# Introduction of Structured Learning

Hung-yi Lee

# Structured Learning

- We need a more powerful function $f$
  - Input and output are both objects with structures
  - *Object*: sequence, list, tree, bounding box …

$$f : X \rightarrow Y$$

**X** is the space of one kind of object

**Y** is the space of another kind of object

In the previous lectures, the input and output are both vectors.

# Introduction of Structured Learning

## Unified Framework

# Unified Framework

**Training**

- Find a function F

$$F: X \times Y \to R$$

- F(x,y): evaluate how compatible the objects x and y is
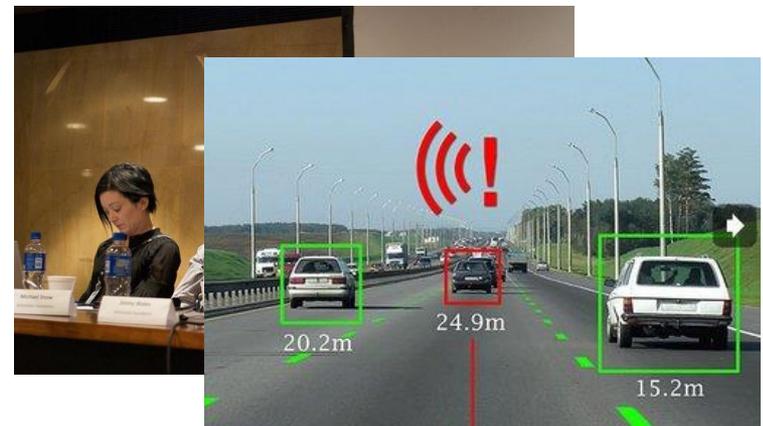
**Inference (Testing)**

- Given an object x

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

$$f : X \to Y \quad \blacktriangleright \quad f(x) = \tilde{y} = \arg\max_{y \in Y} F(x, y)$$

# Unified Framework – Object Detection

- Task description
  - Using a bounding box to highlight the position of a certain object in an image
  - E.g. A detector of Haruhi

$X$ : Image $\longrightarrow$ $Y$ : Bounding Box
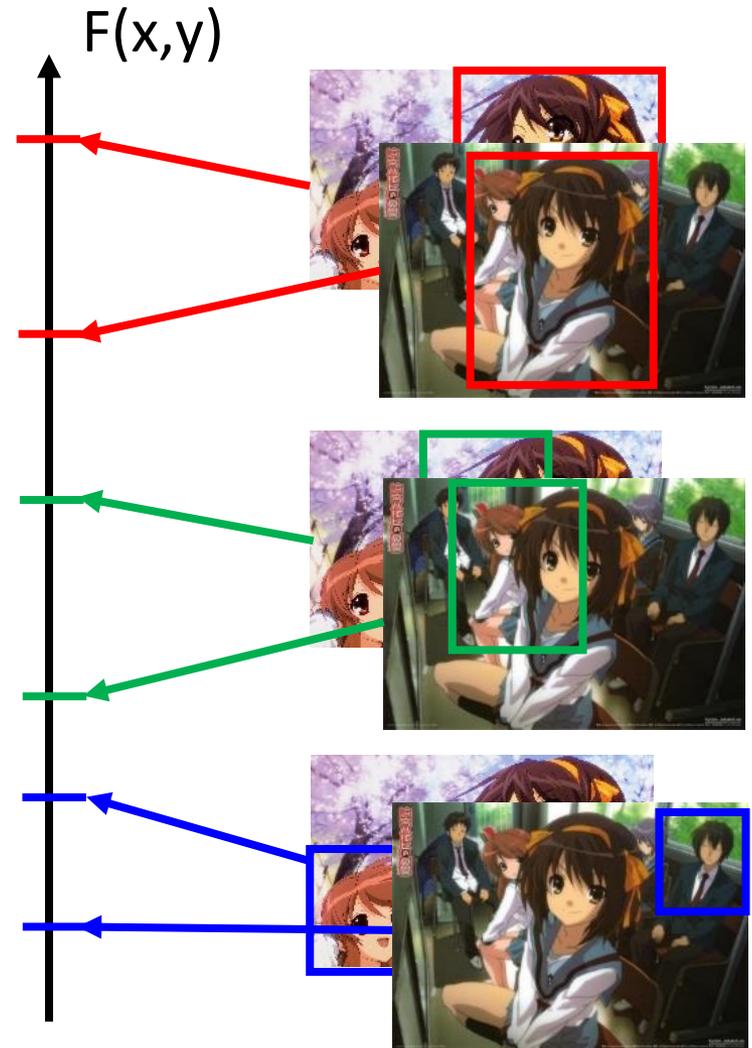
Haruhi

(the girl with yellow ribbon)

# Unified Framework – Object Detection

**Training**

- Find a function F

$$F : X \times Y \to R$$

- F(x,y): evaluate how compatible the objects x and y is

x: Image ➡

y: Bounding Box ➡

F(x,y) ➡ F( )

the correctness of taking range of y in x as "Haruhi"

F(x,y)

# Unified Framework – Object Detection

$\tilde{y}$ **(output result)**

F(x,y)

- Find a function F
$$F : X \times Y \rightarrow R$$
- F(x,y): evaluate how compatible the objects x and y is

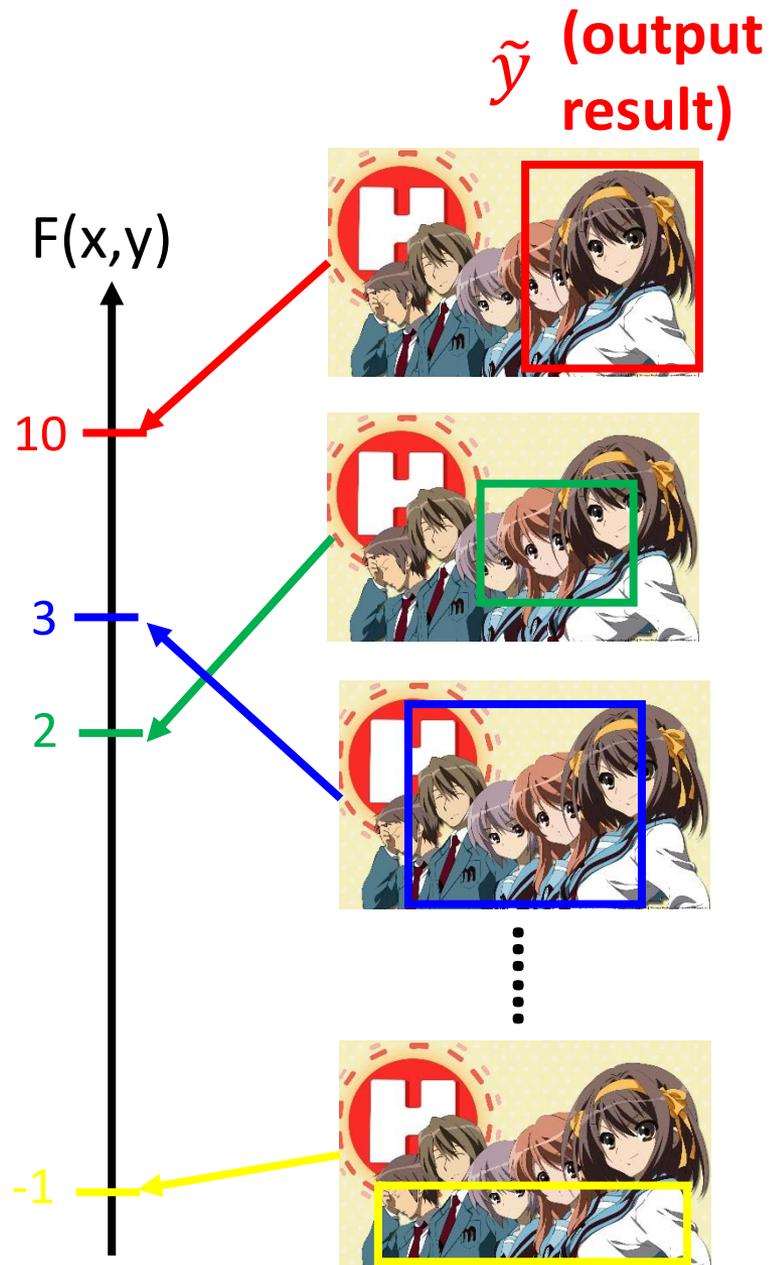**Inference (Testing)**

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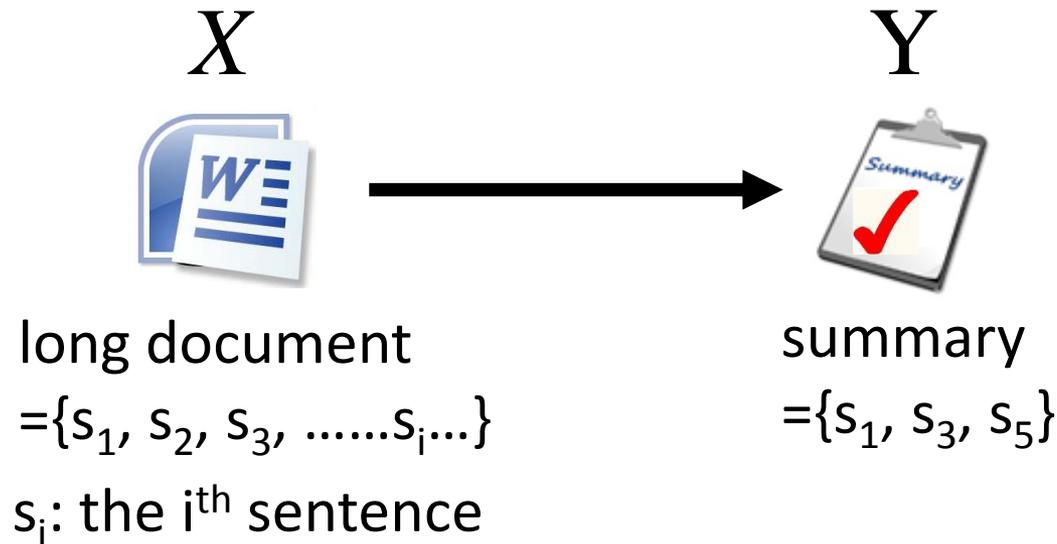$$\tilde{y} = \arg \max_{y \in Y} F(x, y)$$

