# Structured Linear Model

Hung-yi Lee

# Structured Linear Model

**Problem 1: Evaluation**

- What does F(x,y) look like? ➡️ in a specific form

**Problem 2: Inference**

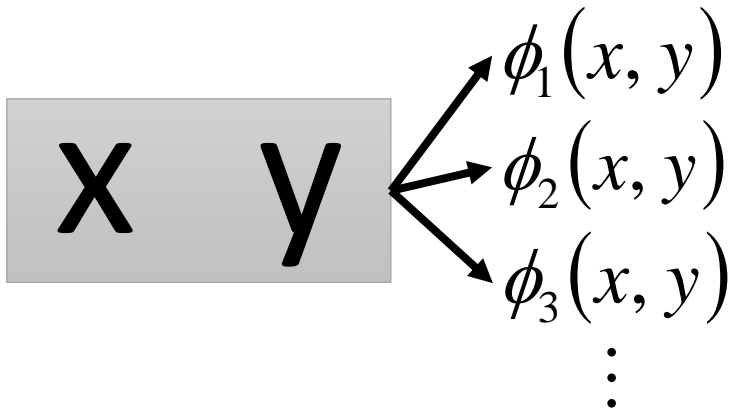- How to solve the "arg max" problem

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

**Problem 3: Training**

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

# Structured Linear Model: Problem 1

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

Characteristics

$$\phi_1(x, y)$$
$$\phi_2(x, y)$$
$$\phi_3(x, y)$$
$$\vdots$$

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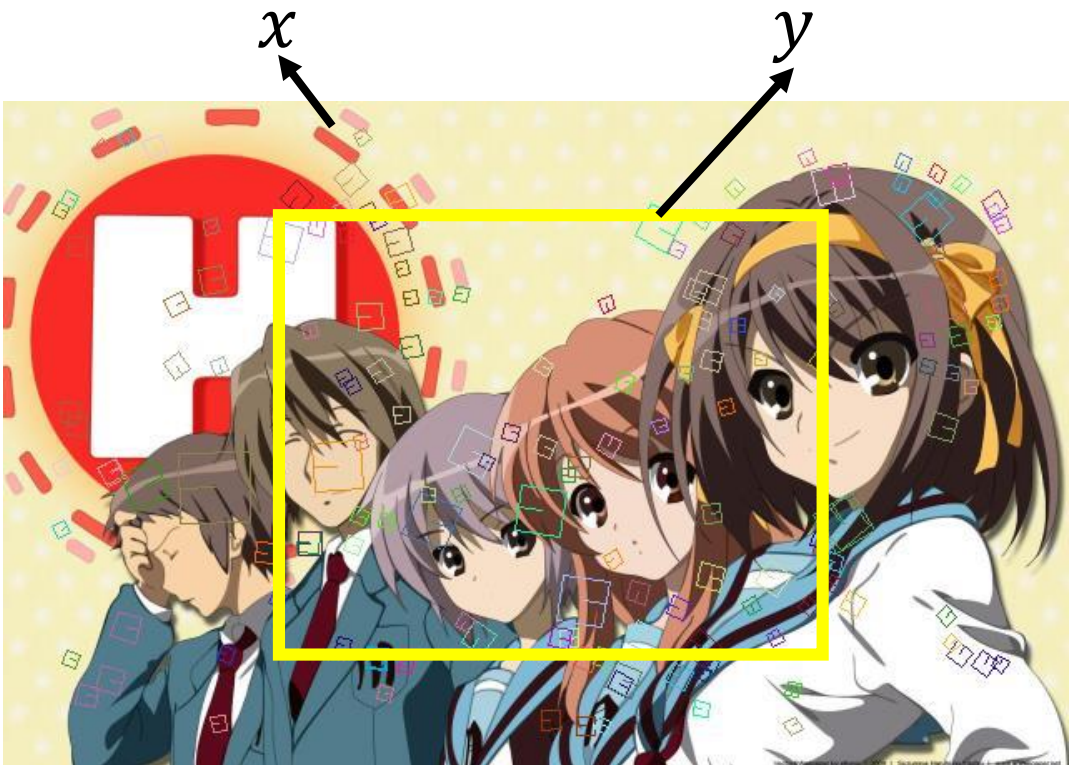
$$\text{F}(x, y) = w_1 \cdot \phi_1(x, y) + w_2 \cdot \phi_2(x, y) + w_3 \cdot \phi_3(x, y)\dots$$

