

Three Problems

Problem 1: Evaluation

- What does $F(x,y)$ look like?

Problem 2: Inference

- How to solve the “arg max” problem?

$$y = \arg \max_x F(x,y)$$

Problem 3: Training

- Given training data, how to find the best model?

Have you heard the three problems elsewhere?

Hidden Markov Model

• Three Basic Problems for HMMs

Given an observation sequence $\bar{O}=(o_1,o_2,\dots,o_T)$, and an HMM

$\lambda=(A,B,\pi)$

– Problem 1 :

How to *efficiently* compute $P(\bar{O}|\lambda)$?

\Rightarrow *Evaluation problem*

– Problem 2 :

How to choose an optimal state sequence $\mathbf{q}=(q_1,q_2,\dots,q_T)$?

\Rightarrow *Decoding Problem*

– Problem 3 :

Given some observations \bar{O} for the HMM λ , how to adjust the model parameter $\lambda=(A,B,\pi)$ to maximize $P(\bar{O}|\lambda)$?

\Rightarrow *Learning /Training Problem*

HMM其實就是structure learning的一個特例

From 數位語音處理

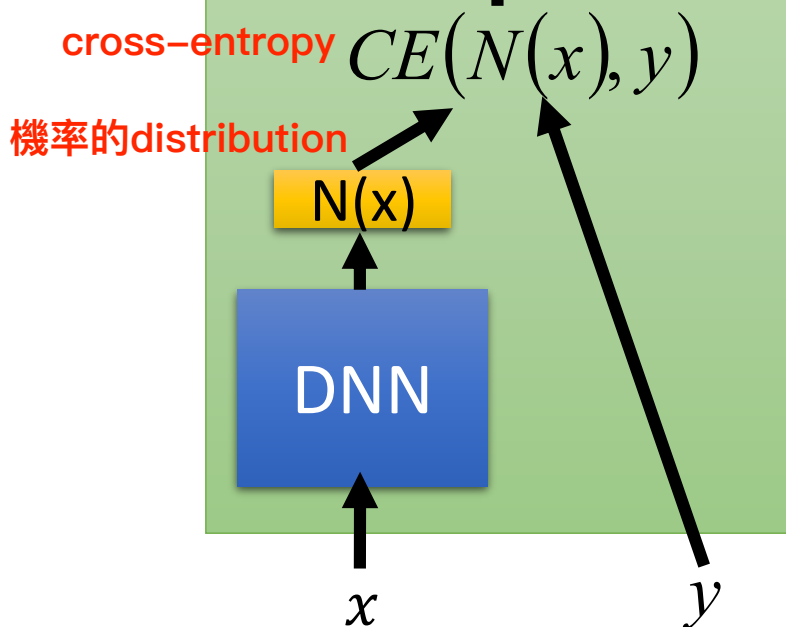
Link to DNN?

The same as what we have learned.

