

# Unsupervised Learning

Hung-yi Lee

# Introduction

- We already learn some machine learning techniques.
- With labelled data, you can do any thing (hopefully).
- Labelling data is expensive.
- What can we do if there is no sufficient training data?
- Unsupervised Learning Approaches
  - Restricted Boltzmann Machine (RBM)
  - Auto-encoder

# Semi-supervised Learning

Labelled  
data



cat



dog

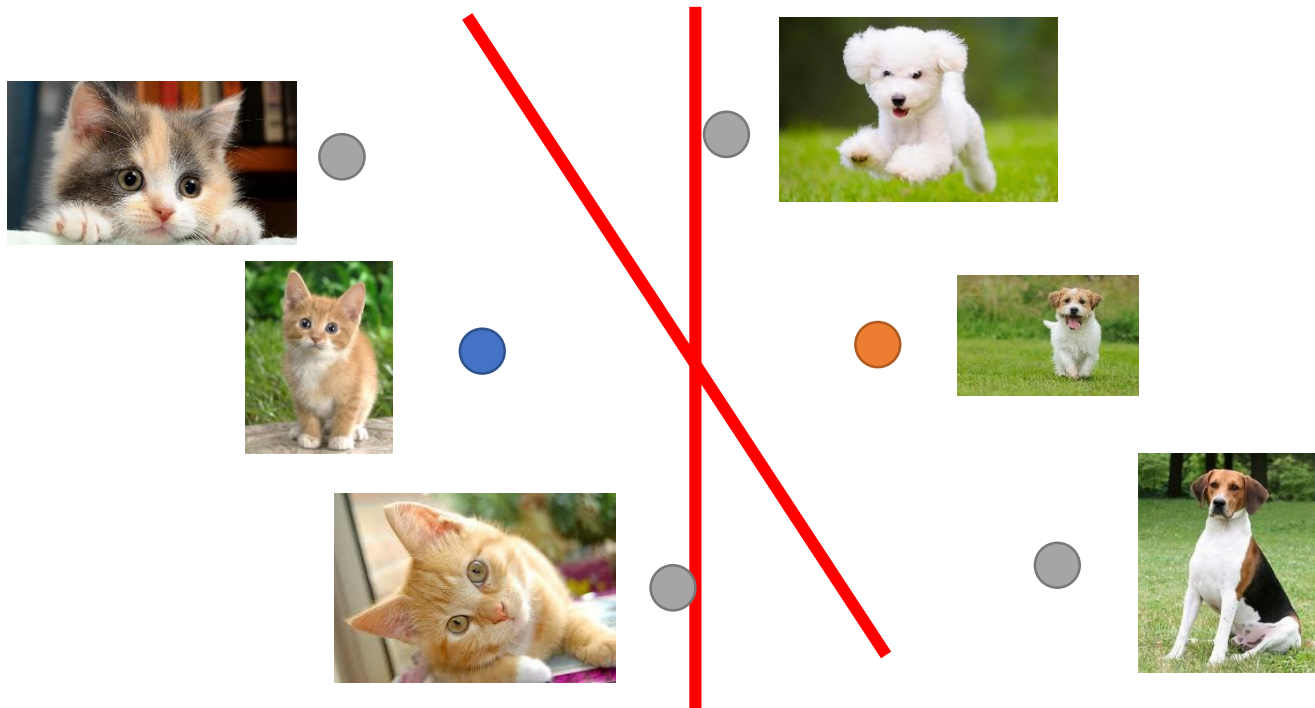
Unlabelled  
data



(Image of cats and dogs without labeling)

# Semi-supervised Learning

- Why semi-supervised learning helps?



The distribution of the unlabeled data tell us something.