input x =

Enumerate all possible bounding box y

10

3

2

-1

# Unified Framework
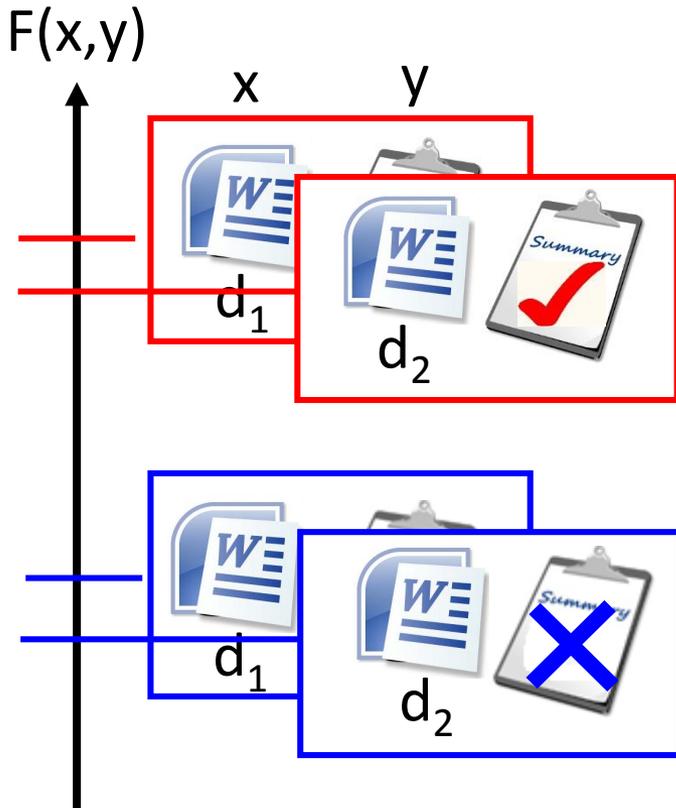# - Summarization

- Task description
  - Given a long document
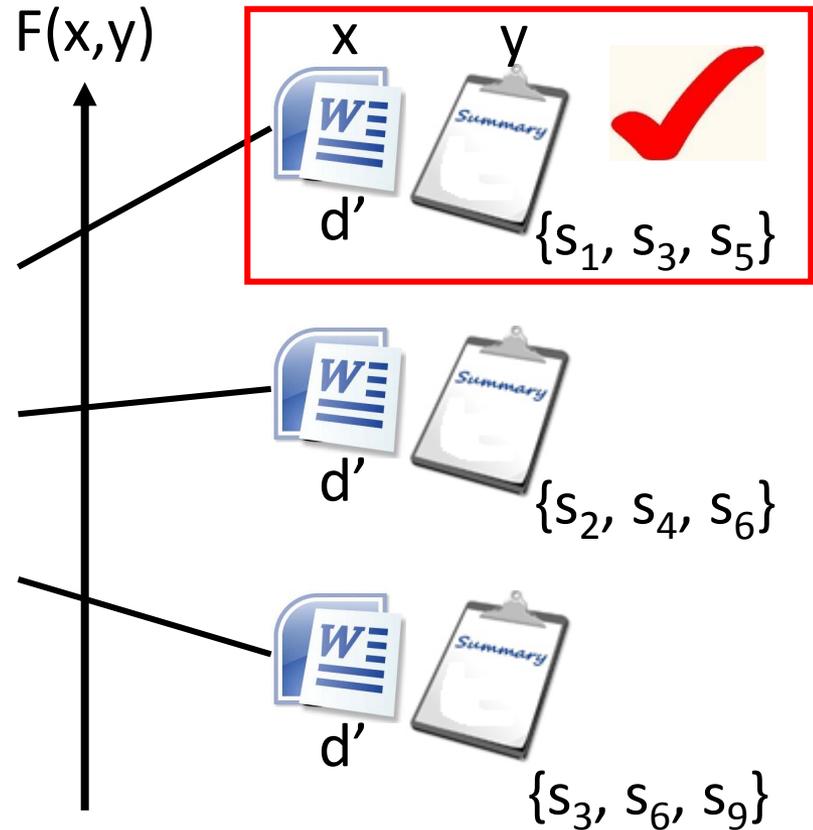  - Select a set of sentences from the document, and cascade the sentences to form a short paragraph

$X$ $\longrightarrow$ Y

long document
$=\{s_1, s_2, s_3, ......s_i...\}$

$s_i$: the $i^{th}$ sentence

summary
$=\{s_1, s_3, s_5\}$

# Unified Framework
## - Summarization



**Training**

**Inference**

F(x,y)

x    y

$d_1$
$d_2$

$d_1$
$d_2$

F(x,y)

x    y

$d'$    $\{s_1, s_3, s_5\}$
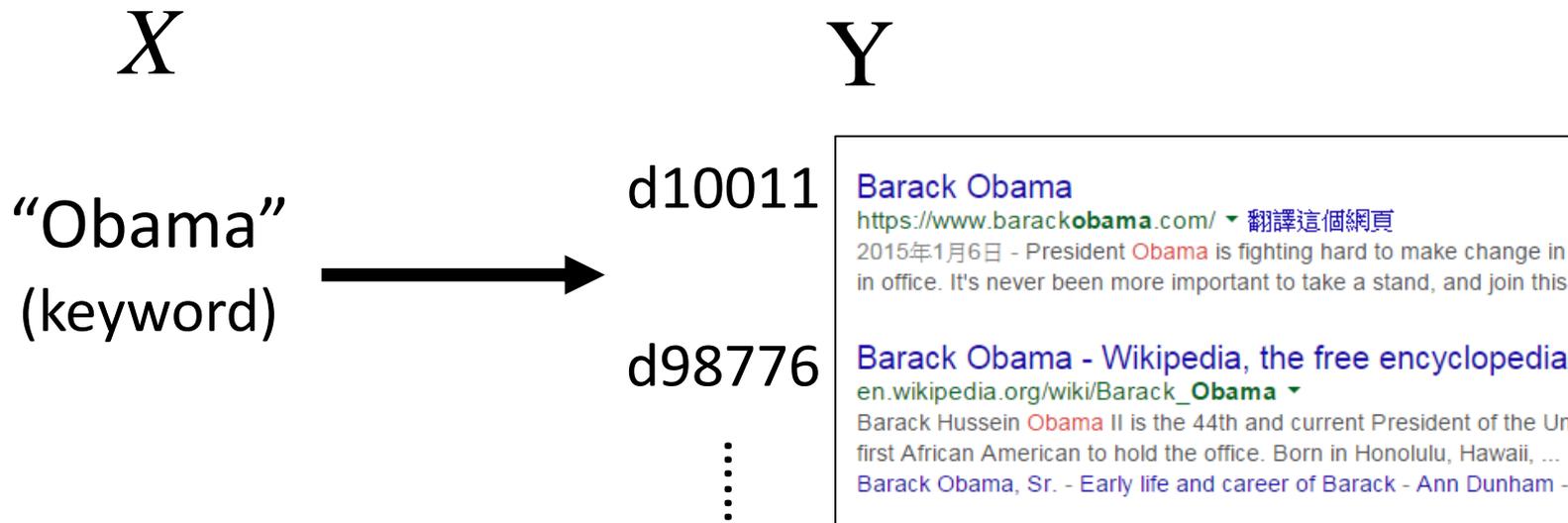
$d'$    $\{s_2, s_4, s_6\}$

$d'$    $\{s_3, s_6, s_9\}$

# Unified Framework
# - Retrieval

- Task description
    - User input a keyword Q
    - System returns a *list* of web pages

$X$

$Y$

"Obama"
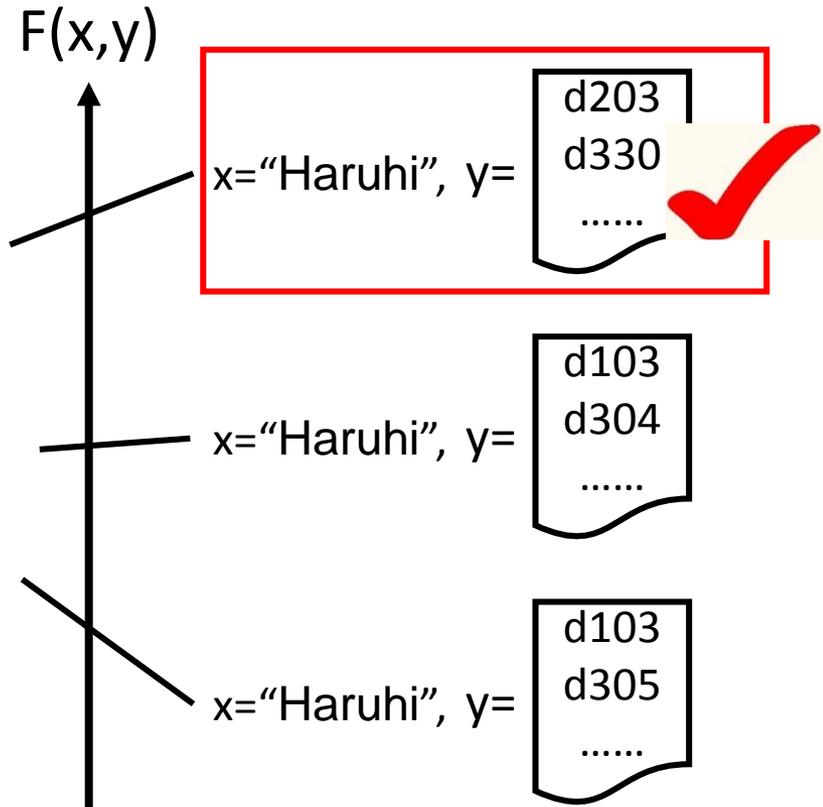(keyword)

→

d10011



Barack Obama
https://www.barackobama.com/ ▼ 翻譯這個網頁
2015年1月6日 - President Obama is fighting hard to make change in
in office. It's never been more important to take a stand, and join this

d98776

Barack Obama - Wikipedia, the free encyclopedia
en.wikipedia.org/wiki/Barack_Obama ▼
Barack Hussein Obama II is the 44th and current President of the Un
first African American to hold the office. Born in Honolulu, Hawaii, ...
Barack Obama, Sr. - Early life and career of Barack - Ann Dunham -

A list of web pages (Search Result)

# Unified Framework
## - Retrieval



Training

Inference

F(x,y)

F(x,y)

x="Trump", y= d103 d300 ...... ✔

x="Trump", y= d133 d220 ✖ ......

d666

d103

x="Haruhi", y= d203 d330 ...... ✔

x="Haruhi", y= d103 d304 ......

x="Haruhi", y= d103 d305 ......