Training

$$F : X \times Y \rightarrow \mathbb{R}$$
$$F(x, y) = -CE(N(x), y)$$

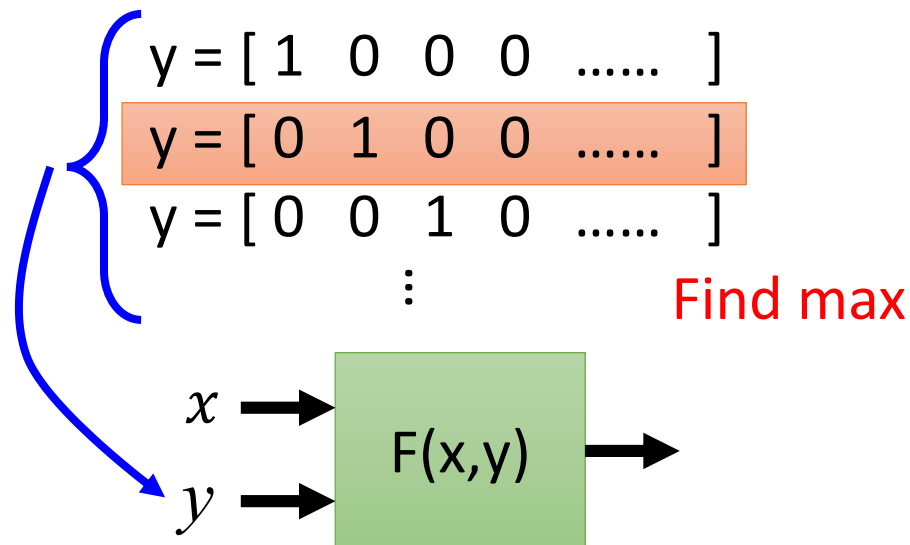
CE要越小越好



Inference

$$\tilde{y} = \arg \max_{y \in Y} F(x, y)$$

In handwriting digit classification, there are only 10 possible y .



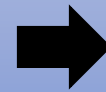
Introduction of Structured Learning Linear Model

Structured Linear Model

做出一個限制：假設evaluation function是linear

Problem 1: Evaluation

- What does $F(x, y)$ look like?



in a specific form

Problem 2: Inference

- How to solve the “arg max” problem

$$y = \arg \max_{y \in Y} F(x, y)$$

Problem 3: Training

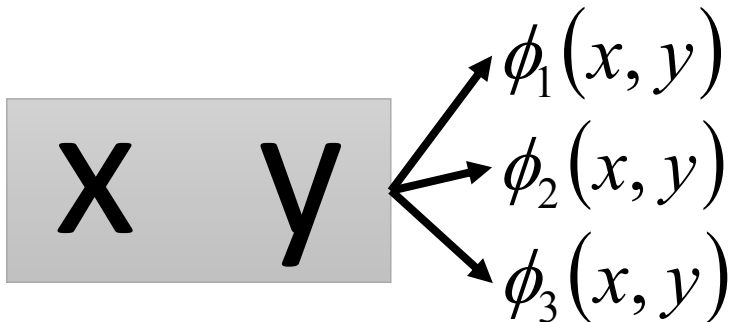
- Given training data, how to find $F(x, y)$

Structured Linear Model: Problem 1

限制假設為linear

- **Evaluation:** What does $F(x, y)$ look like?

Characteristics



抽取不同的feature做比較

$$F(x, y) = \boxed{w_1} \cdot \phi_1(x, y) \\ + \boxed{w_2} \cdot \phi_2(x, y) \\ + \boxed{w_3} \cdot \phi_3(x, y) \dots$$

Learning
from data

$$F(x, y) = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w \end{bmatrix} \cdot \begin{bmatrix} \phi_1(x, y) \\ \phi_2(x, y) \\ \phi_3(x, y) \\ \vdots \\ \phi(x, y) \end{bmatrix}$$


