


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

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

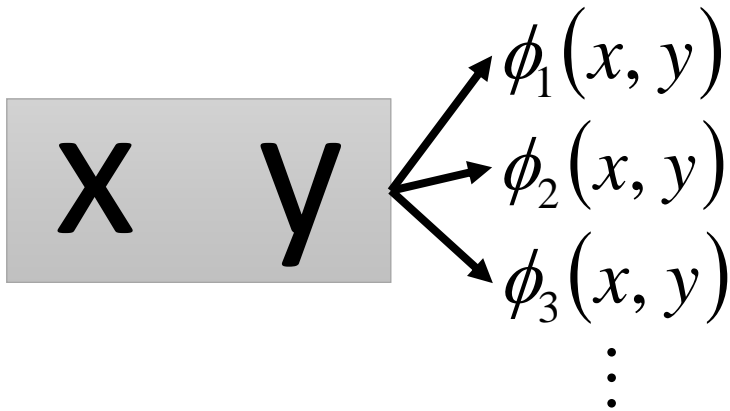
Problem 3: Training

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

Structured Linear Model: Problem 1

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

Characteristics



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

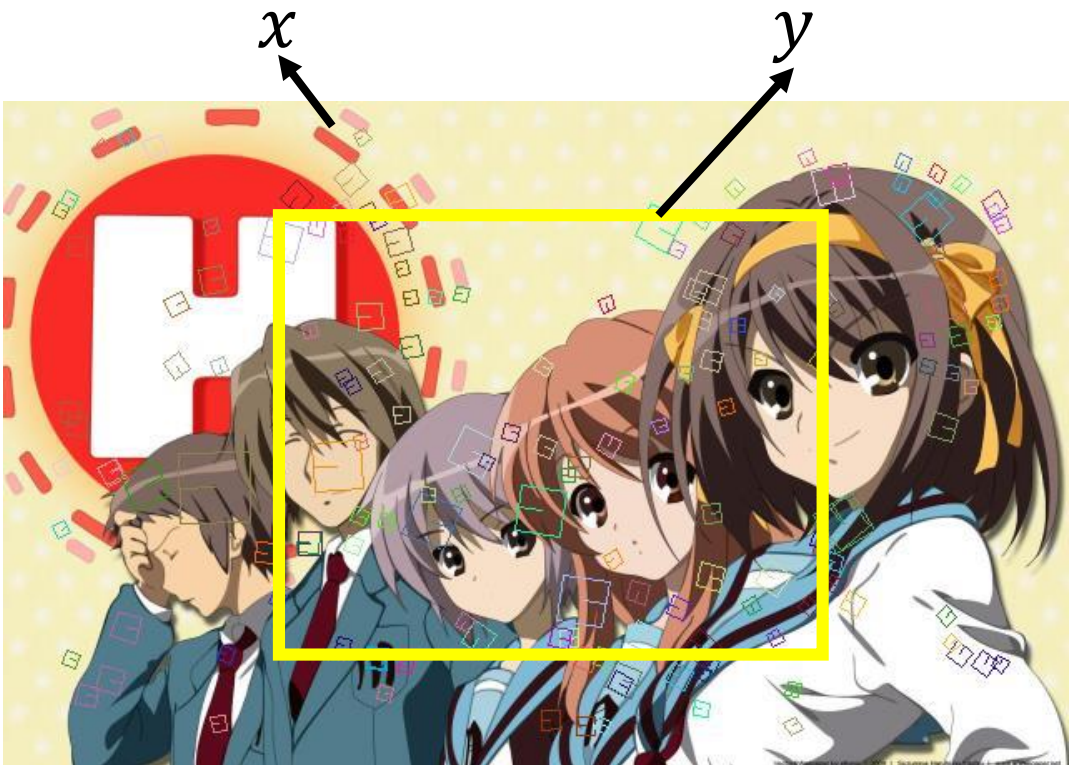
Learning
from data

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

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

Structured Linear Model: Problem 1

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

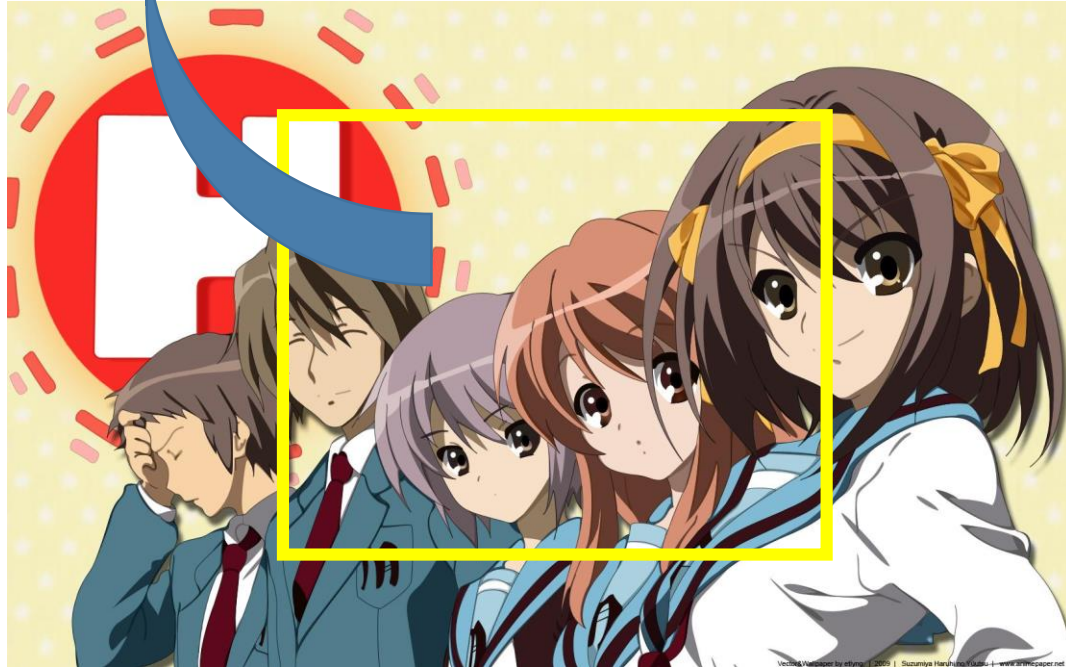
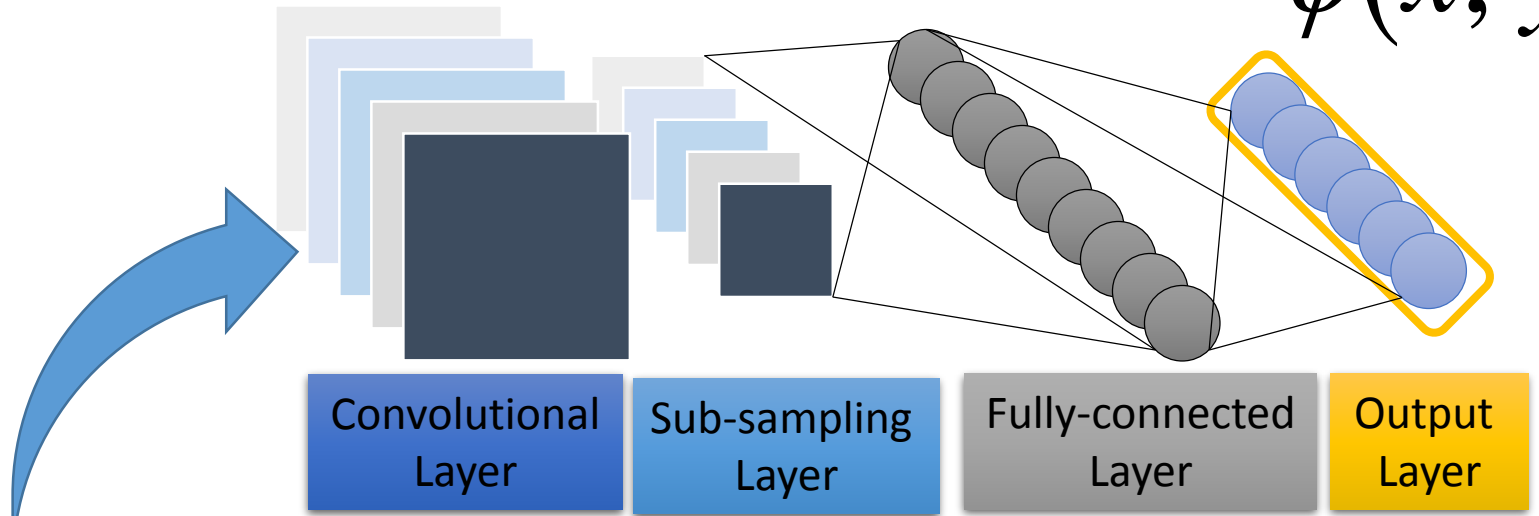