Learning from data

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

# Structured Linear Model: Problem 1

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

- Example: ***Object Detection***



$x$

$y$

$\phi($  $) = \begin{bmatrix} \text{percentage of color red in box y} \\ \text{percentage of color green in box y} \\ \text{percentage of color blue in box y} \\ \text{percentage of color red out of box y} \\ \ldots\ldots \\ \text{area of box y} \\ \text{number of specific patterns in box y} \\ \ldots\ldots \end{bmatrix}$

$\phi(x, y)$



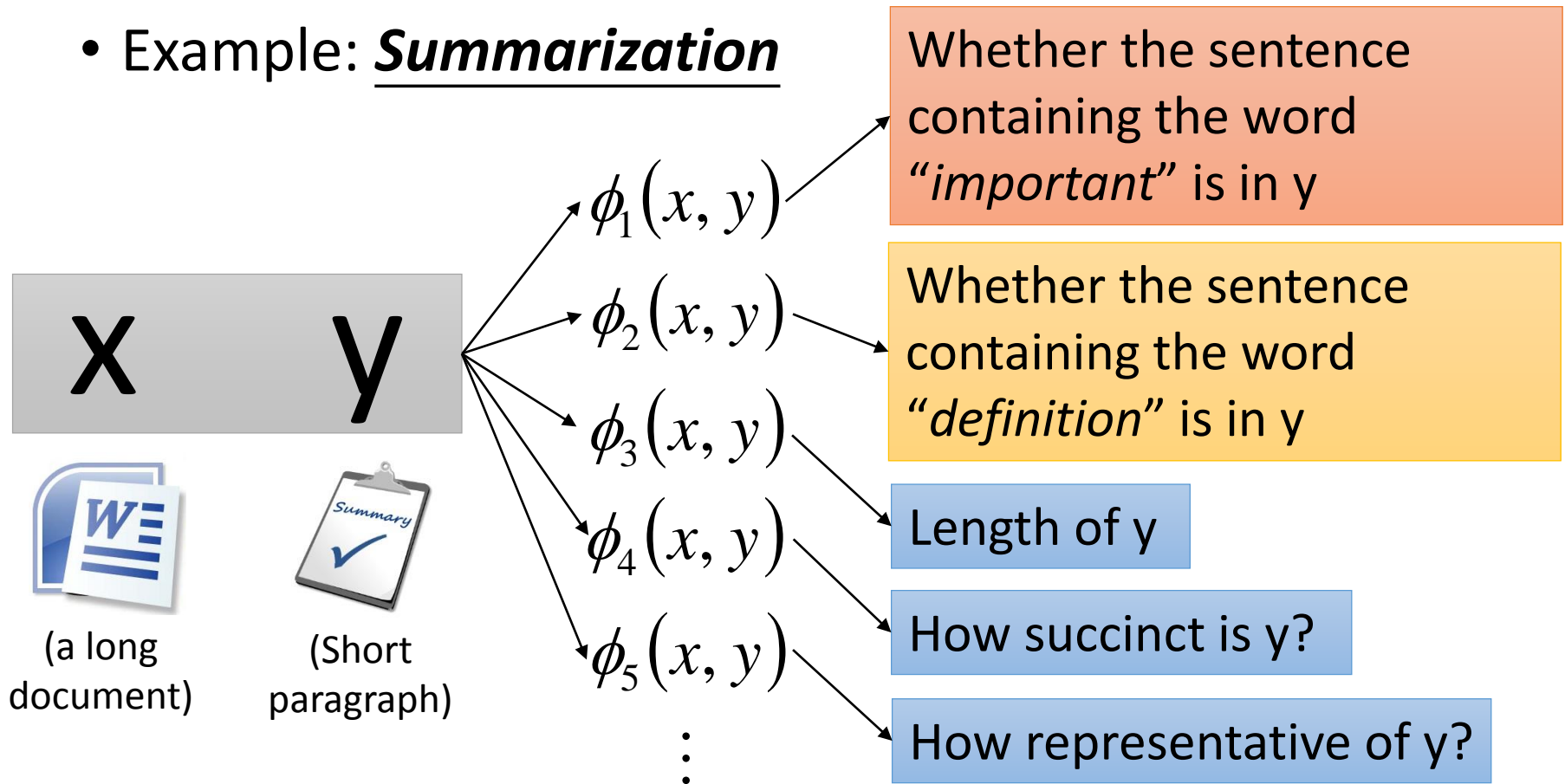| Convolutional Layer | Sub-sampling Layer | Fully-connected Layer | Output Layer |

$\phi($                                                    $)$

# Structured Linear Model: Problem 1

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



$\phi_1(x, y)$ — Whether the sentence containing the word "*important*" is in y

$\phi_2(x, y)$ — Whether the sentence containing the word "*definition*" is in y

$\phi_3(x, y)$

$\phi_4(x, y)$ — Length of y

$\phi_5(x, y)$ — How succinct is y?