↓

$$F(x, y) = w \cdot \phi(x, y)$$

Structured Linear Model: Problem 1

- **Evaluation**: What does $F(x,y)$ look like?
- Example: **Object Detection**

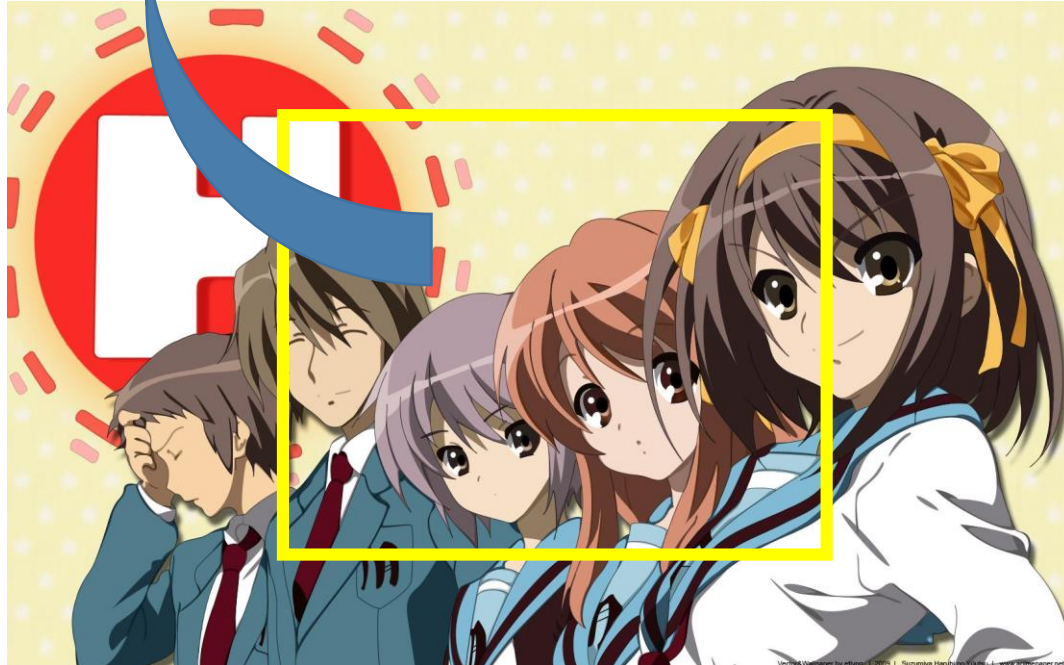
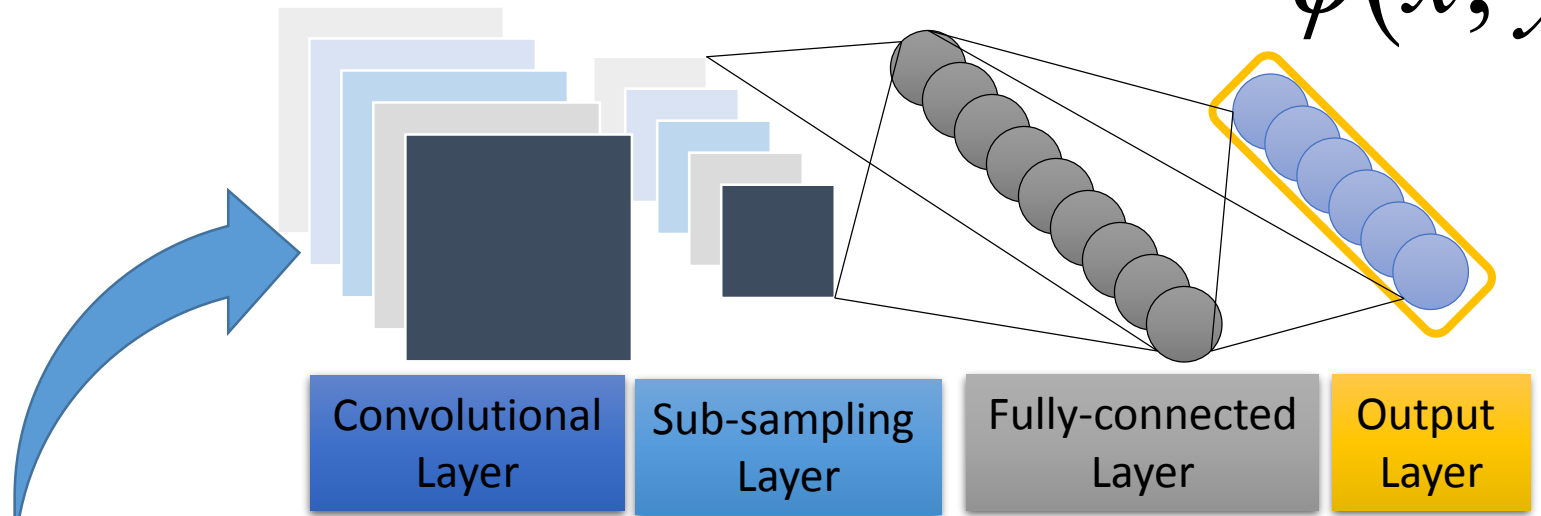
$\phi($



$) =$

- percentage of color red in box y
- percentage of color green in box y
- percentage of color blue in box y
- percentage of color red out of box y
-
- area of box y
- number of specific patterns in box y
-

$$\phi(x, y)$$



$\phi($)

Structured Linear Model:

Problem 2

- **Inference:** How to solve the “arg max” problem

$$y = \arg \max_{y \in Y} F(x, y)$$

$$F(x, y) = w \cdot \phi(x, y) \Rightarrow y = \arg \max_{y \in Y} w \cdot \phi(x, y)$$

- Assume we have solved this question.