$\phi($

) =

- 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
-

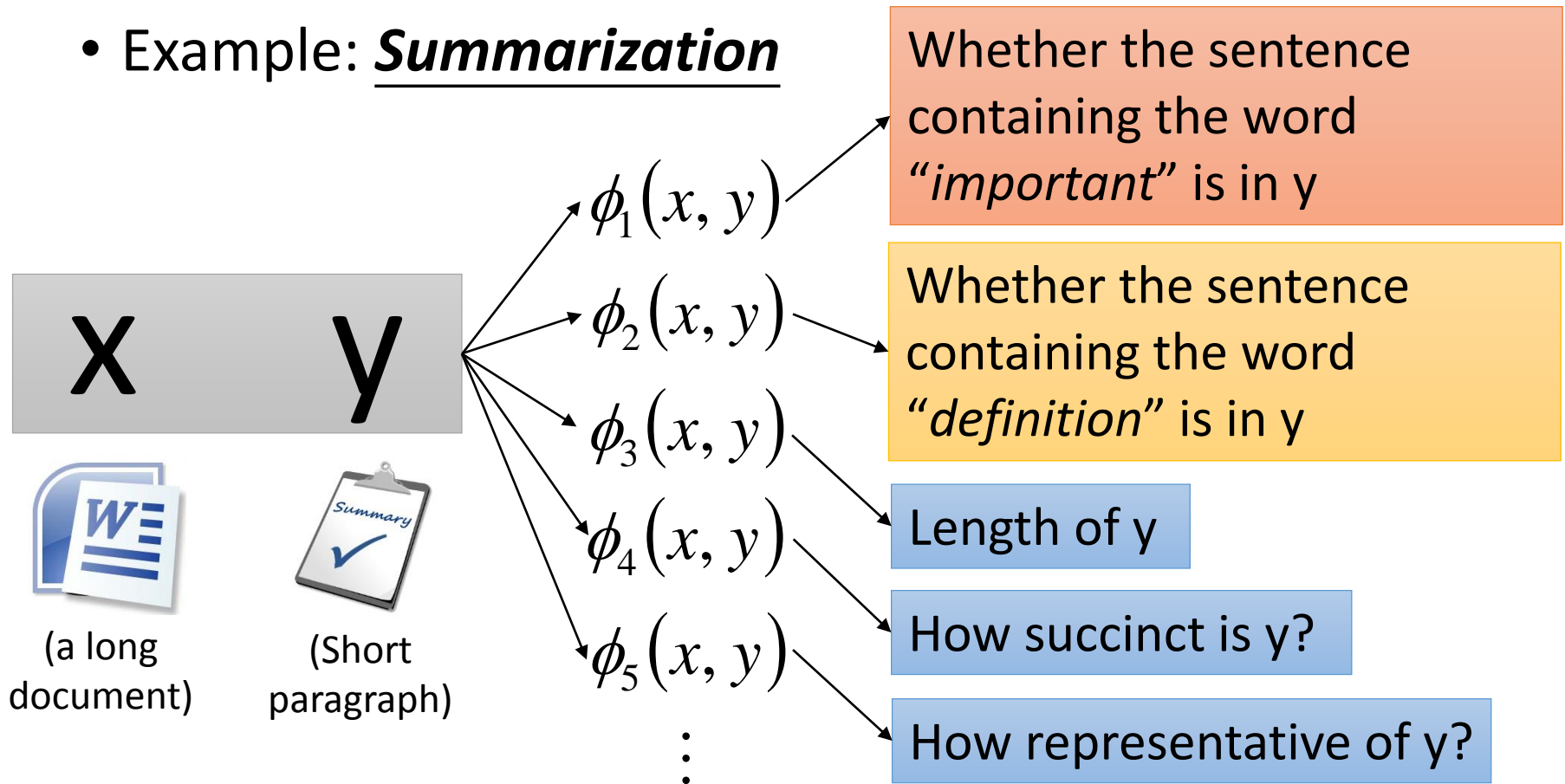
$$\phi(x, y)$$



$\phi($)

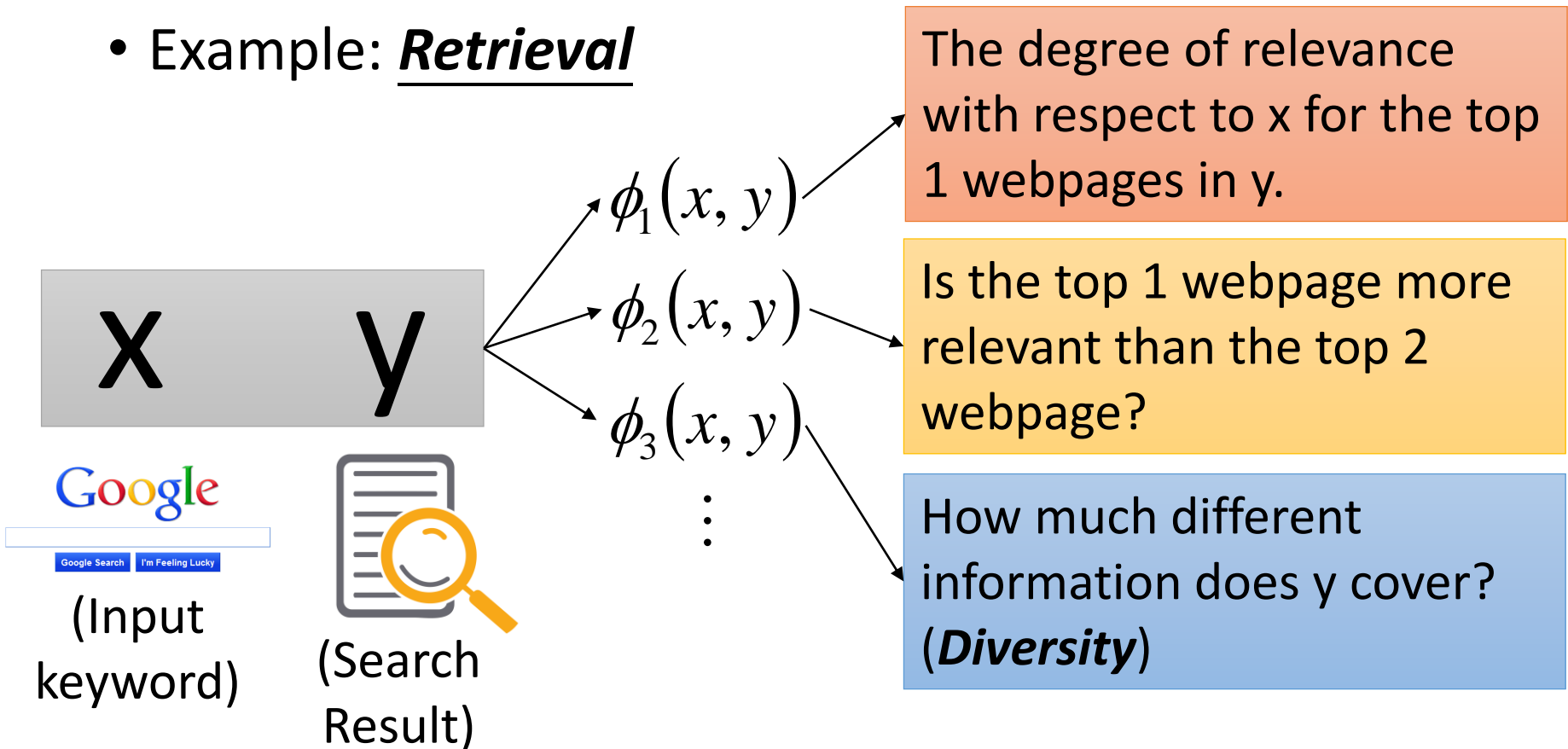
Structured Linear Model: Problem 1

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



Structured Linear Model: Problem 1

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



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.