# Transfer Learning

Labelled  
data



cat



dog

Labeled  
data



elephant



elephant



tiger

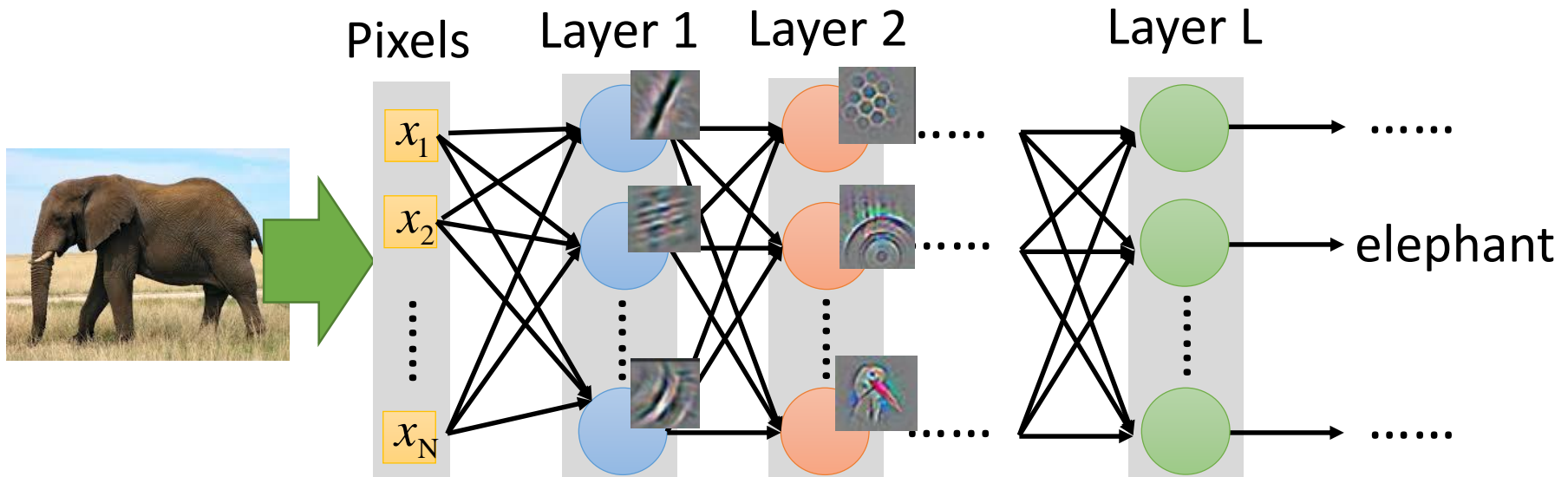


tiger

Not related to the task considered

# Transfer Learning

- Widely used on image processing
  - Using sufficient labeled data to learn a CNN
  - Using this CNN as feature extractor