# *Statistics*
# *Unified Framework*

## Training

- Find a function F
$$F : X \times Y \to \mathrm{R}$$
- F(x,y): evaluate how compatible the objects x and y is

## Inference

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

$$F(x, y) = P(x, y)?$$

## Training

- Estimate the probability P(x,y)
$$P : X \times Y \to [0,1]$$

## Inference

- Given an object x
$$\tilde{y} = \arg\max_{y \in Y} P(y \mid x)$$
$$= \arg\max_{y \in Y} \frac{P(x, y)}{P(x)}$$
$$= \arg\max_{y \in Y} P(x, y)$$

# Statistics

## Unified Framework

$$F(x, y) = P(x, y)?$$

### Drawback for probability

- Probability cannot explain everything
- 0-1 constraint is not necessary

### Strength for probability

- Meaningful

Energy-based Model: http://www.cs.nyu.edu/~yann/research/ebm/

## Training

- Estimate the probability P(x,y)

$$P : X \times Y \to [0,1]$$

## Inference

- Given an object x

$$\tilde{y} = \arg\max_{y \in Y} P(y \mid x)$$

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

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

# Unified Framework

That's it!?

## Training

- Find a function F

$$F: X \times Y \rightarrow R$$

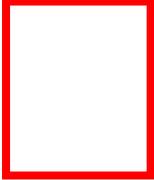- F(x,y): evaluate how compatible the objects x and y is

## Inference (Testing)

- Given an object x

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

There are three problems in this framework.

# Problem 1

- ***Evaluation***: What does F(x,y) look like?
  - How F(x,y) compute the "compatibility" of objects x and y

***Object Detection:*** F(x=        , y=        )

***Summarization:*** F(x=        , y=        )

(a long document)        (a short paragraph)

***Retrieval:*** F(x= "Obama"        , y=        )

(keyword)

(Search Result)

# Problem 2

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

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

The space *Y* can be extremely large!

**_Object Detection:_**  *Y*=All possible bounding box (maybe tractable)

**_Summarization:_**  *Y*=All combination of sentence set in a document …

**_Retrieval:_** *Y*=All possible webpage ranking ….

# Problem 3

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

## *Principle*

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 F(x,y) such that ……

# Three Problems

## Problem 1: Evaluation

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

## Problem 2: Inference

- How to solve the "arg max" problem

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

## Problem 3: Training

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

From 數位語音處理

# Link to DNN?

**Training**

$$F : X \times Y \to \mathrm{R}$$

$$F(x, y) = -CE(N(x), y)$$

$$CE(N(x), y)$$
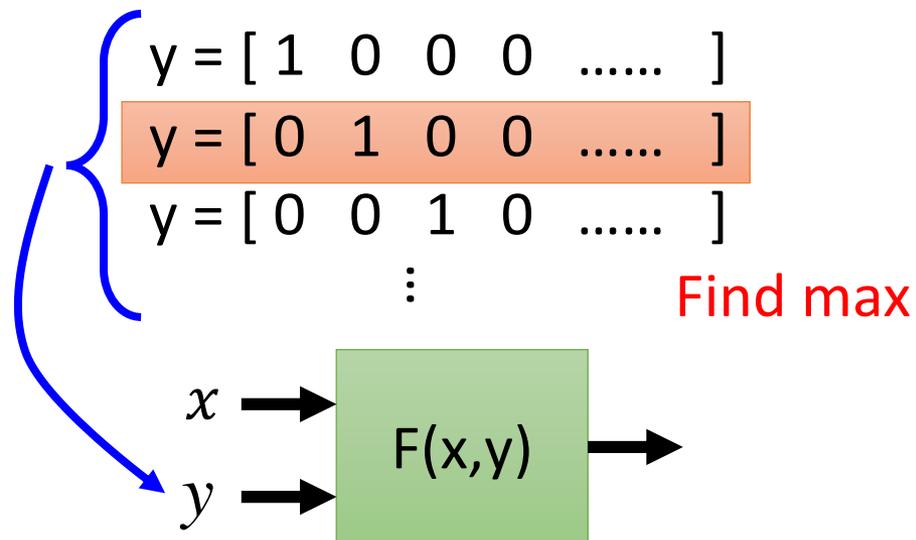
N(x)

DNN

$x$       $y$

**Inference**

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

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

y = [ 1  0  0  0  ……  ]
y = [ 0  1  0  0  ……  ]
y = [ 0  0  1  0  ……  ]
⋮

**Find max**
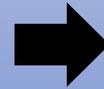
$x$ →

$y$ →   F(x,y) →

# Introduction of Structured Learning

## Linear Model

# Structured Linear Model

**Problem 1: Evaluation**

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

**Problem 2: Inference**

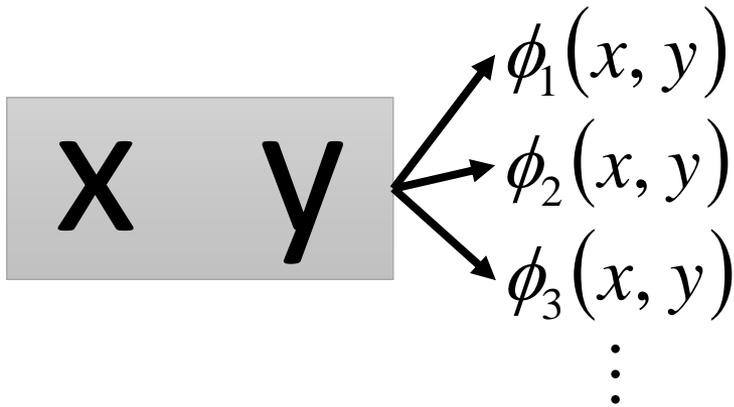- How to solve the "arg max" problem

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

**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

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

X    y

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

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

Learning from data

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

# Structured Linear Model: Problem 1

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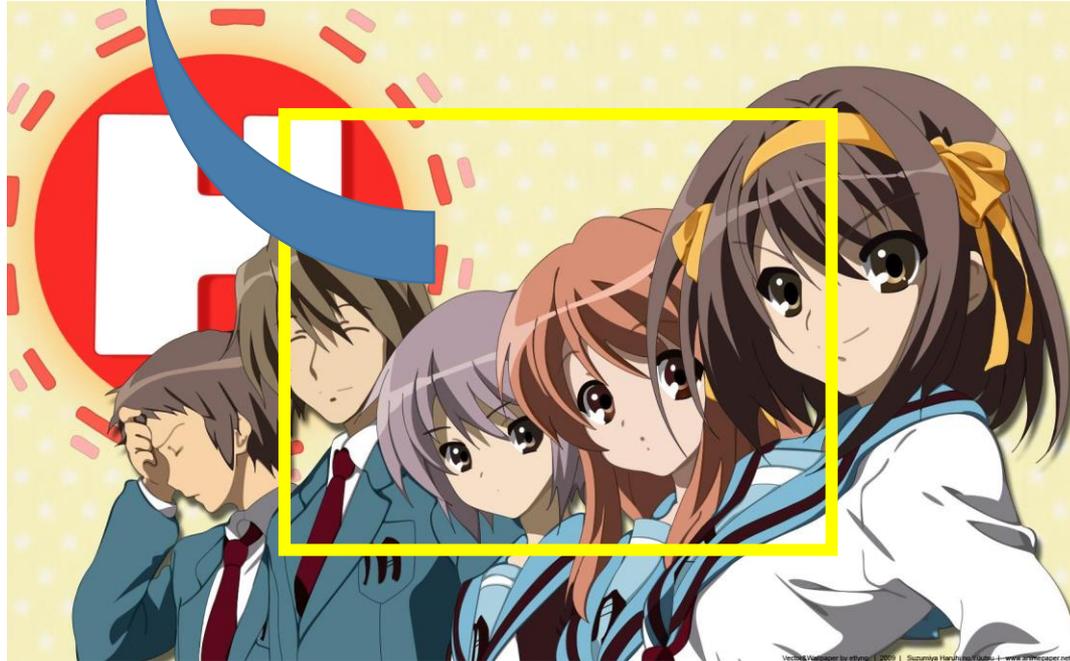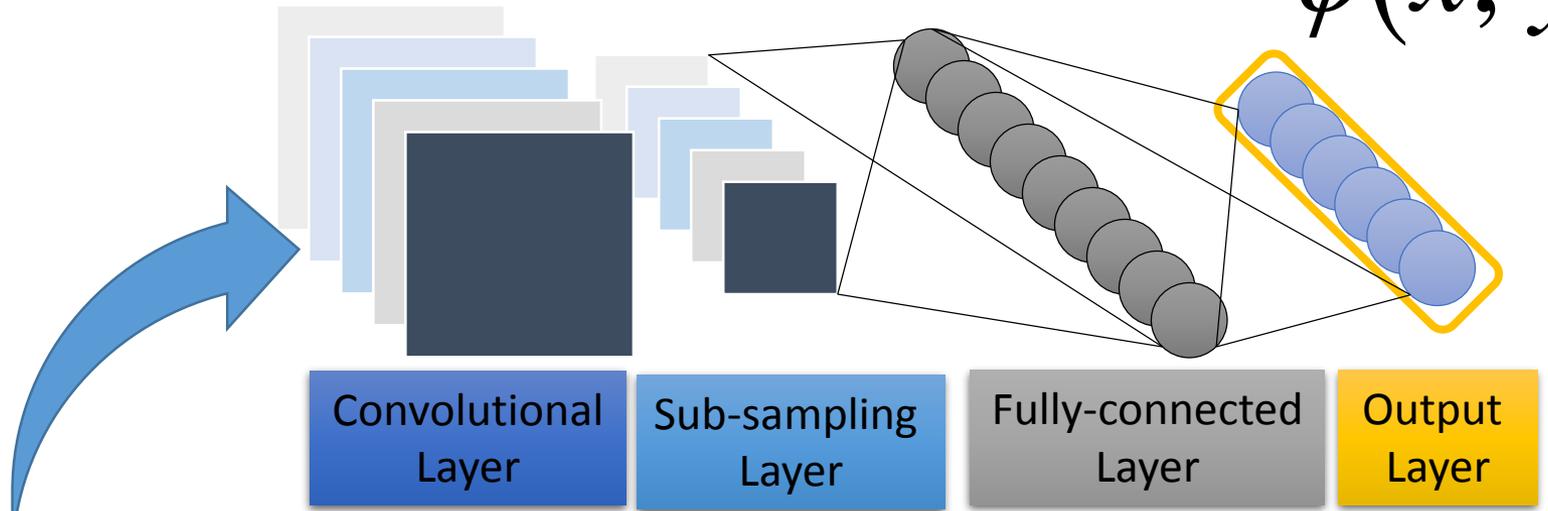- Example: ***Object Detection***