⋮ — How representative of y?

X Y

(a long document)    (Short paragraph)

# Structured Linear Model: Problem 1

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

- Example: ***Retrieval***



$\phi_1(x, y)$ — The degree of relevance with respect to x for the top 1 webpages in y.

$\phi_2(x, y)$ — Is the top 1 webpage more relevant than the top 2 webpage?

$\phi_3(x, y)$

$\vdots$

How much different information does y cover? (***Diversity***)

X (Input keyword)

y (Search Result)

# Structured Linear Model: Problem 2

- **Inference**: How to solve the "arg max" problem

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

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

● Assume we have solved this question.

# Structured Linear Model: Problem 3

- Training: Given training data, how to learn F(x,y)
  - F(x,y) = w·ϕ(x,y), so what we have to learn is w

Training data: $\left\{ \left( x^1, \hat{y}^1 \right), \left( x^2, \hat{y}^2 \right), \ldots, \left( x^r, \hat{y}^r \right), \ldots \right\}$
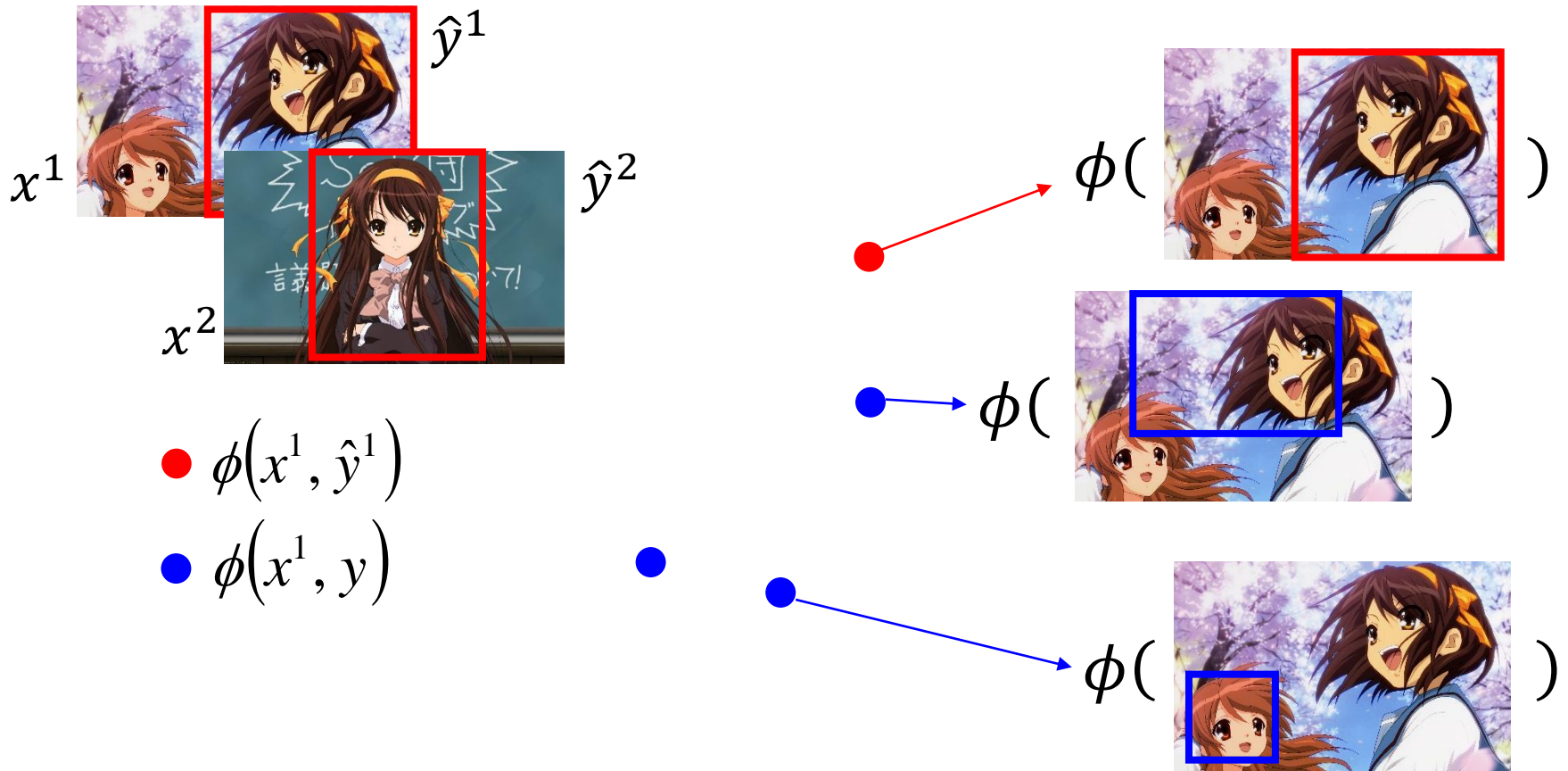
We should find w such that