假設第一步是linear第二步也解出來則第三步是很簡單的

Structured Linear Model:

Problem 3

- Training: Given training data, how to learn $F(x,y)$
 - $F(x,y) = w \cdot \phi(x,y)$, so what we have to learn is w

Training data: $\{(x^1, \hat{y}^1), (x^2, \hat{y}^2), \dots, (x^r, \hat{y}^r), \dots\}$

We should find w such that

$\forall r$ (All training examples)

$\forall y \in Y - \{\hat{y}^r\}$ (All incorrect label
for r-th example)

$$w \cdot \phi(x^r, \hat{y}^r) > w \cdot \phi(x^r, y)$$

Structured Linear Model:

Problem 3



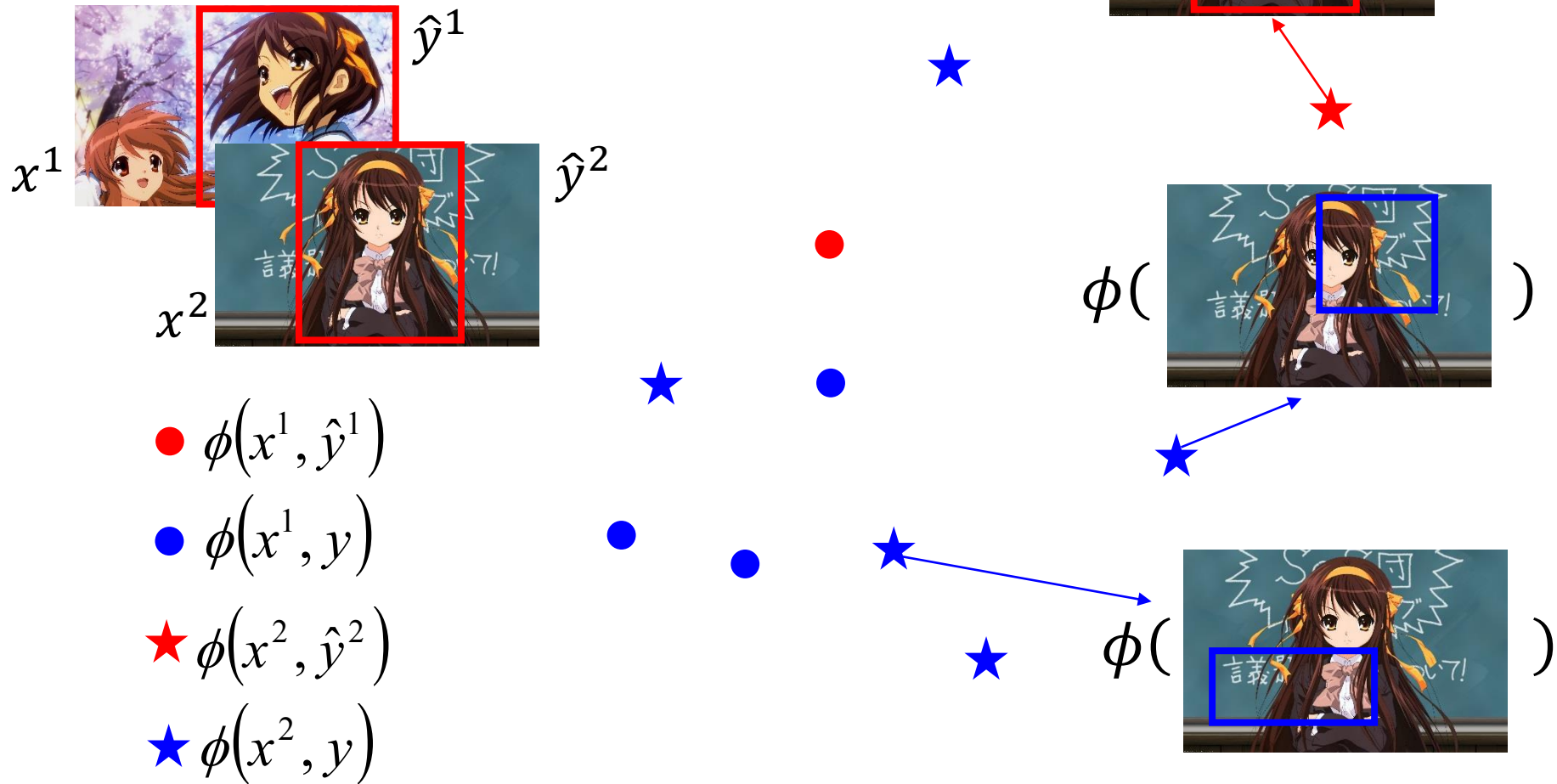