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

$\phi(x, y)$

Convolutional Layer

Sub-sampling Layer

Fully-connected Layer

Output Layer

$\phi(\quad)$

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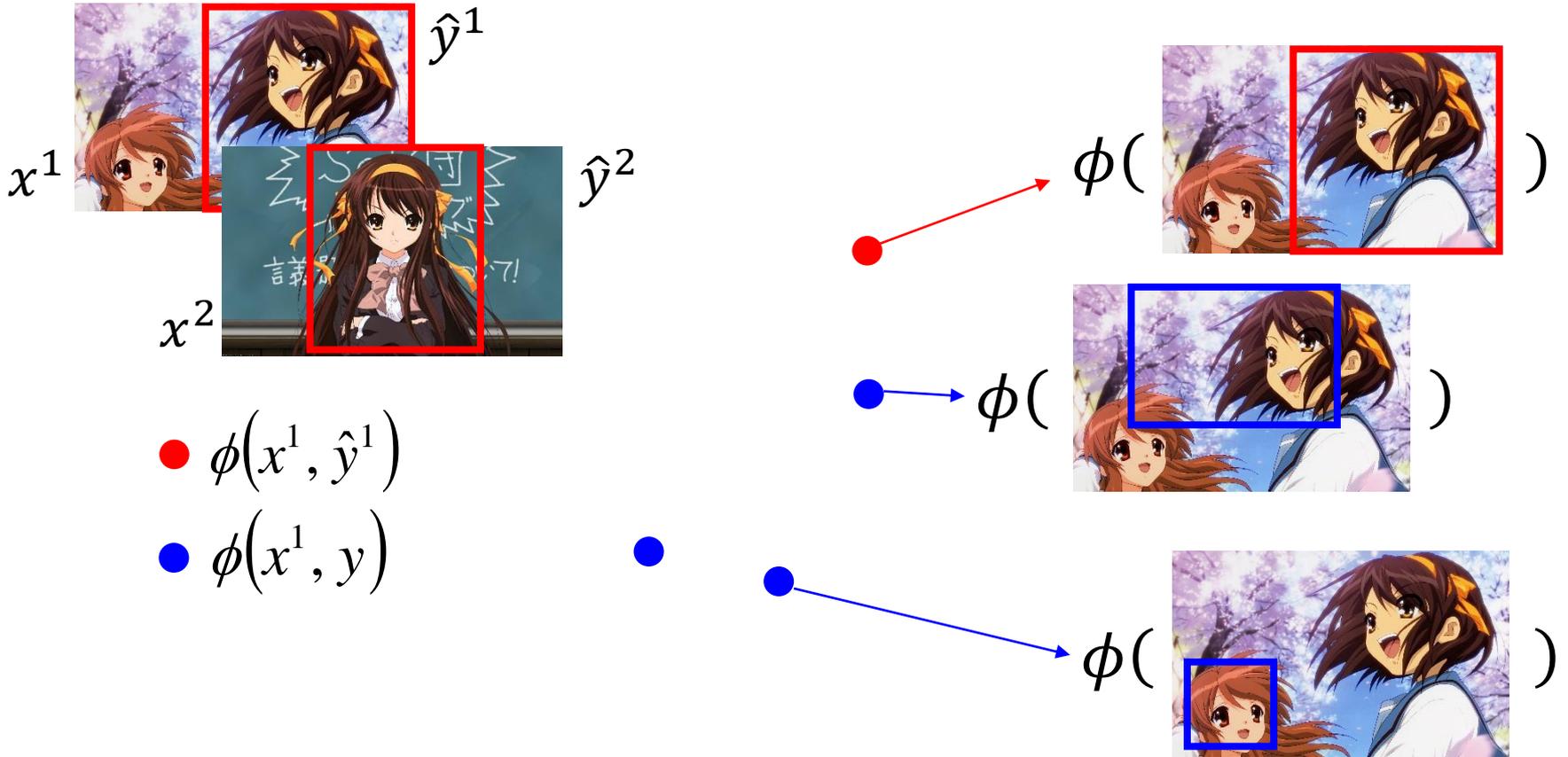
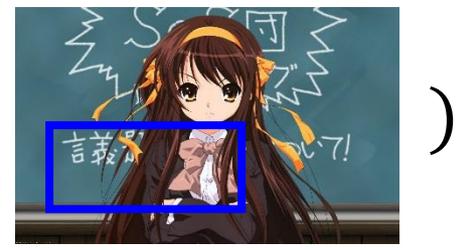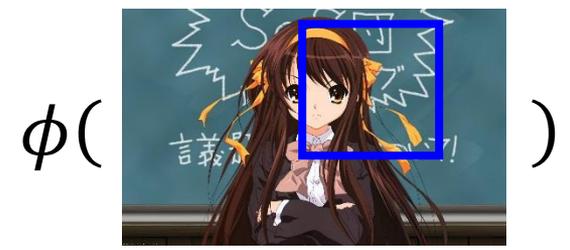We should find w such that

$\forall r$ (All training examples)

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

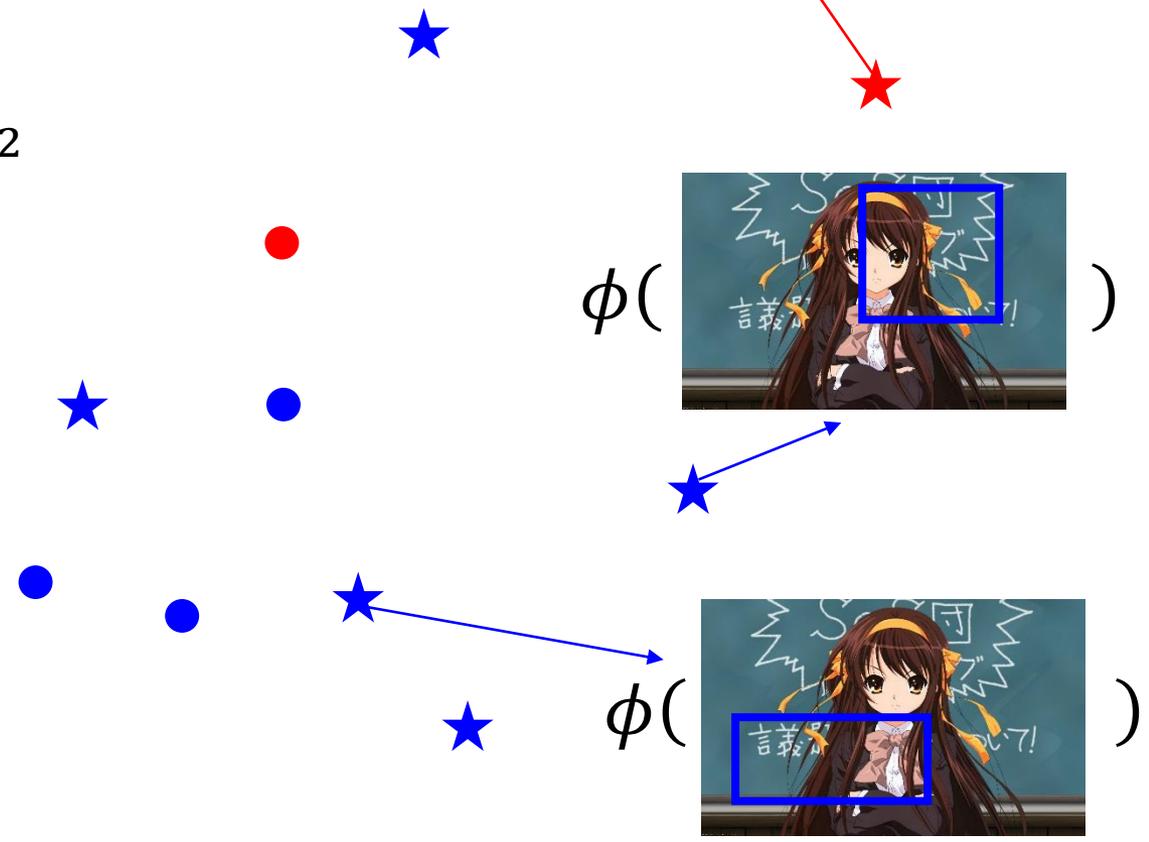$$w \cdot \phi\left(x^r, \hat{y}^r\right) > w \cdot \phi\left(x^r, y\right)$$

# Structured Linear Model:
## Problem 3



$\hat{y}^1$

$x^1$

$\hat{y}^2$

$x^2$

● $\phi(x^1, \hat{y}^1)$

● $\phi(x^1, y)$

$\phi($   $)$

$\phi($   $)$

$\phi($   $)$

# Structured Linear Model:
## Problem 3



$\hat{y}^1$

$x^1$

$x^2$

$\hat{y}^2$

$\phi($ )

$\phi($ )

$\phi($ )