$\forall r$ (All training examples)

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

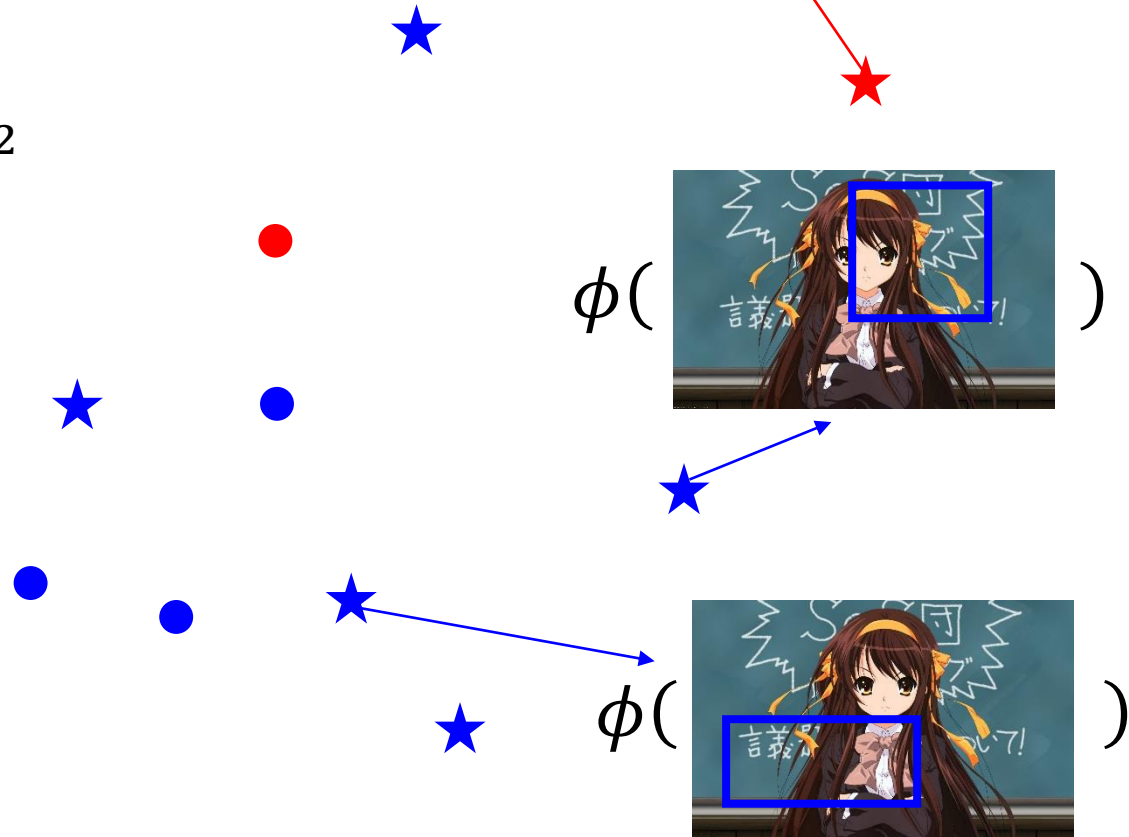$$w \cdot \phi\left( x^r, \hat{y}^r \right) > w \cdot \phi\left( x^r, y \right)$$

# Structured Linear Model:
## Problem 3

# Structured Linear Model:
## Problem 3



$x^1$

$\hat{y}^1$

$\hat{y}^2$

$x^2$

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

$\phi(\quad)$

$\phi(\quad)$

$\phi(\quad)$

# Structured Linear Model:
## Problem 3



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

$$\phi(x^1, y)$$

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

$$\phi(x^2, y)$$

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

# Solution of Problem 3

Difficult?