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

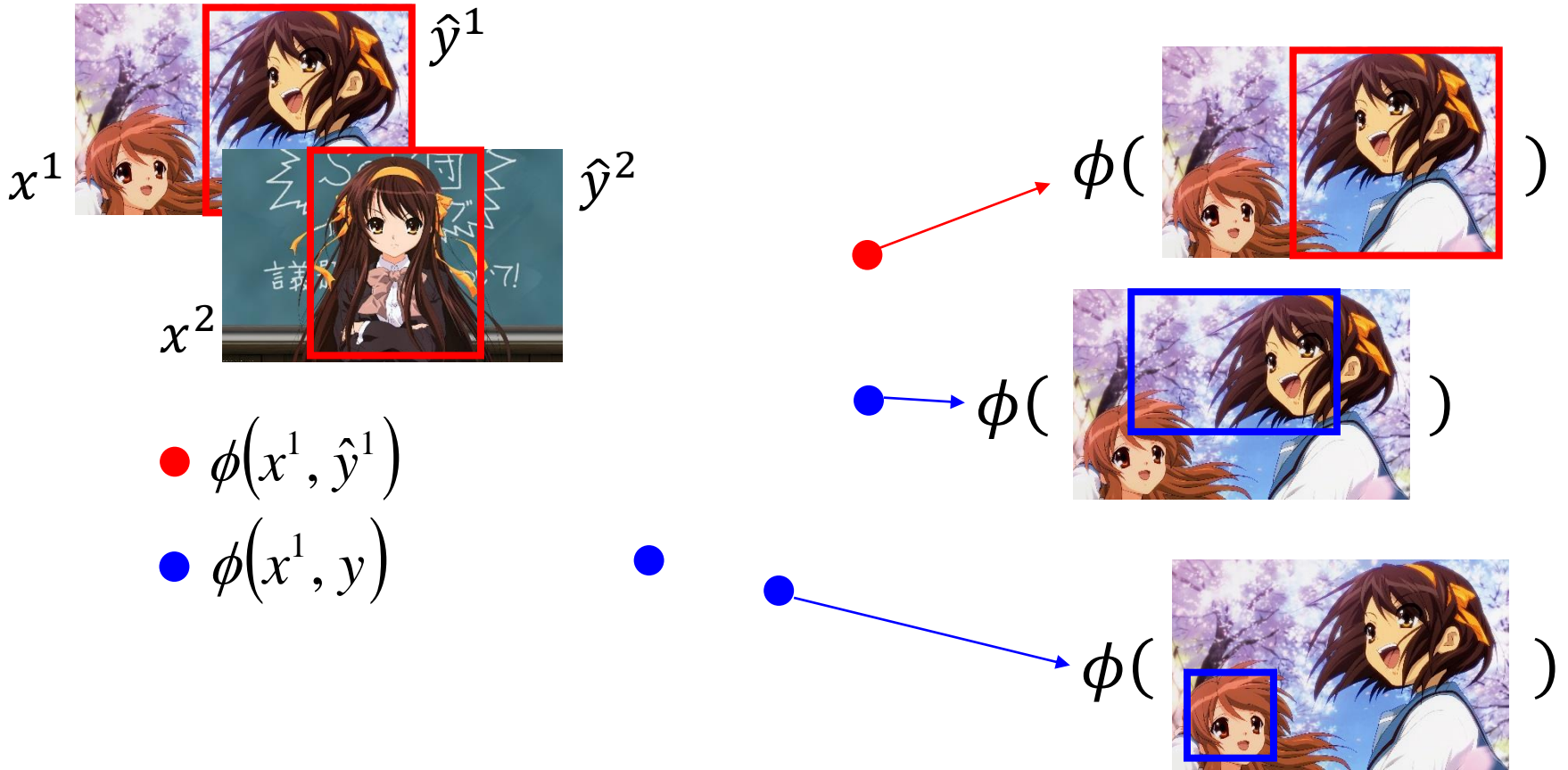
We should find w such that

$\forall r$ (All training examples)