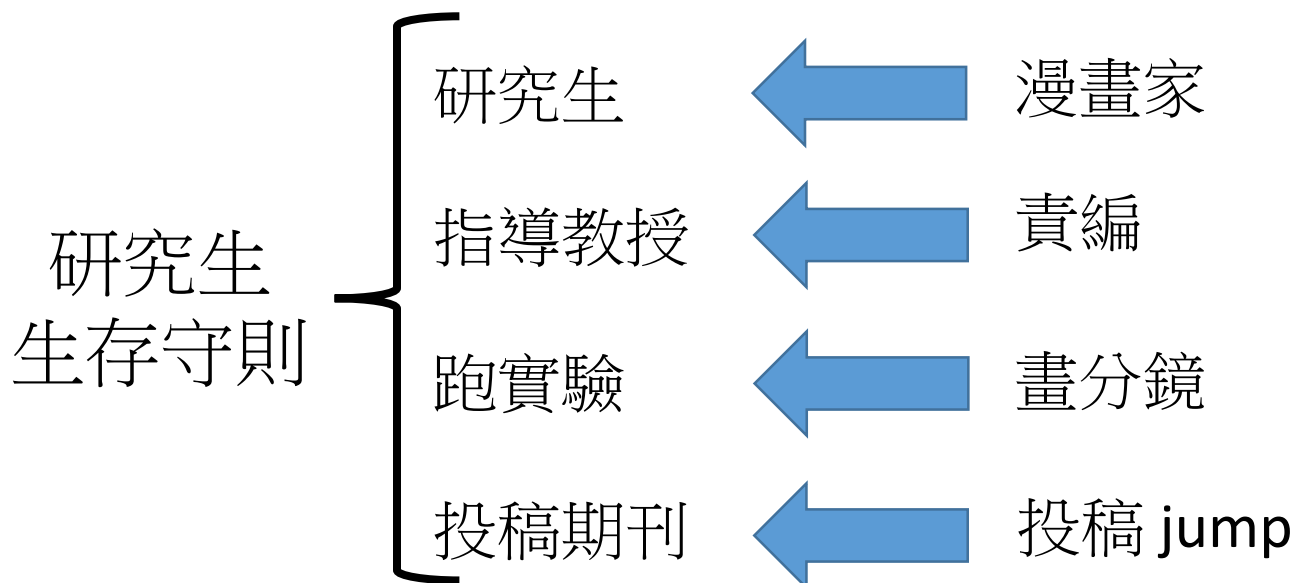


# Transfer Learning

- Example in real life

研究生 online

漫畫家 on-line

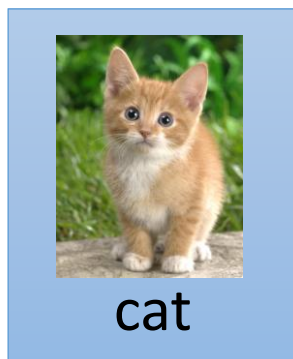


爆漫王

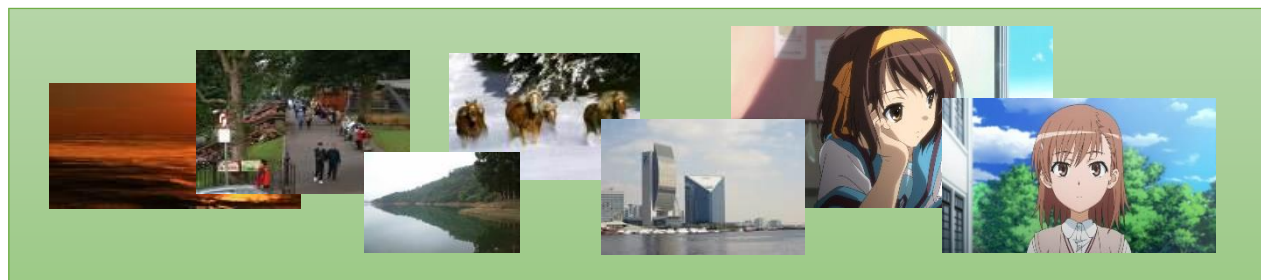
# Self-taught Learning

- Transfer learning with unlabeled data is not related to the task.

Labelled  
data



Unlabeled  
data

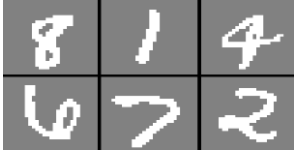


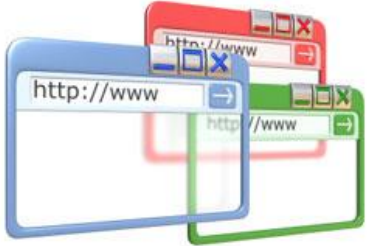
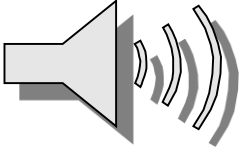



(Just crawl millions of images from the Internet)



# Self-taught Learning

- Sometimes unlabeled data is not related to the task.

	Labelled data	Unlabeled data
Digit Recognition	 Digits	 character
Document Classification	 News	 Webpages
Speech Recognition	 Taiwanese	 English Chinese .....

# Why self-taught learning can work?

- Why Unlabeled and unrelated data can help?
- Find the latent factors controlling the observation