Not as difficult as expected

# Algorithm

- **Input**: training data set $\left\{ \left( x^1, \hat{y}^1 \right), \left( x^2, \hat{y}^2 \right), \ldots, \left( x^r, \hat{y}^r \right), \ldots \right\}$
- **Output**: weight vector w
- **Algorithm**: Initialize w = 0
  - do
    - For each pair of training example $\left( x^r, \hat{y}^r \right)$
      - Find the label $\tilde{y}^r$ maximizing $w \cdot \phi(x^r, y)$

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

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

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

  - until w is not updated ➡ We are done!

# Algorithm - Example

$\hat{y}^1$

$x^1$

$\hat{y}^2$

$x^2$

- 🔴 $\phi(x^1, \hat{y}^1)$
- 🔵 $\phi(x^1, y)$
- ⭐ $\phi(x^2, \hat{y}^2)$
- ⭐ $\phi(x^2, y)$

$w$

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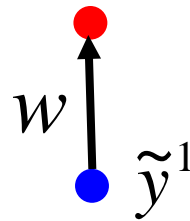$\geq w \cdot \phi(x^2, y)$

# Algorithm - Example

Initialize w = 0

pick $\left(x^1, \hat{y}^1\right)$

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

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

$$w \rightarrow w + \phi\left(x^1, \hat{y}^1\right) - \phi\left(x^1, \widetilde{y}^1\right)$$