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

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

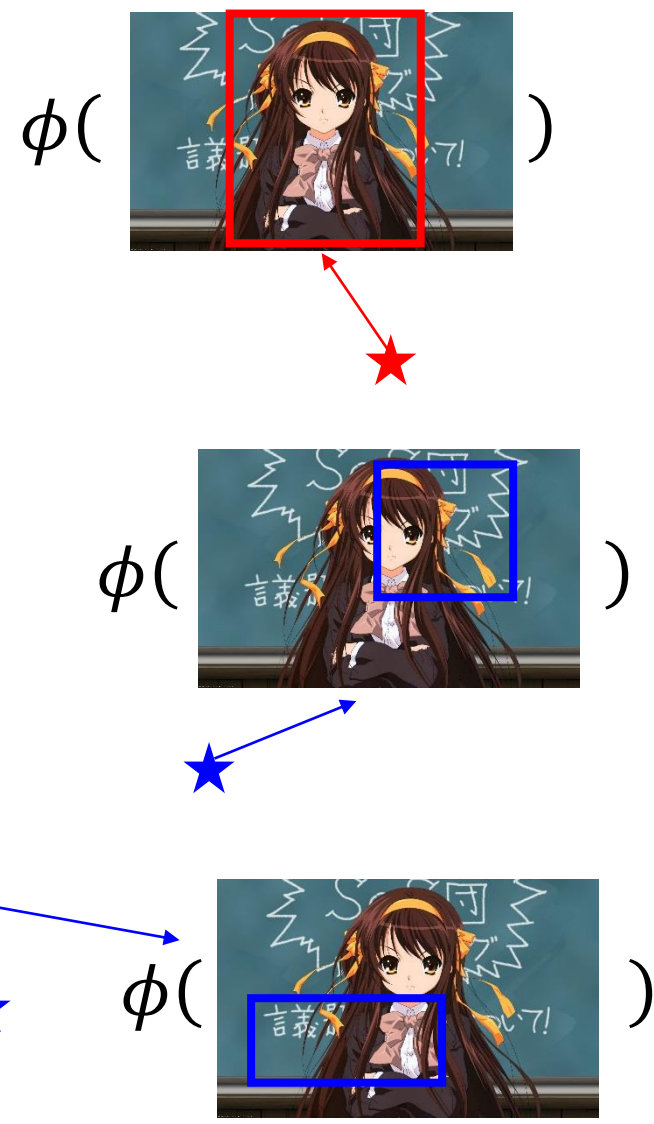
Structured Linear Model: Problem 3



Structured Linear Model: Problem 3



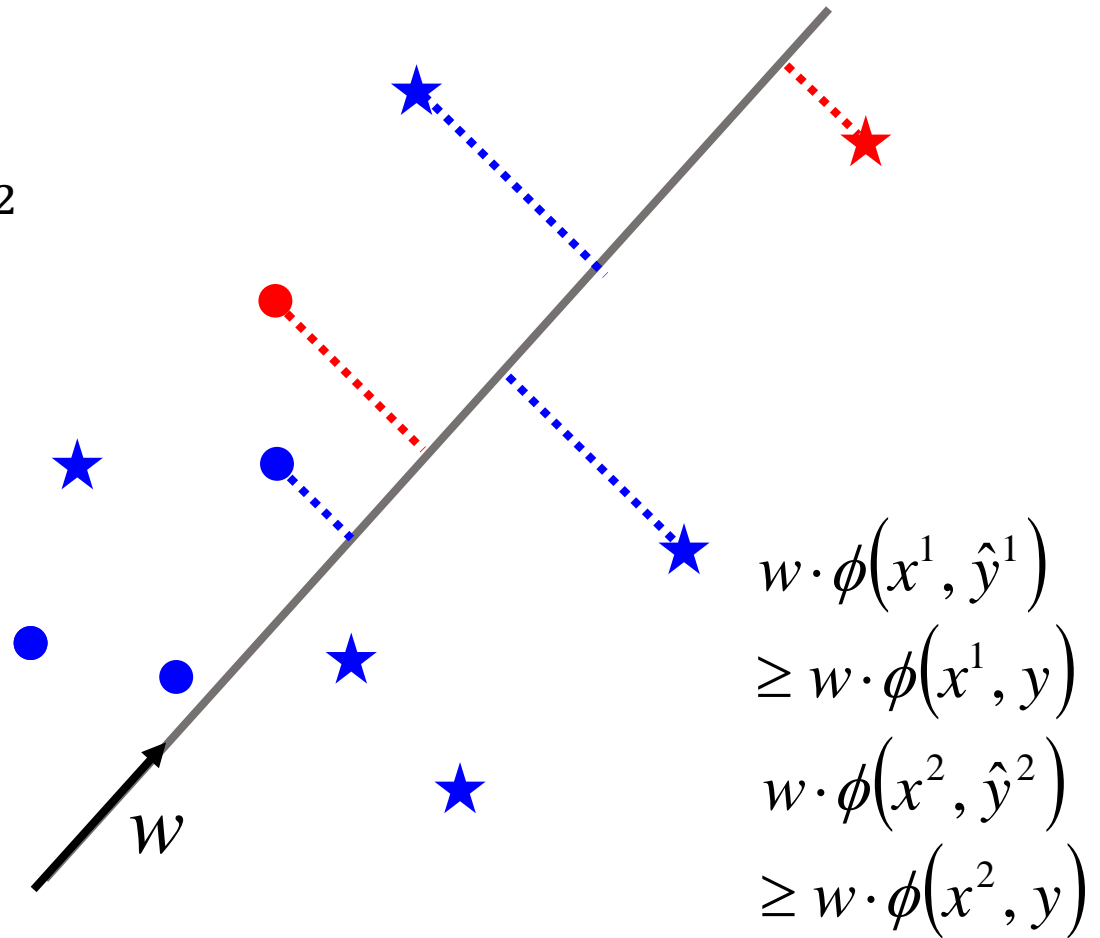
- $\phi(x^1, \hat{y}^1)$
- $\phi(x^1, y)$
- ★ $\phi(x^2, \hat{y}^2)$
- ★ $\phi(x^2, y)$



Structured Linear Model: Problem 3



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




Solution of Problem 3

Difficult?

Not as difficult as expected

Algorithm

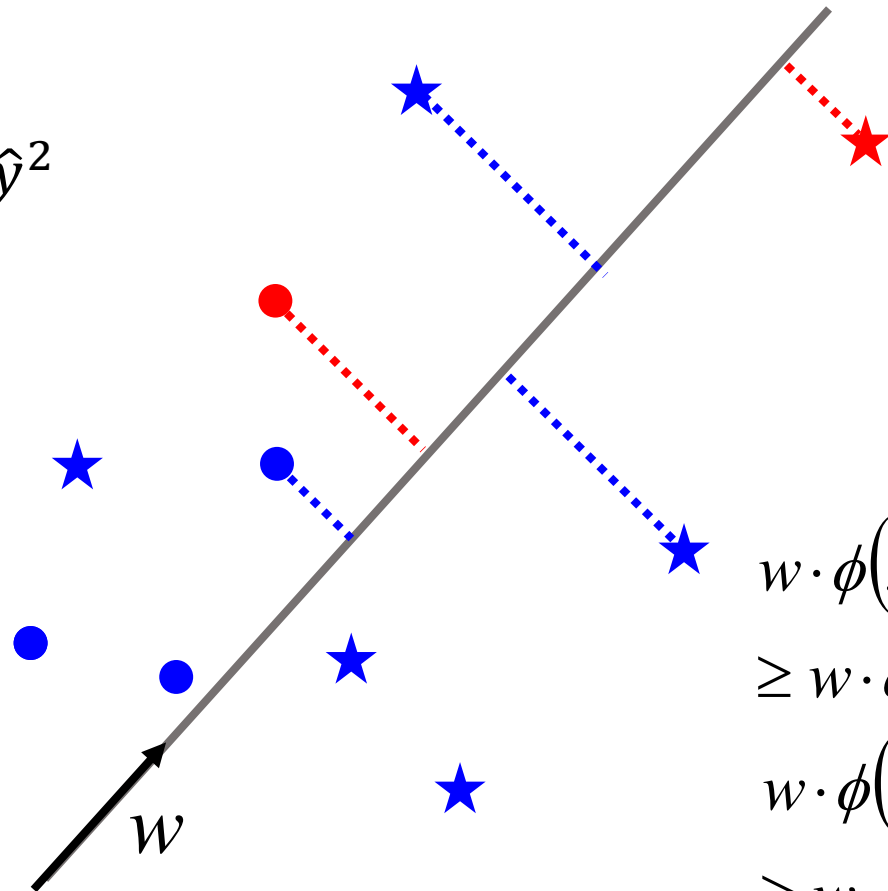
Will it terminate?