正確的 ● $\phi(x^1, \hat{y}^1)$

錯誤的 ● $\phi(x^1, y)$



Structured Linear Model:

Problem 3



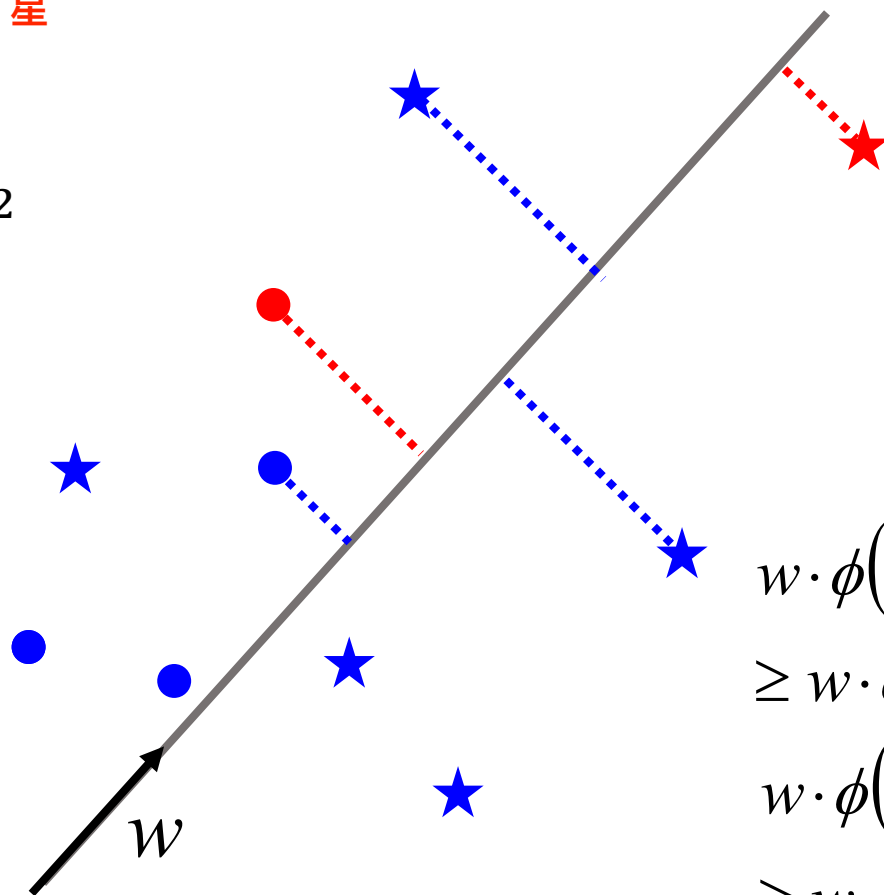
Structured Linear Model: Problem 3



- $\phi(x^1, \hat{y}^1)$
- $\phi(x^1, y)$
- ★ $\phi(x^2, \hat{y}^2)$
- ★ $\phi(x^2, y)$

紅色圈贏過所有藍色
圈色星贏過所有藍色
星

圈圈跟星星間彼此不比較



$$\begin{aligned}
 w \cdot \phi(x^1, \hat{y}^1) &\geq w \cdot \phi(x^1, y) \\
 w \cdot \phi(x^2, \hat{y}^2) &\geq w \cdot \phi(x^2, y)
 \end{aligned}$$

Solution of Problem 3

Difficult?

Not as difficult as expected

Algorithm

Will it terminate?