$\bullet\ \phi\left(x^1, \hat{y}^1\right)$

$\bullet\ \phi\left(x^1, y\right)$

$\star\ \phi\left(x^2, \hat{y}^2\right)$

$\star\ \phi\left(x^2, y\right)$

$w$   $\widetilde{y}^1$

Because w=0 at this time, $\phi(x^1, y)$ always 0

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

# Algorithm - Example

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

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

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

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

$\bullet$ $\phi\left(x^1, \hat{y}^1\right)$

$\bullet$ $\phi\left(x^1, y\right)$

$\star$ $\phi\left(x^2, \hat{y}^2\right)$

$\star$ $\phi\left(x^2, y\right)$

$\widetilde{y}^2$

$w$

$w$

# Algorithm - Example

$\bullet \, \phi(x^1, \hat{y}^1)$

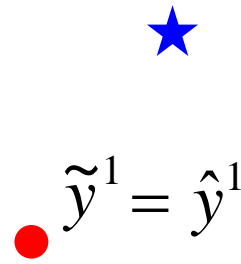$\bullet \, \phi(x^1, y)$

$\star \, \phi(x^2, \hat{y}^2)$

$\star \, \phi(x^2, y)$

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

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

$\widetilde{y}^1 = \hat{y}^1$ ➡ do not update w

$\widetilde{y}^1 = \hat{y}^1$

$\widetilde{y}^2 = \hat{y}^2$

$w$

$w \cdot \phi(x^1, \hat{y}^1)$

$\geq w \cdot \phi(x^1, y)$

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

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

$w \cdot \phi(x^2, \hat{y}^2)$

$\geq w \cdot \phi(x^2, y)$

$\widetilde{y}^2 = \hat{y}^2$ ➡ do not update w

So we are done

# 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)

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

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

# Assumption: Separable

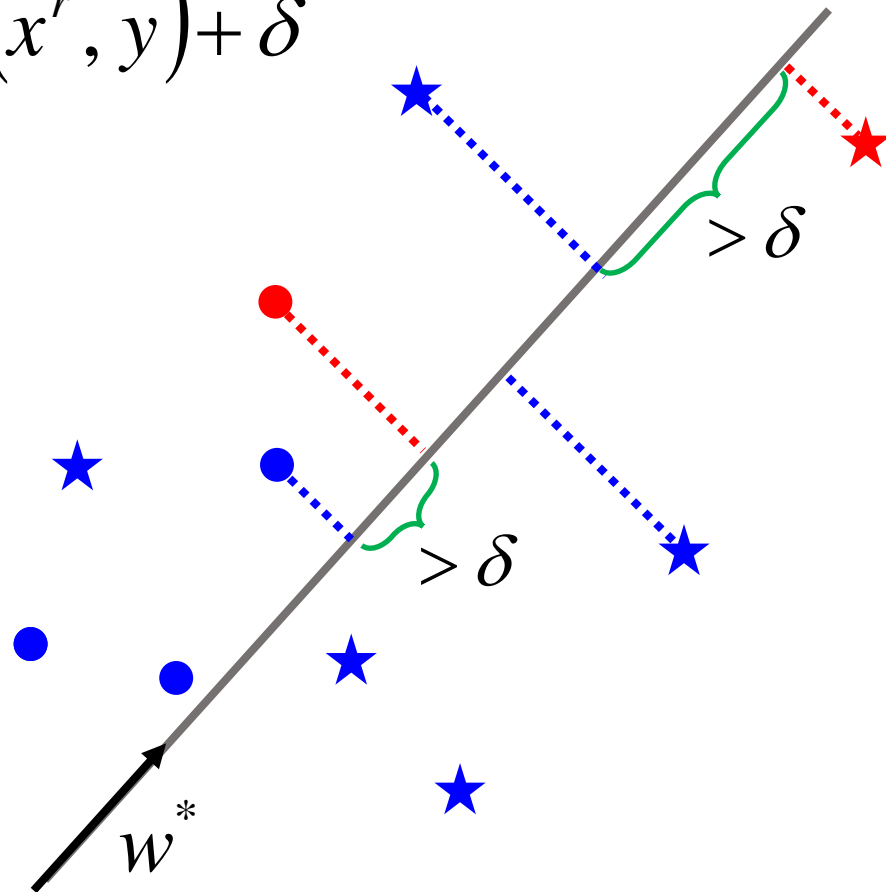$$\hat{w} \cdot \phi\left(x^r, \hat{y}^r\right) \geq \hat{w} \cdot \phi\left(x^r, y\right) + \delta$$

- 🔴 $\phi\left(x^1, \hat{y}^1\right)$
- 🔵 $\phi\left(x^1, y\right)$
- ⭐ $\phi\left(x^2, \hat{y}^2\right)$
- ⭐ $\phi\left(x^2, y\right)$
- ......

$> \delta$

$> \delta$

$w^*$

# Proof of Termination

w is updated <span style="color:red">once it sees a mistake</span>

$$w^0 = 0 \rightarrow w^1 \rightarrow w^2 \rightarrow \ldots\ldots \rightarrow w^k \rightarrow w^{k+1} \rightarrow \ldots\ldots$$

$$w^k = w^{k-1} + \phi\left(x^n, \hat{y}^n\right) - \phi\left(x^n, \tilde{y}^n\right)$$ <span style="color:blue">(the relation of $w^k$ and $w^{k-1}$)</span>

Proof that: The angle $\rho_k$ between $\hat{w}$ and $w_k$ is smaller as k increases

Analysis $\cos\rho_k$ <span style="color:blue">(larger and larger?)</span> $\quad \cos\rho_k = \dfrac{\hat{w}}{\|\hat{w}\|} \cdot \dfrac{w^k}{\|w^k\|}$

$$\hat{w} \cdot w^k = \hat{w} \cdot \left(w^{k-1} + \phi\left(x^n, \hat{y}^n\right) - \phi\left(x^n, \tilde{y}^n\right)\right)$$

$$= \hat{w} \cdot w^{k-1} + \underline{\hat{w} \cdot \phi\left(x^n, \hat{y}^n\right) - \hat{w} \cdot \phi\left(x^n, \tilde{y}^n\right)} \geq \hat{w} \cdot w^{k-1} + \delta$$

$$\geq \delta \;\; \text{(Separable)}$$

# Proof of Termination

w is updated <span style="color:red">once it sees a mistake</span>

$$w^0 = 0 \rightarrow w^1 \rightarrow w^2 \rightarrow \ldots\ldots \rightarrow w^k \rightarrow w^{k+1} \rightarrow \ldots\ldots$$

$$w^k = w^{k-1} + \phi\left(x^n, \hat{y}^n\right) - \phi\left(x^n, \tilde{y}^n\right) \text{ (the relation of } w^k \text{ and } w^{k-1}\text{)}$$

Proof that: The angle $\rho_k$ between $\hat{w}$ and $w_k$ is smaller as k increases

Analysis $\cos \rho_k$ (larger and larger?) $\quad \cos \rho_k = \dfrac{\hat{w} \quad w^k}{\|\hat{w}\| \cdot \|w^k\|}$

$$\hat{w} \cdot w^k \geq \hat{w} \cdot w^{k-1} + \delta$$

<span style="color:red">=0</span> $\qquad\qquad\qquad$ <span style="color:red">≥δ</span>

$$\hat{w} \cdot w^1 \geq \hat{w} \cdot w^0 + \delta \qquad \hat{w} \cdot w^2 \geq \hat{w} \cdot w^1 + \delta \quad \cdots\cdots$$

$$\hat{w} \cdot w^1 \geq \delta \qquad\qquad \hat{w} \cdot w^2 \geq 2\delta \qquad \cdots\cdots$$

$$\hat{w} \cdot w^k \geq k\delta$$

(so what)