- **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 $w \cdot \phi(x^r, y)$
$$\tilde{y}^r = \arg \max_{y \in Y} w \cdot \phi(x^r, y) \text{ (question 2)}$$
 - If $\tilde{y}^r \neq \hat{y}^r$, update w
$$w \rightarrow w + \phi(x^r, \hat{y}^r) - \phi(x^r, \tilde{y}^r)$$
 - until w is not updated  We are done!

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

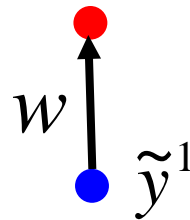
Initialize $w = 0$

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

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

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

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



● $\phi(x^1, \hat{y}^1)$

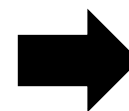
● $\phi(x^1, y)$

★ $\phi(x^2, \hat{y}^2)$

★ $\phi(x^2, y)$



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



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

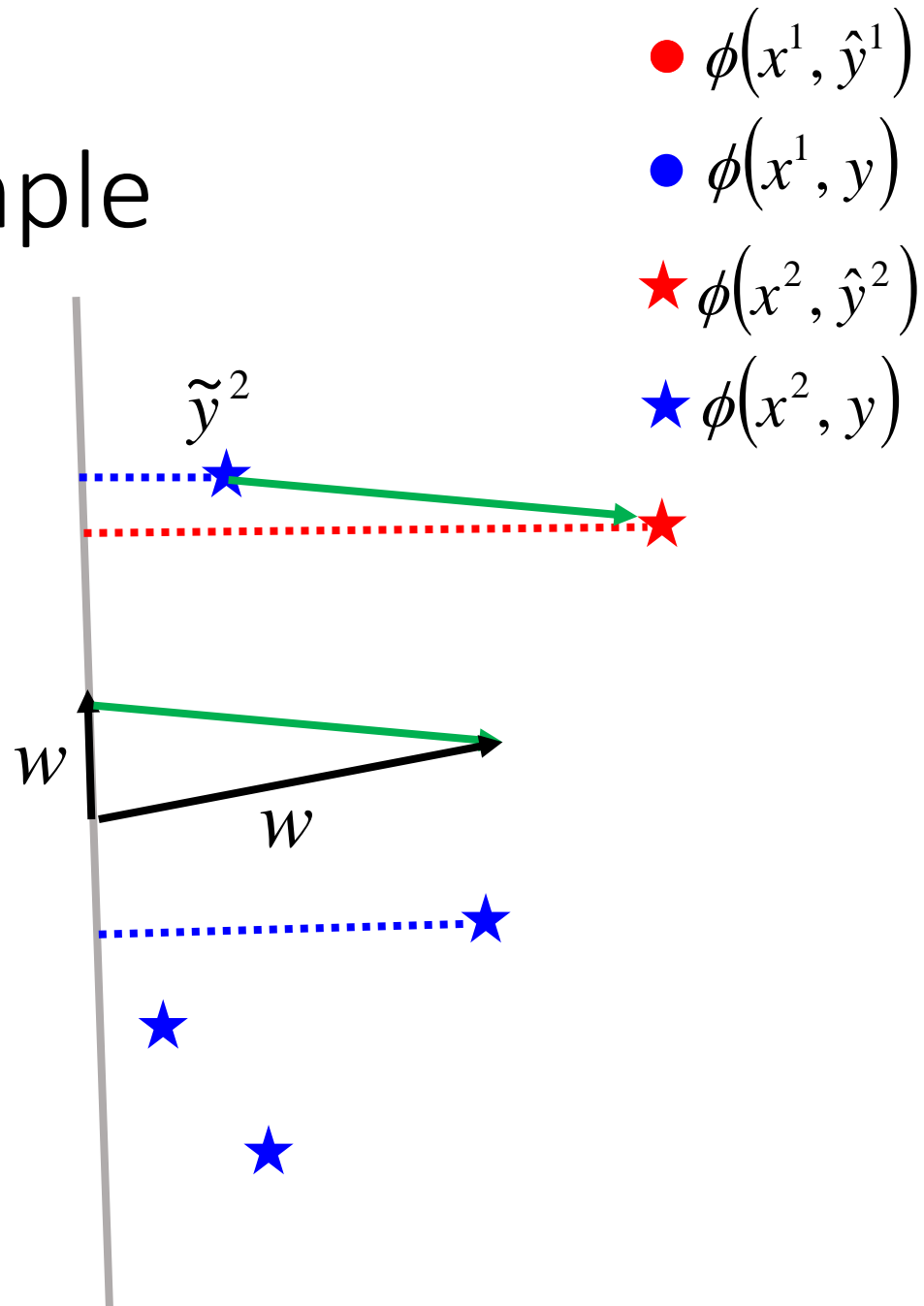
Algorithm - Example

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

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



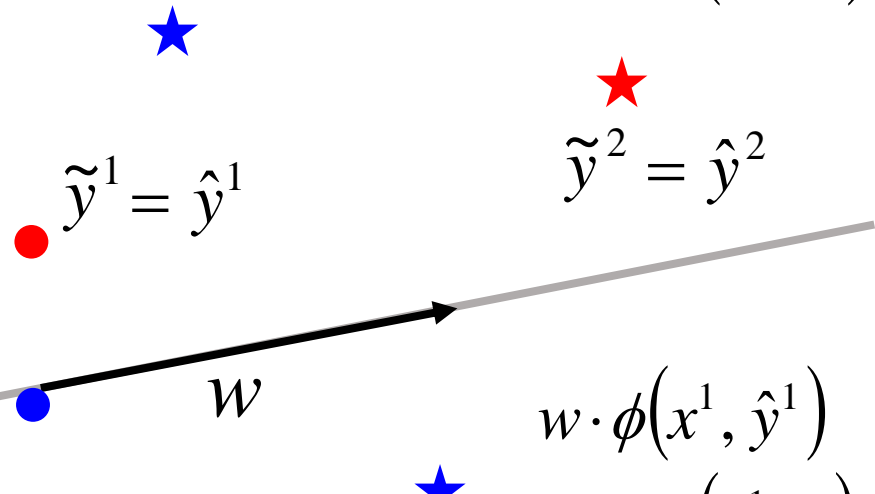
Algorithm - Example

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

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

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

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



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

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

$\tilde{y}^2 = \hat{y}^2$ ➡ 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

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)


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

Assumption: Separable

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

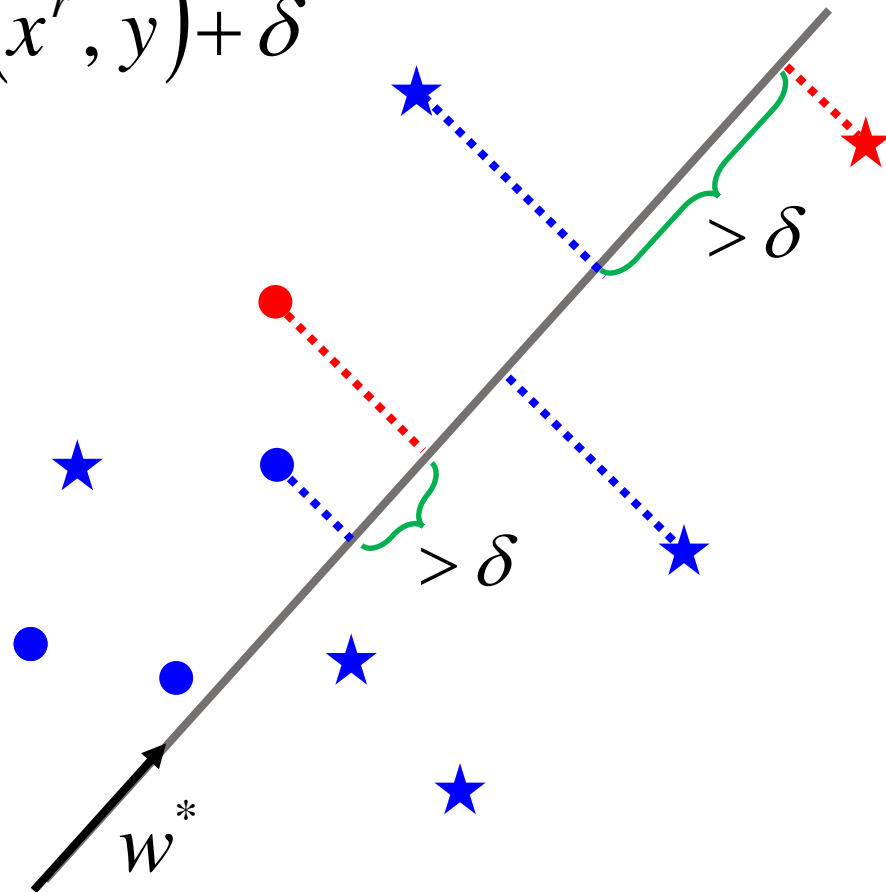