- **Input**: training data set $\{(x^1, \hat{y}^1), (x^2, \hat{y}^2), \dots, (x^r, \hat{y}^r), \dots\}$
- **Output**: weight vector w
- **Algorithm**: Initialize $w = 0$

- do

- For each pair of training example (x^r, \hat{y}^r)
 - Find the label \tilde{y}^r maximizing $w \cdot \phi(x^r, y)$

解這個optimization的problem $\tilde{y}^r = \arg \max_{y \in Y} w \cdot \phi(x^r, y)$ (question 2)

- If $\tilde{y}^r \neq \hat{y}^r$, update w 如果找出來不是正確答案，則做更新
如果找出來是正確的則這個式子會不增不減

$$w \rightarrow w + \phi(x^r, \hat{y}^r) - \phi(x^r, \tilde{y}^r)$$

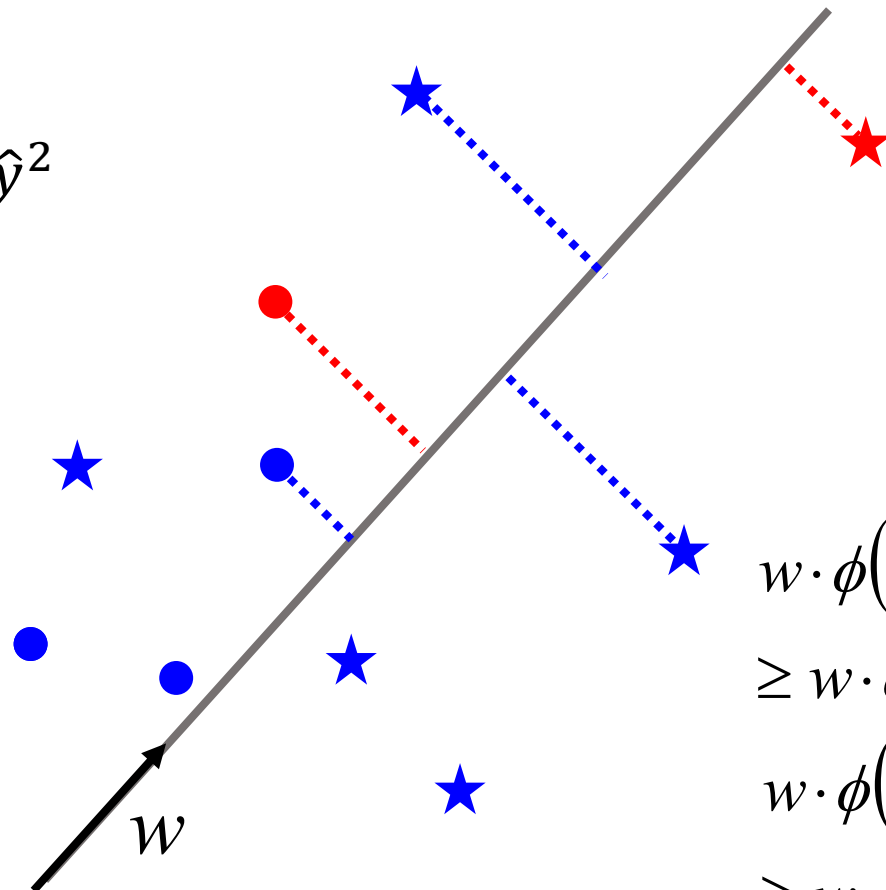
- until w is not updated  We are done!

對所有train example，每個pair都是正確的話則結束

Algorithm - Example



- $\phi(x^1, \hat{y}^1)$
- $\phi(x^1, y)$
- ★ $\phi(x^2, \hat{y}^2)$
- ★ $\phi(x^2, y)$



$$\begin{aligned}
 w \cdot \phi(x^1, \hat{y}^1) &\geq w \cdot \phi(x^1, y) \\
 w \cdot \phi(x^2, \hat{y}^2) &\geq w \cdot \phi(x^2, y)
 \end{aligned}$$

Algorithm - Example

Initialize $w = 0$

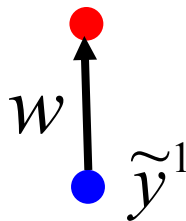
pick (x^1, \hat{y}^1) 隨機選一個
example pair

$$\tilde{y}^1 = \arg \max_{y \in Y} w \cdot \phi(x^1, y)$$

If $\tilde{y}^1 \neq \hat{y}^1$, update w

$$w \rightarrow w + \boxed{\phi(x^1, \hat{y}^1) - \phi(x^1, \tilde{y}^1)}$$

vector



- $\phi(x^1, \hat{y}^1)$
- $\phi(x^1, y)$
- ★ $\phi(x^2, \hat{y}^2)$
- ★ $\phi(x^2, y)$