# Proof of Termination

$$\cos \rho_k = \frac{\hat{w}}{\|\hat{w}\|} \cdot \frac{w^k}{\boxed{\|w^k\|}} \qquad w^k = w^{k-1} + \phi(x^n, \hat{y}^n) - \phi(x^n, \tilde{y}^n)$$

$$\left\|w^k\right\|^2 = \left\|w^{k-1} + \phi(x^n, \hat{y}^n) - \phi(x^n, \tilde{y}^n)\right\|^2$$

$$= \left\|w^{k-1}\right\|^2 + \underbrace{\left\|\phi(x^n, \hat{y}^n) - \phi(x^n, \tilde{y}^n)\right\|^2}_{> 0} + \underbrace{2w^{k-1} \cdot \left(\phi(x^n, \hat{y}^n) - \phi(x^n, \tilde{y}^n)\right)}_{? \; < 0 \; \text{(mistake)}}$$

> Assume the distance between any two feature vector is smaller than R

$$\left\|w^1\right\|^2 \leq \left\|w^0\right\|^2 + \mathbf{R}^2 = \mathbf{R}^2$$

$$\left\|w^2\right\|^2 \leq \left\|w^1\right\|^2 + \mathbf{R}^2 \leq 2\mathbf{R}^2$$

$$\cdots$$

$$\leq \left\|w^{k-1}\right\| + \mathbf{R}^2$$
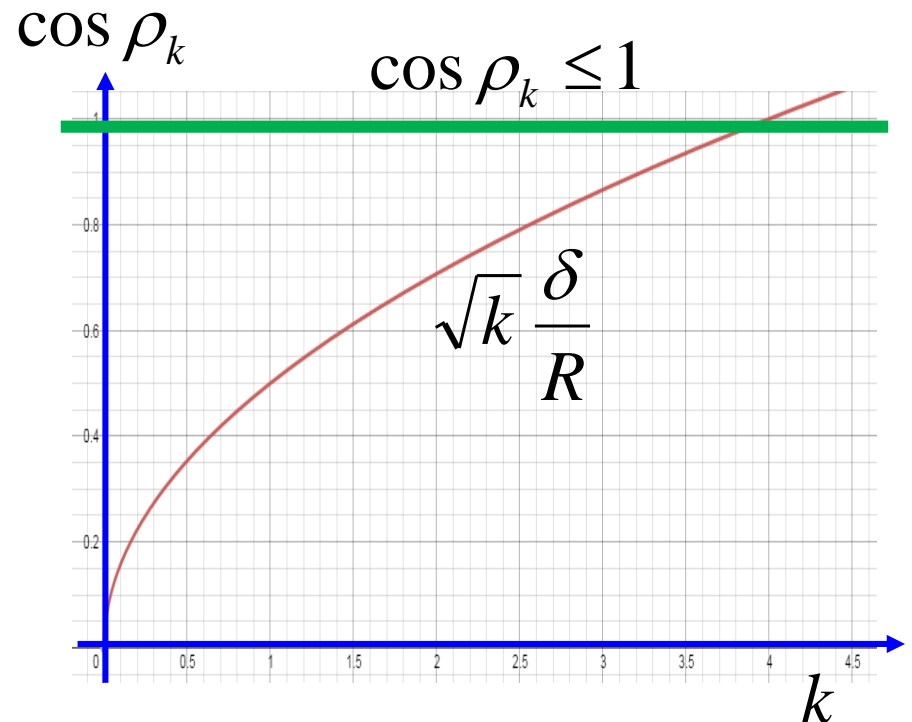
$$\left\|w^k\right\|^2 \leq k\mathbf{R}^2$$