● $\phi\left(x^1, \hat{y}^1\right)$

● $\phi\left(x^1, y\right)$

★ $\phi\left(x^2, \hat{y}^2\right)$

★ $\phi\left(x^2, y\right)$

# Structured Linear Model:
## Problem 3



$x^1$

$\hat{y}^1$

$\hat{y}^2$

$x^2$

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

$w$

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

# Solution of Problem 3
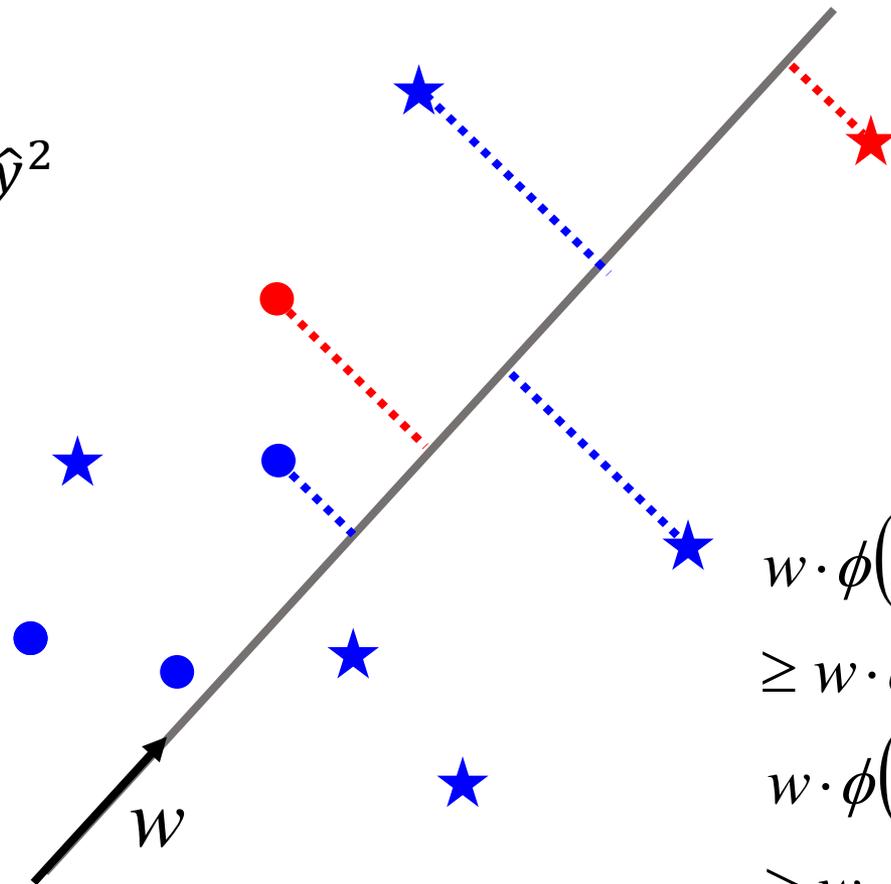
Difficult?

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

$\hat{y}^1$

$x^1$

$\hat{y}^2$

$x^2$

- 🔴 $\phi\left(x^1, \hat{y}^1\right)$

- 🔵 $\phi\left(x^1, y\right)$

- ⭐ $\phi\left(x^2, \hat{y}^2\right)$

- ⭐ $\phi\left(x^2, y\right)$

$w$

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

# Algorithm - Example

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

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

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

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

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

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

$w$

$\widetilde{y}^1$

Because w=0 at this time, $\phi(x^1, y)$ always 0
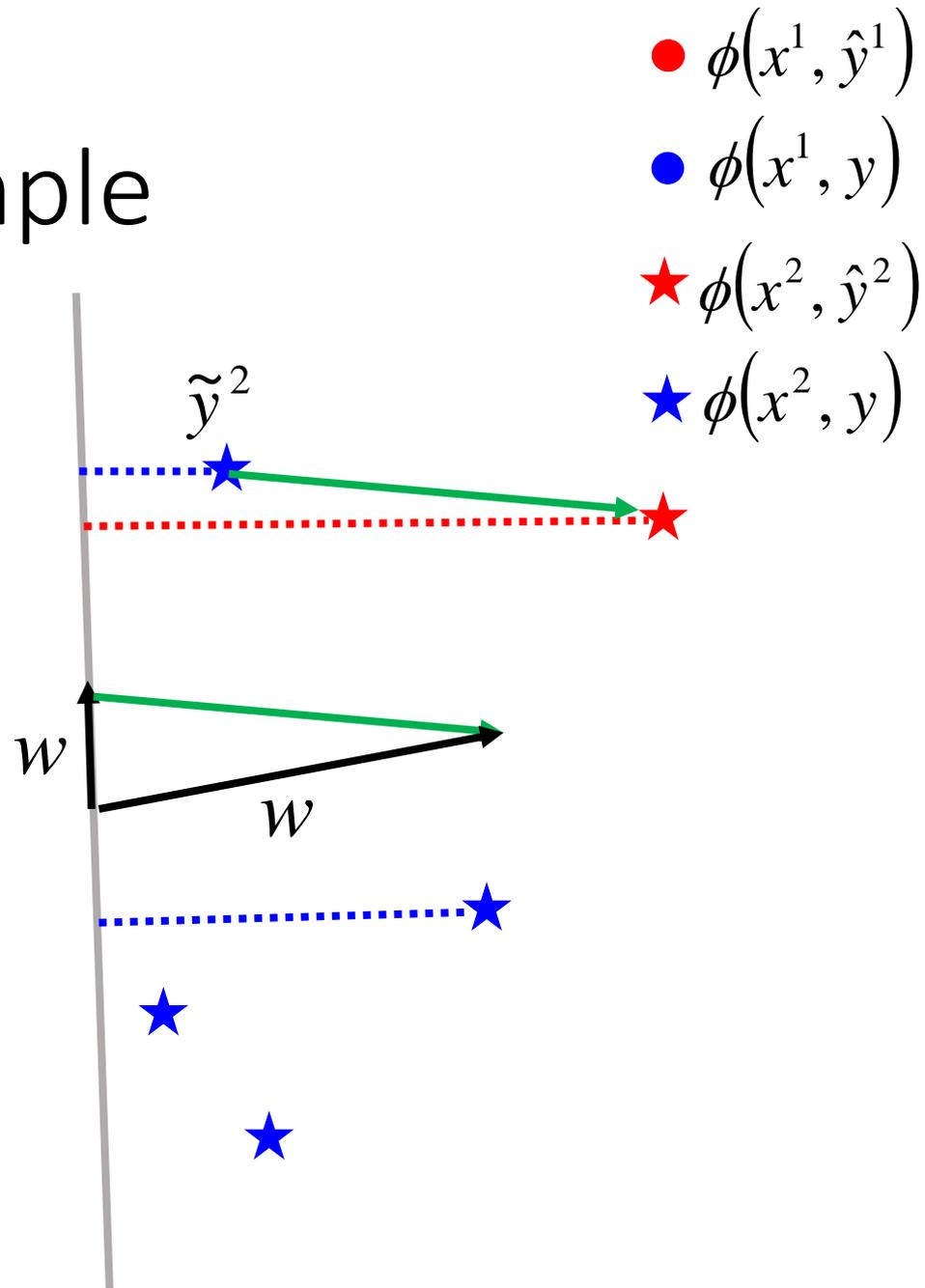
➡ Random pick one point as $\widetilde{y}^r$

# Algorithm - Example



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

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

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

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

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

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

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

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

$\tilde{y}^2$

$w$

$w$

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

$$\widetilde{y}^1 = \arg\max_{y \in Y} w \cdot \phi(x^1, y)$$
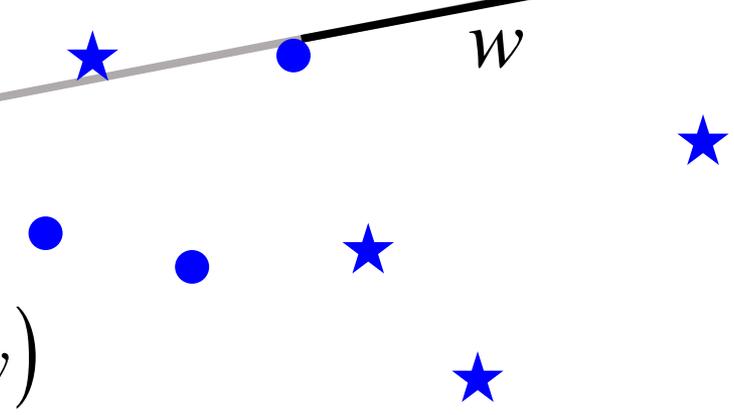
$\widetilde{y}^1 = \hat{y}^1$ ➡ do not update w

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

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

$w$

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

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

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

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

So we are done

# Assumption: Separable

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

$\forall r$ (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)

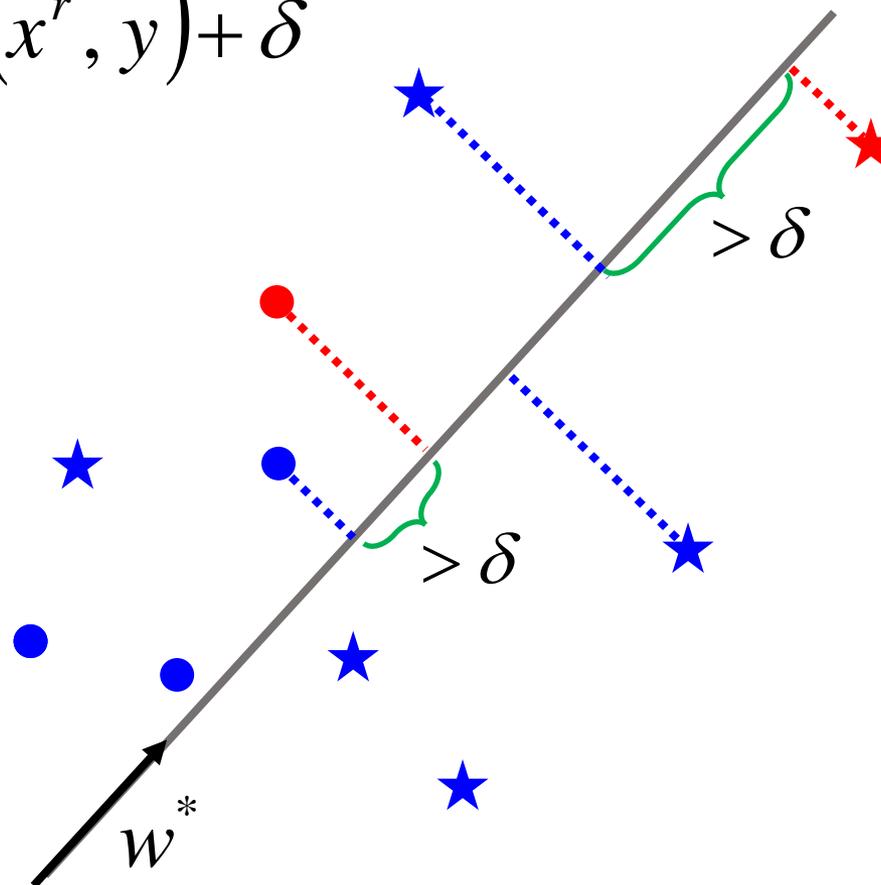$\hat{w} \cdot \phi(x^r, \hat{y}^r) \geq \hat{w} \cdot \phi(x^r, y) + \delta$

# 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>

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

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

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

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

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

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

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

# Proof of Termination

w is updated once it sees a mistake

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

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

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

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

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

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

$$\hat{w} \cdot w^1 \geq \delta \qquad\qquad \hat{w} \cdot w^2 \geq 2\delta \qquad \cdots\cdots \qquad \text{(so what)}$$

# Proof of Termination

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

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

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

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

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

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

$$\cdots$$

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

$$\left\|w^k\right\|^2 \le k\mathbf{R}^2$$

# 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

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

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

$\delta$

$\hat{w}$

$\delta$

$\hat{w}$

$\delta \uparrow$

$R \uparrow$

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

# Graphical Model

A language which describes the evaluation function

# Structured Learning

We also know how to involve hidden information.