Because $w=0$ at this time, $\phi(x^1, y)$ always 0
因為 w 初始為零，因此隨機選一個 y

➡ Random pick
one point as \tilde{y}^r

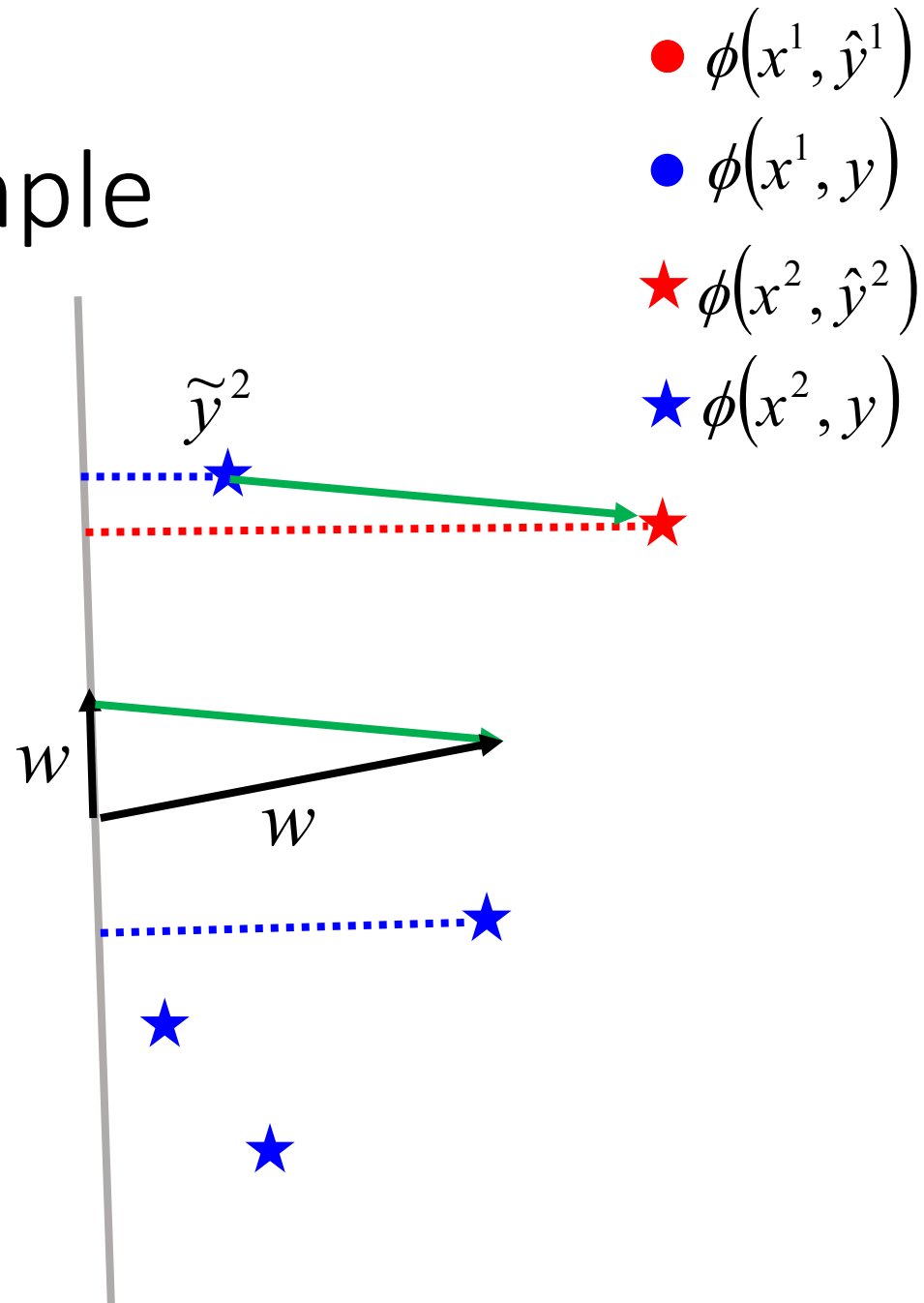
Algorithm - Example

pick (x^2, \hat{y}^2)

