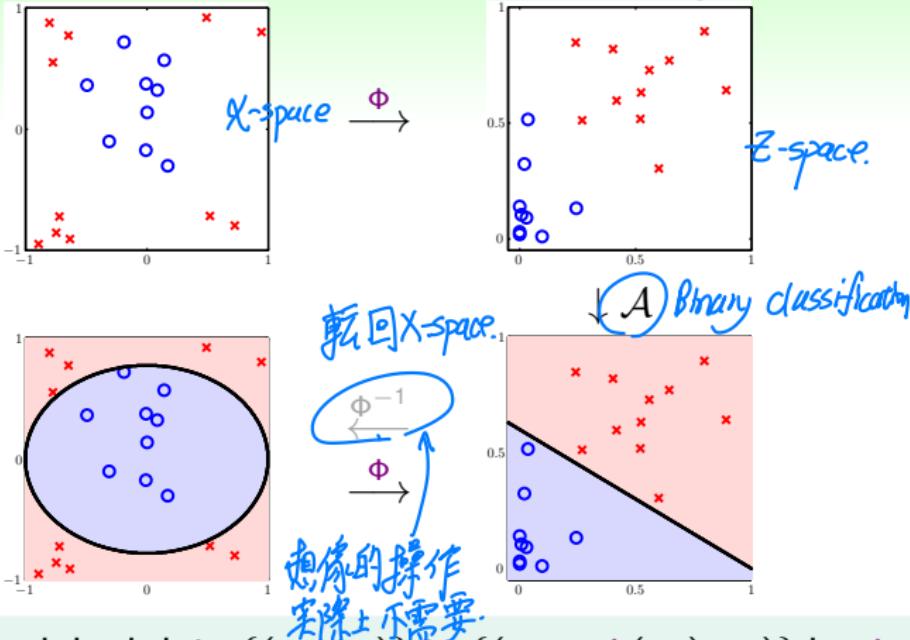
\mathcal{X} -space
quadratic hypotheses
good quadratic hypothesis
separating quadratic hypothesis



- want: get **good perceptron** in \mathcal{Z} -space ← 在 \mathcal{Z} space 找一條好的 line.
- known: get **good perceptron** in \mathcal{X} -space with data $\{(\mathbf{x}_n, y_n)\}$

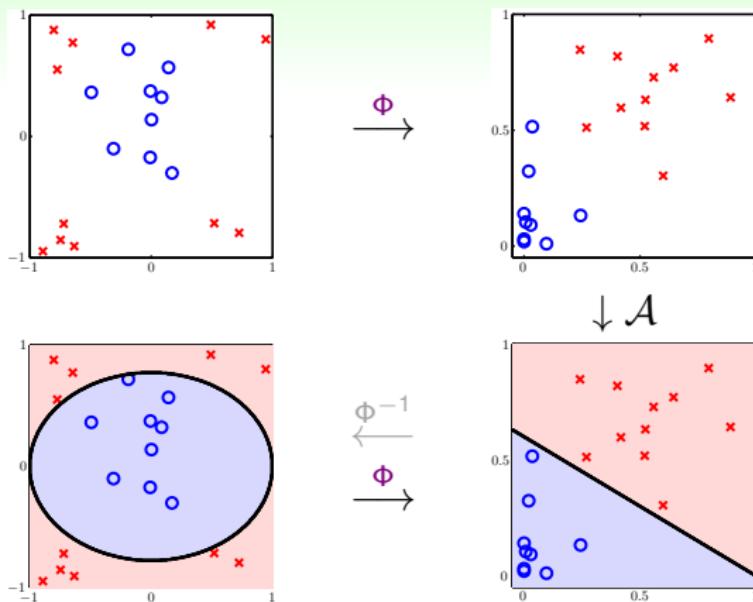
todo: get **good perceptron** in \mathcal{Z} -space with data $\{(\mathbf{z}_n = \Phi_2(\mathbf{x}_n), y_n)\}$