## Problem 1: Evaluation

- What does F(x,y) look like? $F(x,y) = w \cdot \phi(x,y)$

## 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 SVM, etc.

# Difficulties

**_Difficulty 1. Evaluation_** ➡️ Graphical Model

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



$\phi(x, y)$ | 1.2 | 2.6 | 2.7 | 2.3 | 1.5 | 2.5 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 |

Hard to figure out? Hard to interpret the meaning?

**_Difficulty 2. Inference_** ➡️ Gibbs Sampling

We can use Viterbi algorithm to deal with sequence labeling. How about other cases?

# Graphical Model

$$F(x, y) \longleftrightarrow \boxed{\text{Graph}}$$

- Define and describe your evaluation function F(x,y) by a graph

- There are three kinds of graphical model.
  - *Factor graph, Markov Random Field* (MRF) and *Bayesian Network* (BN)
  - Only *factor graph* and *MRF* will be briefly mentioned today.

# Decompose F(x,y)

- $F(x, y)$ is originally a ***global*** function
  - Define over the whole x and y
- Based on graphical model, $F(x, y)$ is the composition of some ***local*** functions
  - x and y are decomposed into smaller components
  - Each local function defines on only a few related components in x and y
  - Which components are related → defined by Graphical model

# Decomposable x and y

- x and y are decomposed into smaller components

**_POS Tagging_**

x: 
| $x_1$ | $x_2$ | $x_3$ | $x_4$ |
| John | saw | the | saw. |

y:
| $y_1$ | $y_2$ | $y_3$ | $y_4$ |
| PN | V | D | N |

x: $x_1$ $x_2$ $x_3$ $x_4$   {word}

y: $y_1$ $y_2$ $y_3$ $y_4$   {tags}

# Factor Graph

Each factor influences some components.

Each factor corresponds to a local function.



factor a
$$f_a(x_1, y_1)$$

factor b
$$f_b(x_2, y_1, y_2)$$

factor c
$$f_d(y_2)$$

Larger value means more compatible.

$$F(x, y) = f_a(x_1, y_1) + f_b(x_2, y_1, y_2) + f_c(y_2)$$

You only have to define the factors.

The local functions of the factors are learned from data.

# Factor Graph - Example

- ***Image De-noising***

Each pixel is one component

Noisy image x

Clean image y

$\{-1,1\}$ $x_1$ $x_2$ $x_3$ $x_4$ $x_5$ $x_6$ $x_7$ $x_8$ $x_9$

$\{-1,1\}$ $y_1$ $y_2$ $y_3$ $y_4$ $y_5$ $y_6$ $y_7$ $y_8$ $y_9$

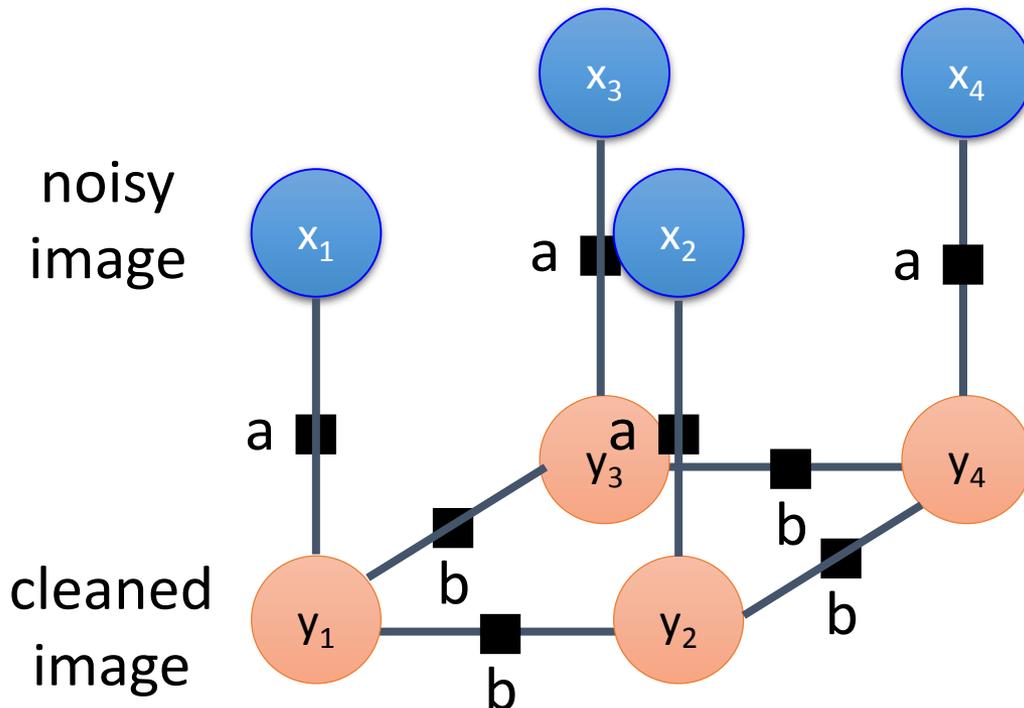http://cs.stanford.edu/people/karpathy/visml/ising_example.html

# Factor Graph - Example

Noisy and clean images are related
➢ **a**: the values of $x_i$ and $y_i$

**_Factor:_**

The colors in the clean image is smooth.
➢ **b**: the values of the neighboring $y_i$

noisy image

cleaned image

$$f_a(x_i, y_i) = \begin{cases} 1 & x_i = y_i \\ -1 & x_i \neq y_i \end{cases}$$

$$f_b(y_i, y_j) = \begin{cases} 2 & y_i = y_j \\ -2 & y_i \neq y_j \end{cases}$$

The weights can be learned from data.