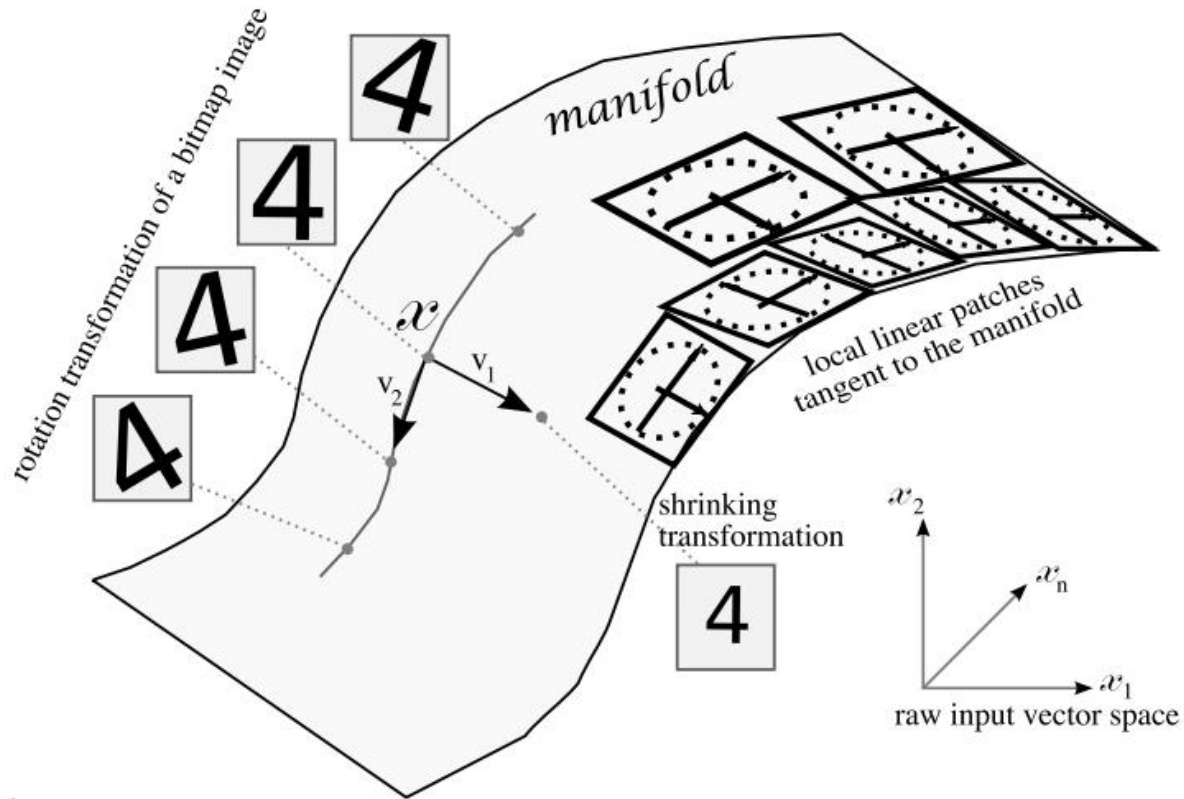
observation

Latent factor



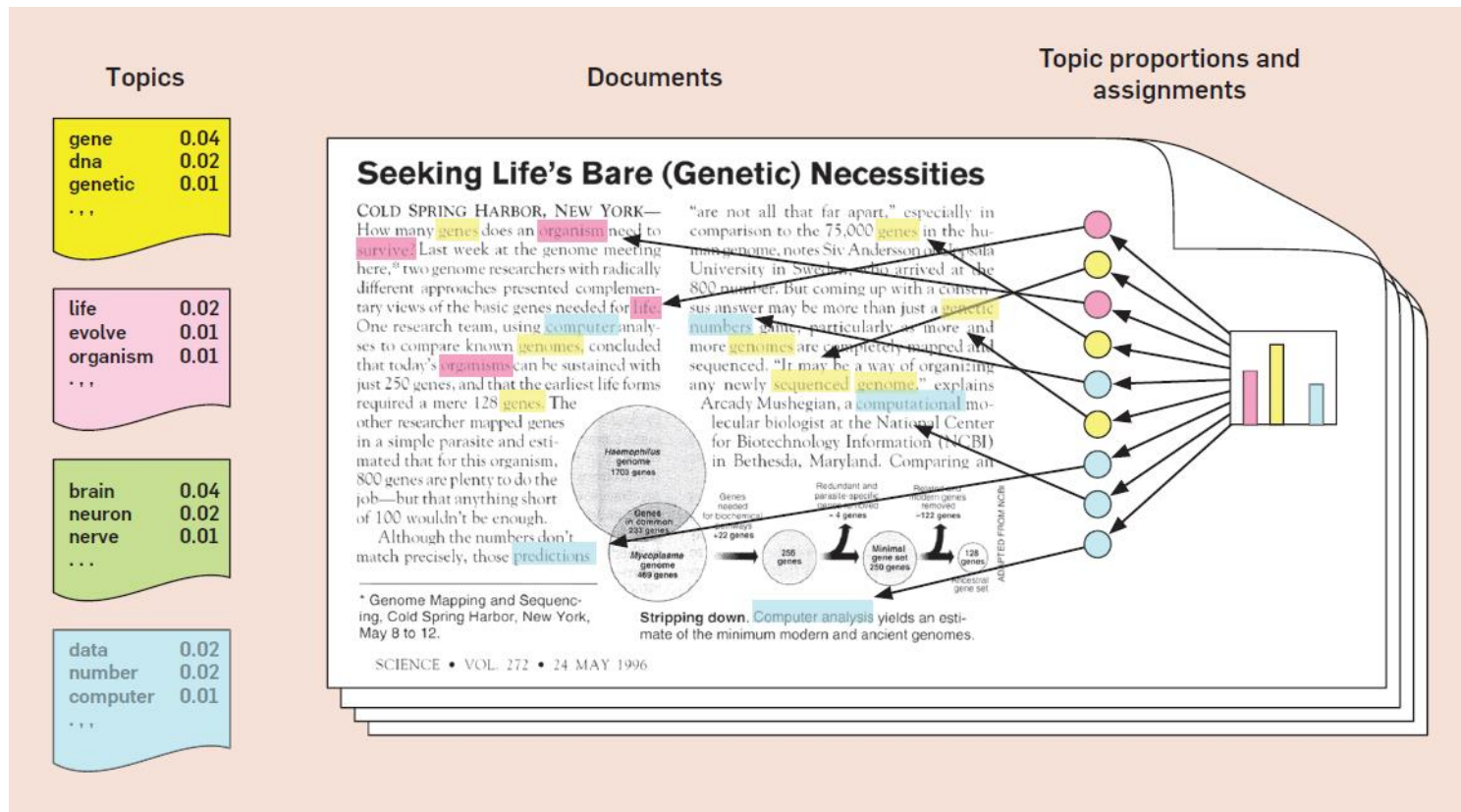
# Latent Factors

- Handwritten Digits



# Latent Factors

- Documents



# Recommendation System

萌傲嬌

A



萌天然呆

B

C



# How to exploit latent factors

- Handwritten Digits