# Proof of Termination

$$\cos \rho_k = \frac{\hat{w}}{\|\hat{w}\|} \cdot \frac{w^k}{\|w^k\|} \qquad \hat{w} \cdot w^k \geq k\delta \qquad \|w^k\|^2 \leq kR^2$$

$$\geq \frac{k\delta}{\sqrt{kR^2}} = \sqrt{k}\,\frac{\delta}{R}$$

$$\sqrt{k}\,\frac{\delta}{R} \leq 1$$

$$k \leq \left(\frac{R}{\delta}\right)^2$$



$\cos \rho_k$

$\cos \rho_k \leq 1$

$\sqrt{k}\,\frac{\delta}{R}$

$k$

# Proof of Termination

$$k \leq \left( \frac{R}{\delta} \right)^2$$

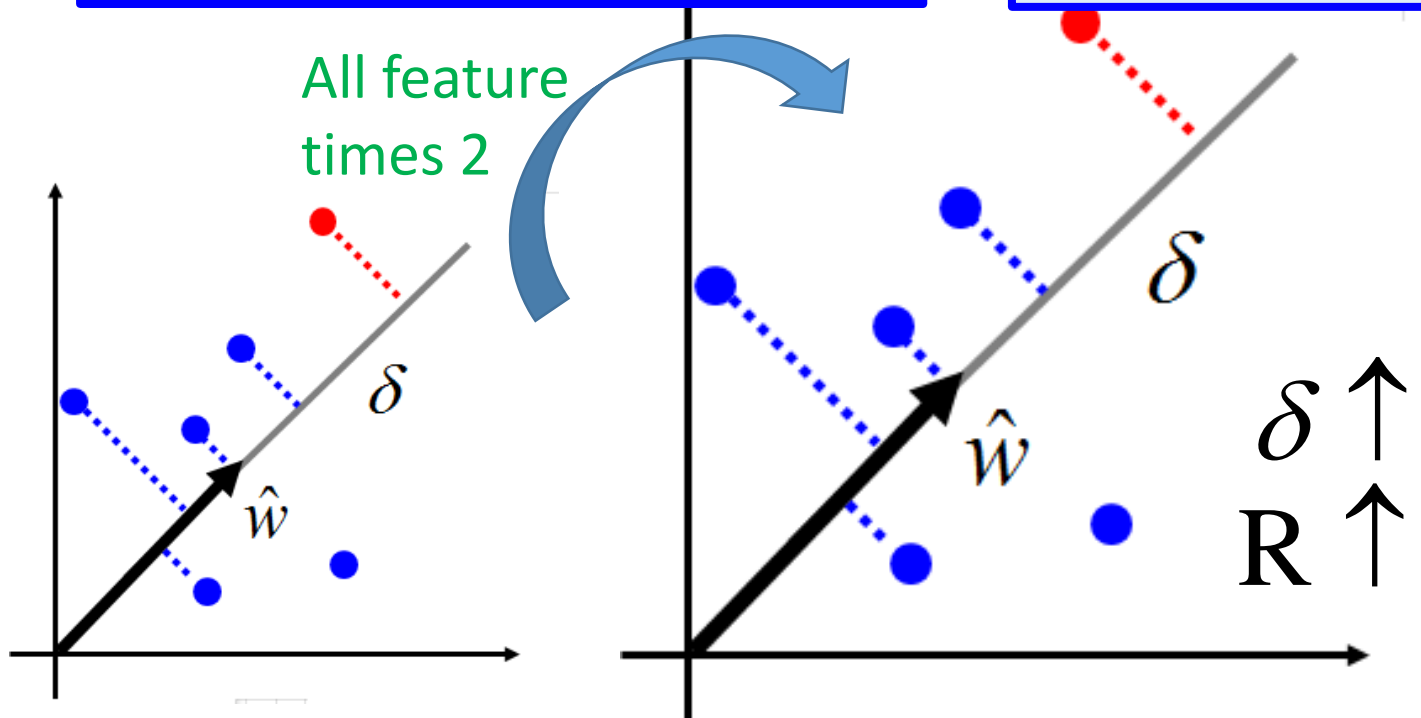The largest distances between features → Normalization

Margin: Is it easy to separable red points from the blue ones → Larger margin, less update

All feature times 2

- 🔴 $\phi(x^r, \hat{y}^r)$
- 🔵 $\phi(x^r, y)$

# Structured Linear Model:
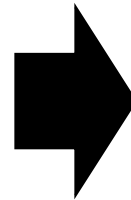## Reduce 3 Problems to 2

**Problem 1: Evaluation**

- How to define $F(x,y)$

**Problem 2: Inference**

- How to find the y with the largest $F(x,y)$

**Problem 3: Training**

- How to learn $F(x,y)$

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

**Problem A: Feature**

- How to define $\phi(x,y)$

**Problem B: Inference**

- How to find the y with the largest $w\cdot\phi(x,y)$