# Factor Graph - Example

**Noisy and clean images are related**

➢ **a**: the values of $x_i$ and $y_i$

***Factor:***

**The colors in the clean image is smooth.**

➢ **b**: the values of the neighboring $y_i$



noisy image

cleaned image

Realize $F(x, y)$ easily from the factor graph

$$F(x, y) = \sum_{i=1}^{4} f_a(x_i, y_i)$$

$$+ f_b(x_1, y_2) + f_b(x_1, y_3)$$
$$+ f_b(x_2, y_4) + f_b(x_3, y_4)$$

# Factor Graph - Example

***Factor:***

➢ **c**: the values of $x_i$ and the values of the neighboring $y_i$

➢ **d**: the values of the neighboring $x_i$ and the values of $y_i$



$$f_c(x_i, y_i, y_{i-1})$$

$$f_d(x_i, x_{i-1}, y_i)$$

$$f_e(x_i, x_{i-1}, y_i, y_{i-1})$$

# Markov Random Field (MRF)

Clique: a set of components connecting to each other

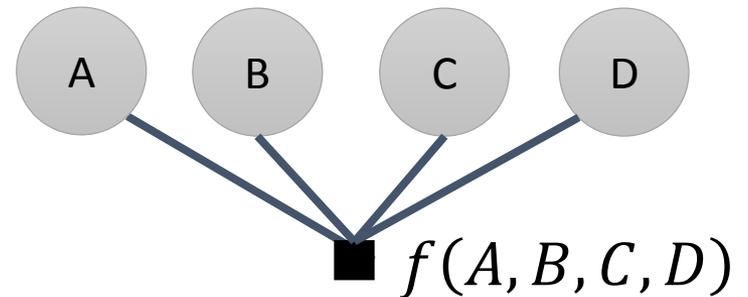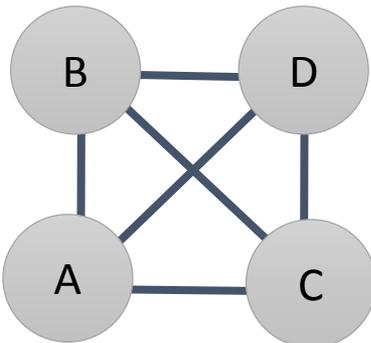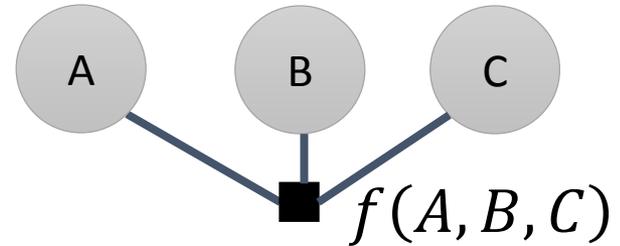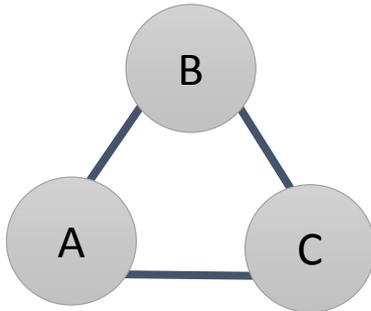Maximum Clique: a clique that is not included by other cliques
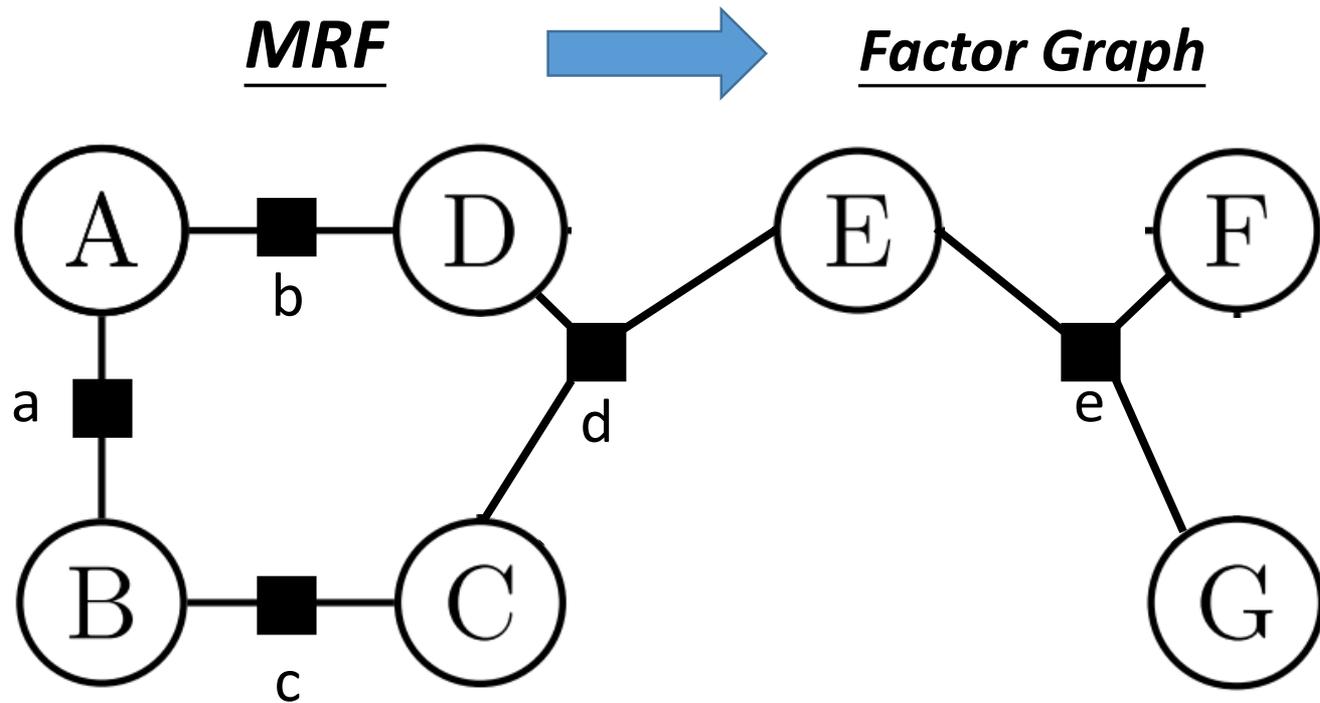
# MRF

Each maximum clique on the graph corresponds to a factor
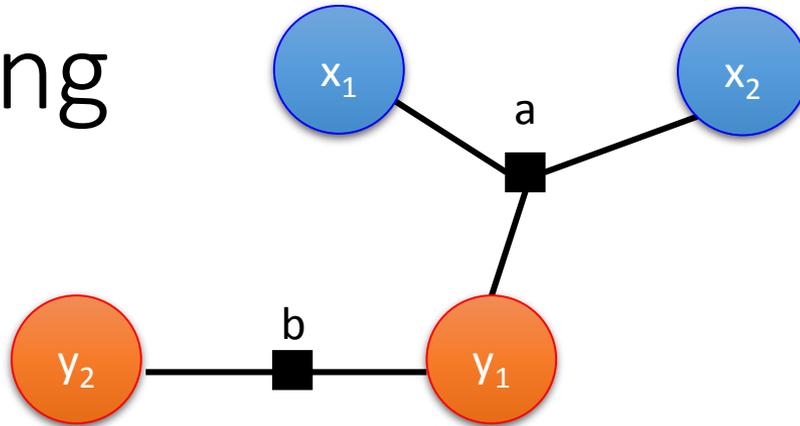
**_MRF_**

**_Factor Graph_**



$f(A,B)$

$f(A,B,C)$

$f(A,B,C,D)$

# MRF

➡ **_Factor Graph_**



**_Evaluation Function_**

$$f_a(A, B) + f_b(A, D) + f_c(B, C) + f_d(C, D, E) + f_e(E, F, G)$$

# Training
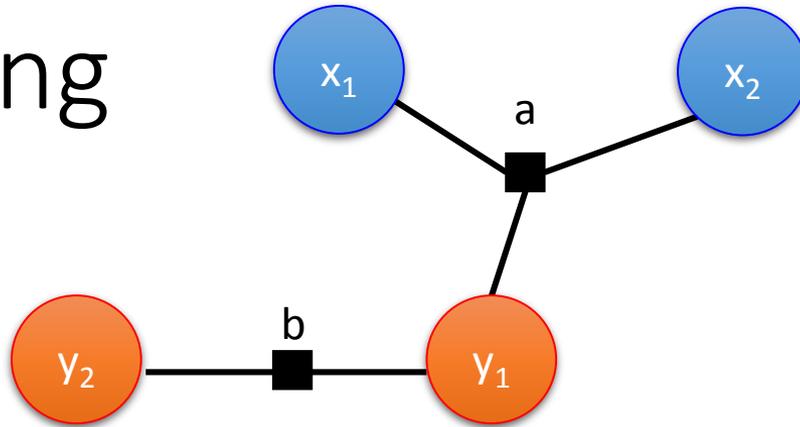


$$F(x, y) = f_a(x_1, x_2, y_1) + f_b(y_1, y_2)$$

$$= w_a \cdot \phi_a(x_1, x_2, y_1) + w_b \cdot \phi_b(y_1, y_2)$$

$$= \begin{bmatrix} w_a \\ w_b \end{bmatrix} \begin{bmatrix} \phi_a(x_1, x_2, y_1) \\ \phi_b(y_1, y_2) \end{bmatrix}$$

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

Simply training by
***structured perceptron***
***or structured SVM***

Max-Margin Markov Networks (M3N)

# Training



$$\phi_b(+1,+1) = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$F(x,y) = f_a(x_1, x_2, y_1) + \underline{f_b(y_1, y_2)}$$

$$= w_a \cdot \phi_a(x_1, x_2, y_1) + \underline{w_b \cdot \phi_b(y_1, y_2)}$$

$$y_1, y_2 \epsilon \{+1, -1\}$$
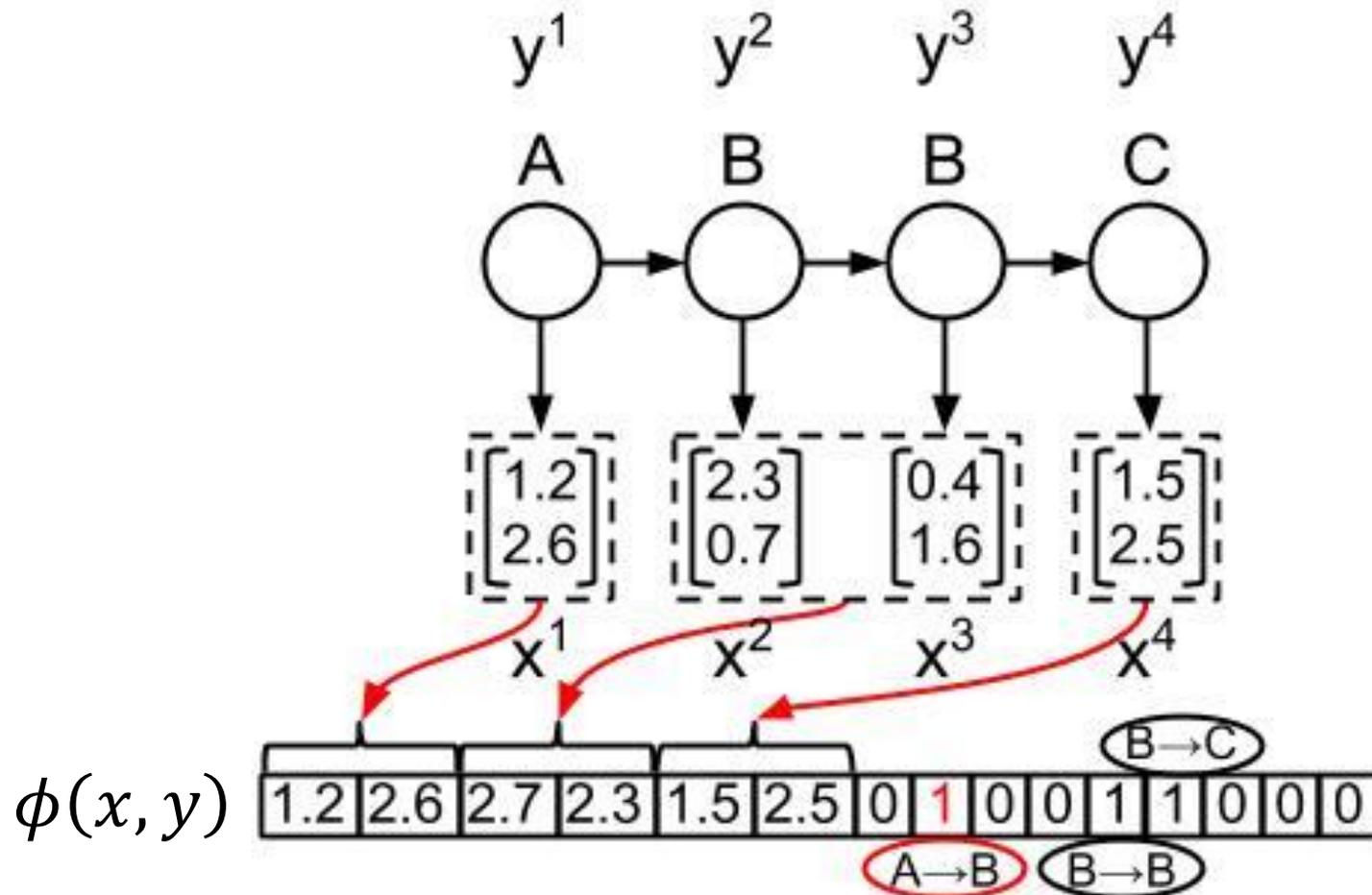
$$\phi_b(+1,-1) = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\phi_b(-1,+1) = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