● $\phi(x^1, \hat{y}^1)$

● $\phi(x^1, y)$

★ $\phi(x^2, \hat{y}^2)$

★ $\phi(x^2, y)$

.....



Proof of Termination

w is updated **once it sees a mistake**

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

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

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

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

$$\begin{aligned} \hat{w} \cdot w^k &= \hat{w} \cdot (w^{k-1} + \phi(x^n, \hat{y}^n) - \phi(x^n, \tilde{y}^n)) \\ &= \hat{w} \cdot w^{k-1} + \underbrace{\hat{w} \cdot \phi(x^n, \hat{y}^n) - \hat{w} \cdot \phi(x^n, \tilde{y}^n)}_{\geq \delta \text{ (Separable)}} \geq \hat{w} \cdot w^{k-1} + \delta \end{aligned}$$

Proof of Termination

w is updated **once it sees a mistake**

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

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

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

Analysis $\cos \rho_k$ (larger and larger?)

$$\cos \rho_k = \frac{\hat{w} \cdot w^k}{\|\hat{w}\| \cdot \|w^k\|}$$

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

$$\hat{w} \cdot w^1 \geq \hat{w} \cdot w^0 + \delta$$

$$\hat{w} \cdot w^2 \geq \hat{w} \cdot w^1 + \delta \quad \dots$$

$$\hat{w} \cdot w^1 \geq \delta$$

$$\hat{w} \cdot w^2 \geq 2\delta$$

$$\dots$$

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

(so what)

Proof of Termination

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

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

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

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

$$\begin{aligned} \|w^1\|^2 &\leq \|w^0\|^2 + R^2 = R^2 \\ \|w^2\|^2 &\leq \|w^1\|^2 + R^2 \leq 2R^2 \\ &\dots \\ \|w^k\|^2 &\leq kR^2 \end{aligned}$$