$$\tilde{y}^2 = \arg \max_{y \in Y} w \cdot \phi(x^2, y)$$

If $\tilde{y}^2 \neq \hat{y}^2$, update w

$$w \rightarrow w + \phi(x^2, \hat{y}^2) - \phi(x^2, \tilde{y}^2)$$



Algorithm - Example

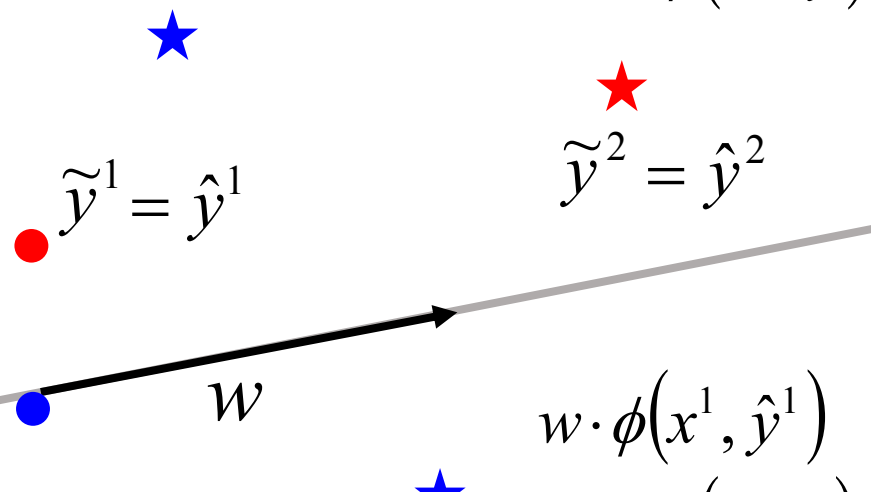
- $\phi(x^1, \hat{y}^1)$
- $\phi(x^1, y)$
- ★ $\phi(x^2, \hat{y}^2)$
- ★ $\phi(x^2, y)$

藍色的星星跟圈圈是有無限多個的，但是只要能解arg max的problem，則依然可以輕易解

pick (x^1, \hat{y}^1) again

$$\tilde{y}^1 = \arg \max_{y \in Y} w \cdot \phi(x^1, y)$$

$\tilde{y}^1 = \hat{y}^1 \Rightarrow$ do not update w



pick (x^2, \hat{y}^2) again

$$\tilde{y}^2 = \arg \max_{y \in Y} w \cdot \phi(x^2, y)$$

$\tilde{y}^2 = \hat{y}^2 \Rightarrow$ do not update w

$$\begin{aligned} w \cdot \phi(x^1, \hat{y}^1) &\geq w \cdot \phi(x^1, y) \\ w \cdot \phi(x^2, \hat{y}^2) &\geq w \cdot \phi(x^2, y) \end{aligned}$$

So we are done

證明會收斂


Assumption: Separable

- There exists a weight vector \hat{w} $\|\hat{w}\| = 1$

$\forall r$ (All training examples)

$\forall y \in Y - \{\hat{y}^r\}$ (All incorrect label for an example)

想要達成的目標


$$\hat{w} \cdot \phi(x^r, \hat{y}^r) \geq \hat{w} \cdot \phi(x^r, y) \quad (\text{The target exists})$$
$$\hat{w} \cdot \phi(x^r, \hat{y}^r) \geq \hat{w} \cdot \phi(x^r, y) + \delta$$

Assumption: Separable

$$\hat{w} \cdot \phi(x^r, \hat{y}^r) \geq \hat{w} \cdot \phi(x^r, y) + \delta$$

● $\phi(x^1, \hat{y}^1)$

● $\phi(x^1, y)$

★ $\phi(x^2, \hat{y}^2)$

★ $\phi(x^2, y)$

.....