The Nonlinear Transform Steps



- 1 transform original data $\{(\mathbf{x}_n, y_n)\}$ to $\{(\mathbf{z}_n = \Phi(\mathbf{x}_n), y_n)\}$ by Φ
- 2 get a good perceptron $\tilde{\mathbf{w}}$ using $\{(\mathbf{z}_n, y_n)\}$ and your favorite linear classification algorithm \mathcal{A}
- 3 return $g(\mathbf{x}) = \text{sign}(\tilde{\mathbf{w}}^T \Phi(\mathbf{x}))$

Nonlinear Model via Nonlinear Φ + Linear Models



two choices:

- feature transform
 Φ How to transform.
- linear model \mathcal{A} ,
not just binary classification

Pandora's box :-):

全部都可以做了

can now freely do **quadratic PLA, quadratic regression, cubic regression, . . . , polynomial regression**

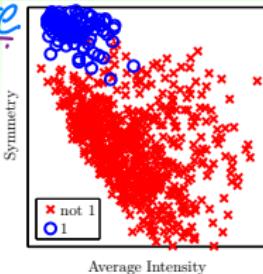
Feature Transform Φ

concert and raw feature

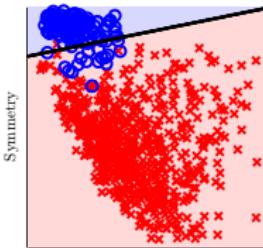
it's kind of concert feature



$$\Phi \rightarrow$$



$$\begin{matrix} \Phi^{-1} \\ \Phi \end{matrix}$$



not new, not just polynomial:

raw (pixels) $\xrightarrow{\text{domain knowledge}}$ concrete (intensity, symmetry)

原力

the force, too good to be true? :-)

Questions?

Computation/Storage Price

Q -th order polynomial transform: $\Phi_Q(\mathbf{x}) = (\quad 1,$

$x_1, x_2, \dots, x_d, \leftarrow -\mathcal{R}$

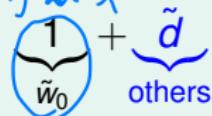
$x_1^2, x_1 x_2, \dots, x_d^2, \leftarrow 2\mathcal{R}$

...,

$Q\mathcal{R}$

$x_1^Q, x_1^{Q-1}x_2, \dots, x_d^Q)$

常數項



d dimensions

d 种不同的項取 Q 个出来

= # ways of $\leq Q$ -combination from d kinds with repetitions

$$= \binom{Q+d}{Q} = \binom{Q+d}{d} = O(Q^d)$$

有幾種取法

計算繁多

= efforts needed for computing/storing $\mathbf{z} = \Phi_Q(\mathbf{x})$ and $\tilde{\mathbf{w}}$

參數變多

Q large \Rightarrow difficult to compute/store

Model Complexity Price

Q -th order polynomial transform: $\Phi_Q(\mathbf{x}) = \begin{pmatrix} 1, \\ x_1, x_2, \dots, x_d, \\ x_1^2, x_1 x_2, \dots, x_d^2, \\ \dots, \\ x_1^Q, x_1^{Q-1} x_2, \dots, x_d^Q \end{pmatrix}$

這是項次的數量 (# of term)

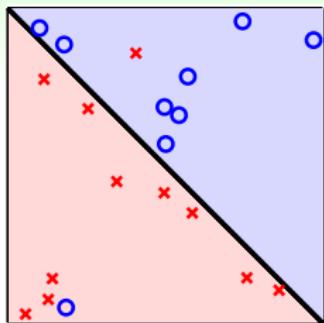
$$\underbrace{1}_{\tilde{w}_0} + \underbrace{\tilde{d}}_{\text{others}} \text{ dimensions} = O(Q^d)$$

↓ 自由度變大了

- number of free parameters $\tilde{w}_i = \tilde{d} + 1 \approx d_{VC}(\mathcal{H}_{\Phi_Q})$
- $d_{VC}(\mathcal{H}_{\Phi_Q}) \leq \tilde{d} + 1$, why?
any $\tilde{d} + 2$ inputs not shattered in \mathcal{Z}
 \implies any $\tilde{d} + 2$ inputs not shattered in \mathcal{X}