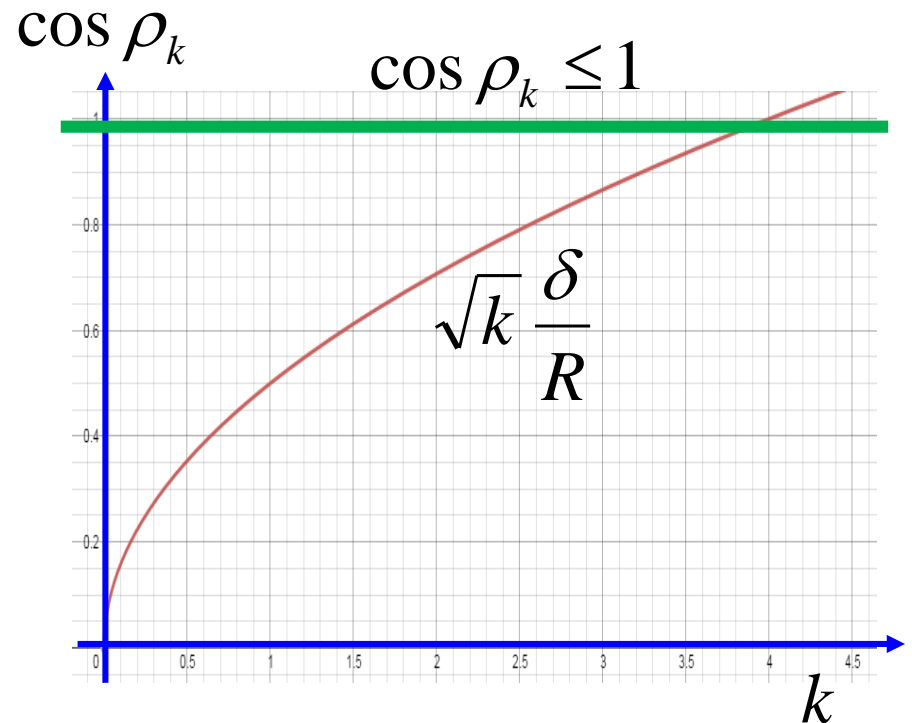
Proof of Termination

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

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

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

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



Proof of Termination

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

The largest distances between features

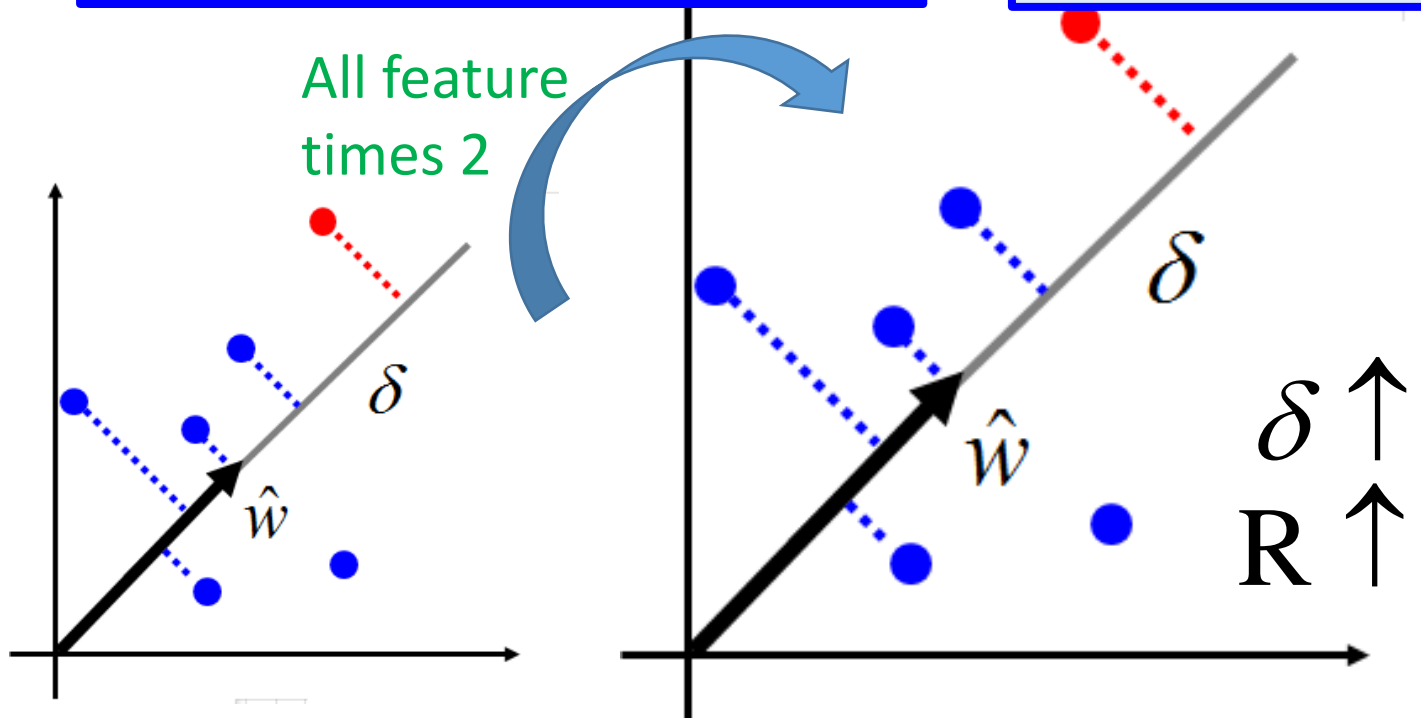
Normalization

Margin: Is it easy to separate red points from the blue ones

Larger margin, less update

• $\phi(x^r, \hat{y}^r)$

• $\phi(x^r, y)$



Structured Linear Model: Reduce 3 Problems to 2

Problem 1: Evaluation

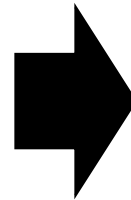
- How to define $F(x,y)$

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