| $y_1$ | $y_2$ | $f_b(y_1, y_2)$ |
|-------|-------|-----------------|
| +1 | +1 | $w_1$ |
| +1 | -1 | $w_2$ |
| -1 | +1 | $w_3$ |
| -1 | -1 | $w_4$ |

$$w_b = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix}$$

$$\phi_b(-1,-1) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

# Now can you interpret this?

# Probability Point of View

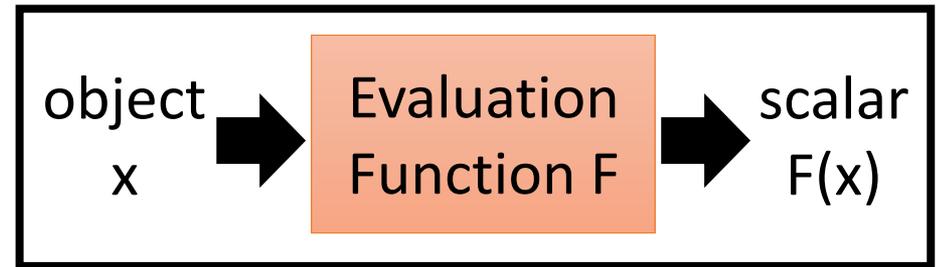- $F(x, y)$ can be any real number

- If you like probability

Between 0 and 1

$$P(x, y) = \frac{e^{F(x,y)}}{\sum_{x',y'} e^{F(x',y')}}$$

→ To be positive

→ normalization
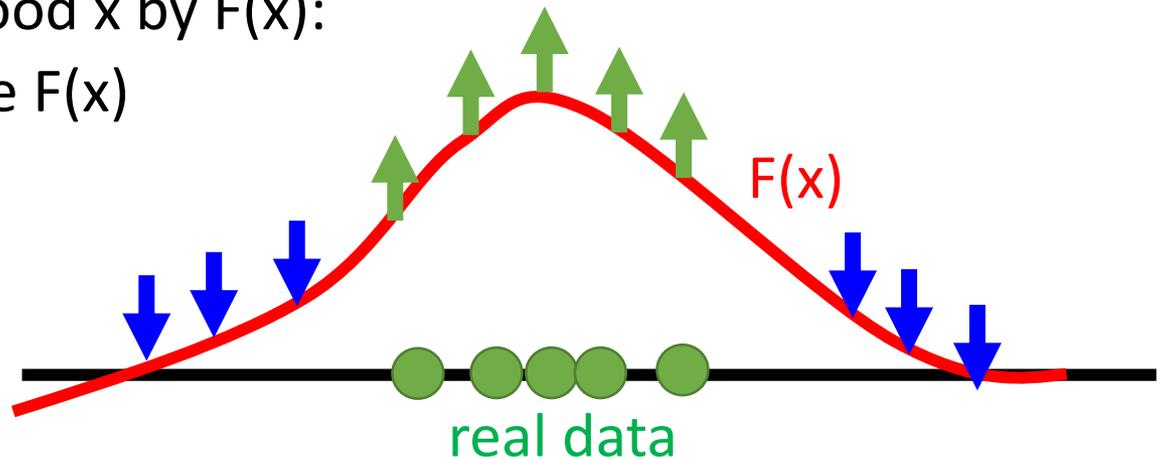
$$P(y|x) = \frac{P(x, y)}{P(x)}$$

$$= \frac{P(x, y)}{\sum_{y''} P(x, y'')} = \frac{\dfrac{e^{F(x,y)}}{\sum_{x',y'} e^{F(x',y')}}}{\sum_{y''} \dfrac{e^{F(x,y'')}}{\sum_{x',y'} e^{F(x',y')}}} = \frac{e^{F(x,y)}}{\sum_{y''} e^{F(x,y'')}}$$

# Evaluation Function

- We want to find an evaluation function F(x)
  - Input: object x, output: scalar F(x) (how "good" the object is)
  - E.g. x are images
    - Real x has high F(x)
  - F(x) can be a network

- We can generate good x by F(x):
  - Find x with large F(x)

- How to find F(x)?

In practice, you cannot decrease all the x other than real data.

object x → Evaluation Function F → scalar F(x)



F(x)

real data

# Evaluation Function
## - Structured Perceptron

- **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

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

  - do

    - For each pair of training example $\left( x^r, \hat{y}^r \right)$

      - Find the label $\tilde{y}^r$ maximizing $F(x^r, y)$

Can be an issue $\Longrightarrow$ $\tilde{y}^r = \arg \max_{y \in Y} F\left( x^r, y \right)$

      - If $\tilde{y}^r \neq \hat{y}^r$, update w
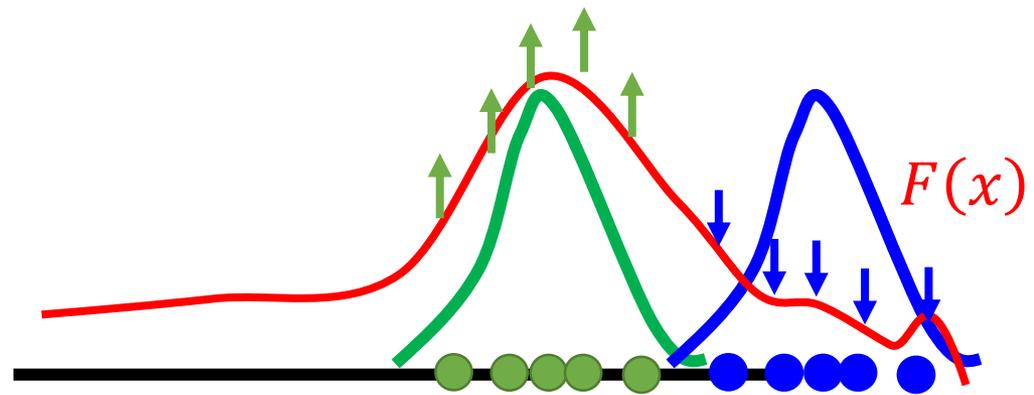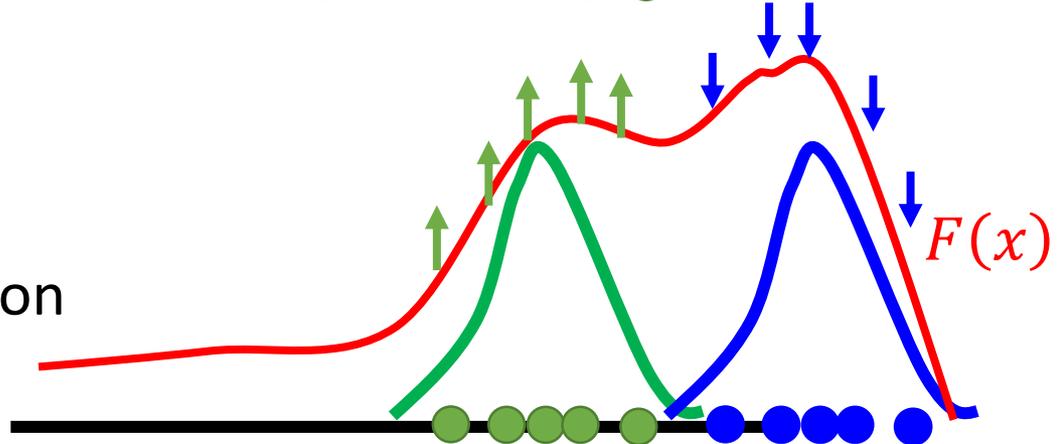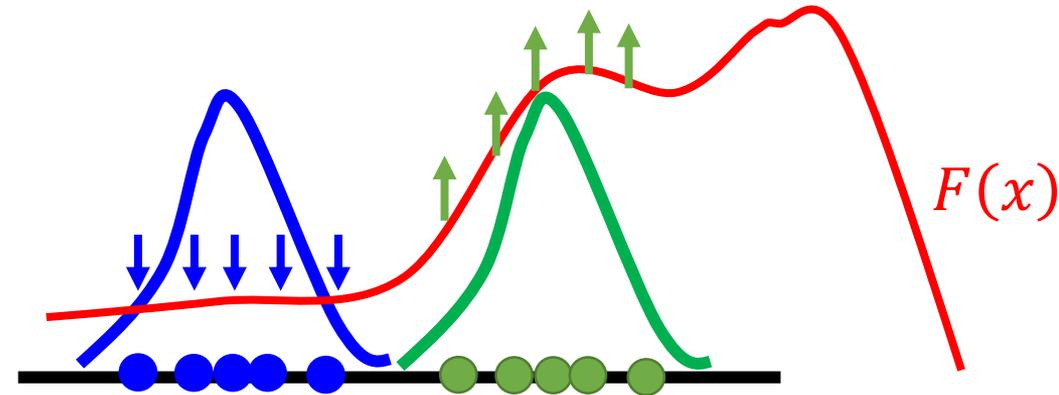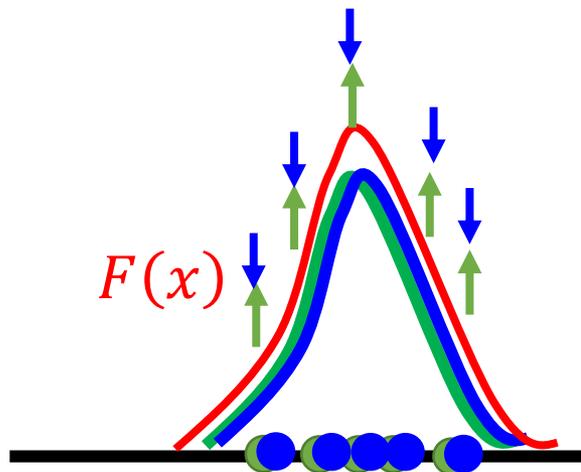
Increase $F(x^r, \hat{y}^r)$, decrease $F(x^r, \tilde{y}^r)$

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

  - until w is not updated $\Longrightarrow$ We are done!

# How about GAN?

- Generator is an intelligent way to find the negative examples.

"Experience replay", parameters from last iteration

In the end ……

# Where are we?



Restricted Boltzmann Machine

Boltzmann Machine

Undirected Graph
(MRF, factor graph, etc.)

Graphical Model

Structured Learning