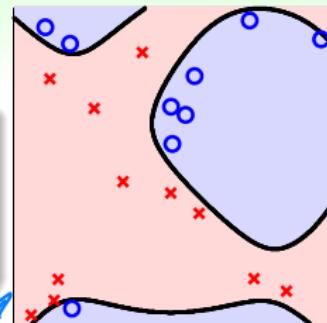
Q large \implies large d_{VC}

Generalization Issue

 Φ_1 (original \mathbf{x})

which one do you prefer? :-)

- Φ_1 'visually' preferred
- Φ_4 : $E_{in}(g) = 0$ but overkill



overfitting? $\Phi_4 \quad Q=4$

- ① can we make sure that $E_{out}(g)$ is close enough to $E_{in}(g)$?
- ② can we make $E_{in}(g)$ small enough?

又来了!

trade-off:

$\tilde{d}(Q)$	1	2
higher	:-)	:-D
lower	-D	:-)

how to pick Q ? visually, maybe?

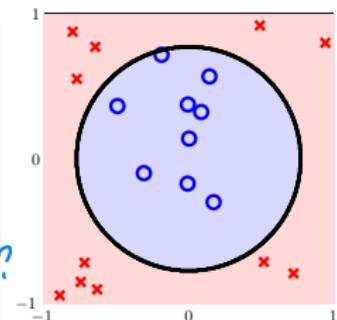
Danger of Visual Choices

first of all, can you really ‘visualize’ when $\mathcal{X} = \mathbb{R}^{10}$? (well, I can’t :-))

Visualize $\mathcal{X} = \mathbb{R}^2$

- full Φ_2 : $\mathbf{z} = (1, x_1, x_2, x_1^2, x_1x_2, x_2^2)$, $d_{VC} = 6$
- or $\mathbf{z} = (1, x_1^2, x_2^2)$, $d_{VC} = 3$, after visualizing?
- or better $\mathbf{z} = (1, x_1^2 + x_2^2)$, $d_{VC} = 2$?
- or even better $\mathbf{z} = (\text{sign}(0.6 - x_1^2 - x_2^2))$? $d_{VC} = 1$?

↑ careful about your brain’s ‘model complexity’
 ↗ 是 human learning



做出來的低, d_{VC} 之後, 但真的好嗎?

for VC-safety, Φ shall be
 decided without ‘peeking’ data

↑ 人類偷看

Questions?

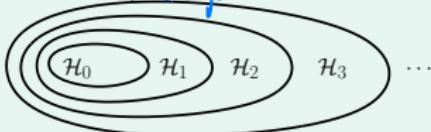
Polynomial Transform Revisited *recusssive*

$$\begin{aligned}
 \Phi_0(\mathbf{x}) &= \left(\underset{0>\mathbb{R}}{1}, \underset{0>\mathbb{R}}{\Phi_0(\mathbf{x})} \right), & x_1, x_2, \dots, x_d \\
 \Phi_1(\mathbf{x}) &= \left(\Phi_1(\mathbf{x}), \quad x_1^2, x_1x_2, \dots, x_d^2 \right) \\
 \Phi_2(\mathbf{x}) &= \left(\Phi_2(\mathbf{x}), \quad x_1^3, x_1^2x_2, \dots, x_d^3 \right) \\
 &\dots & \dots \\
 \Phi_Q(\mathbf{x}) &= \left(\Phi_{Q-1}(\mathbf{x}), \quad x_1^Q, x_1^{Q-1}x_2, \dots, x_d^Q \right)
 \end{aligned}$$

$$\mathcal{H}_{\Phi_0} \subset \mathcal{H}_{\Phi_1} \subset \mathcal{H}_{\Phi_2} \subset \mathcal{H}_{\Phi_3} \subset \dots \subset \mathcal{H}_{\Phi_Q}$$

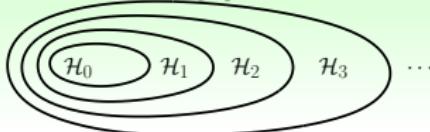
\parallel \parallel \parallel \parallel \parallel
 \mathcal{H}_0 \mathcal{H}_1 \mathcal{H}_2 \mathcal{H}_3 \dots \mathcal{H}_Q

高維 \rightarrow 低維
Included



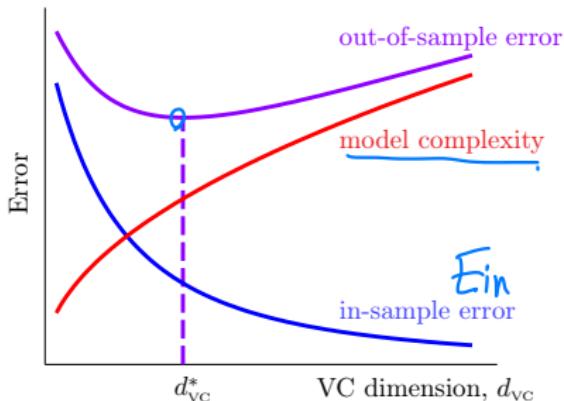
structure: **nested \mathcal{H}_i**

Structured Hypothesis Sets



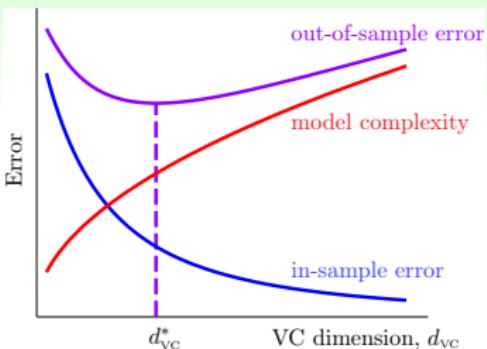
Let $g_i = \operatorname{argmin}_{h \in \mathcal{H}_i} E_{\text{in}}(h)$:

$$\begin{array}{ccccccc} \mathcal{H}_0 & \subset & \mathcal{H}_1 & \subset & \mathcal{H}_2 & \subset & \mathcal{H}_3 & \subset & \dots \\ d_{\text{VC}}(\mathcal{H}_0) & \leq & d_{\text{VC}}(\mathcal{H}_1) & \leq & d_{\text{VC}}(\mathcal{H}_2) & \leq & d_{\text{VC}}(\mathcal{H}_3) & \leq & \dots \\ E_{\text{in}}(g_0) & \geq & E_{\text{in}}(g_1) & \geq & E_{\text{in}}(g_2) & \geq & E_{\text{in}}(g_3) & \geq & \dots \end{array}$$



use \mathcal{H}_{1126} won't be good! :-(

Linear Model First



- tempting sin. use \mathcal{H}_{1126} , low $E_{in}(g_{1126})$ to fool your boss
—really? :-(a dangerous path of no return *dark side of the force*
- safe route: \mathcal{H}_1 first ←先從低維做起
 - if $E_{in}(g_1)$ good enough, live happily thereafter :-)
 - otherwise, move right of the curve
with nothing lost except 'wasted' computation

linear model first: good side of the force
simple, efficient, **safe**, and workable!

Questions?

Summary

① How Can Machines Learn?

Lecture 05: Linear Models

Lecture 06: Beyond Basic Linear Models

- Multiclass via Logistic Regression
predict with maximum estimated $P(k|x)$

- Multiclass via Binary Classification
predict the tournament champion

- Quadratic Hypotheses

linear hypotheses on quadratic-transformed data

- Nonlinear Transform
happy linear modeling after $\mathcal{Z} = \Phi(\mathcal{X})$

- Price of Nonlinear Transform
computation/storage/[model complexity]

- Structured Hypothesis Sets
linear/simpler model first

- next: dark side of the force :-)