The hand written images are composed of strokes.

## *Strokes (Latent Factors)*



No. 1



No. 2



No. 3



No. 4



No. 5

.....

# How to exploit latent factors

## Strokes (Latent Factors)



No. 1

No. 2

No. 3

No. 4

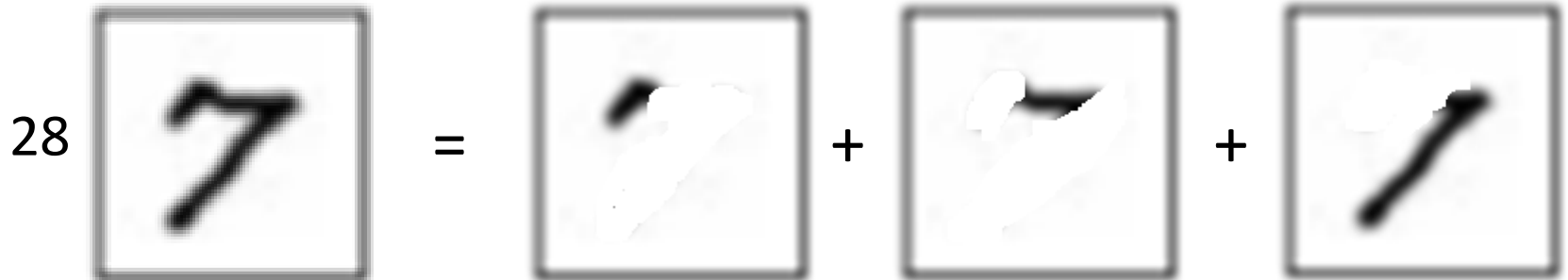
No. 5

28

No. 1

No. 3

No. 5



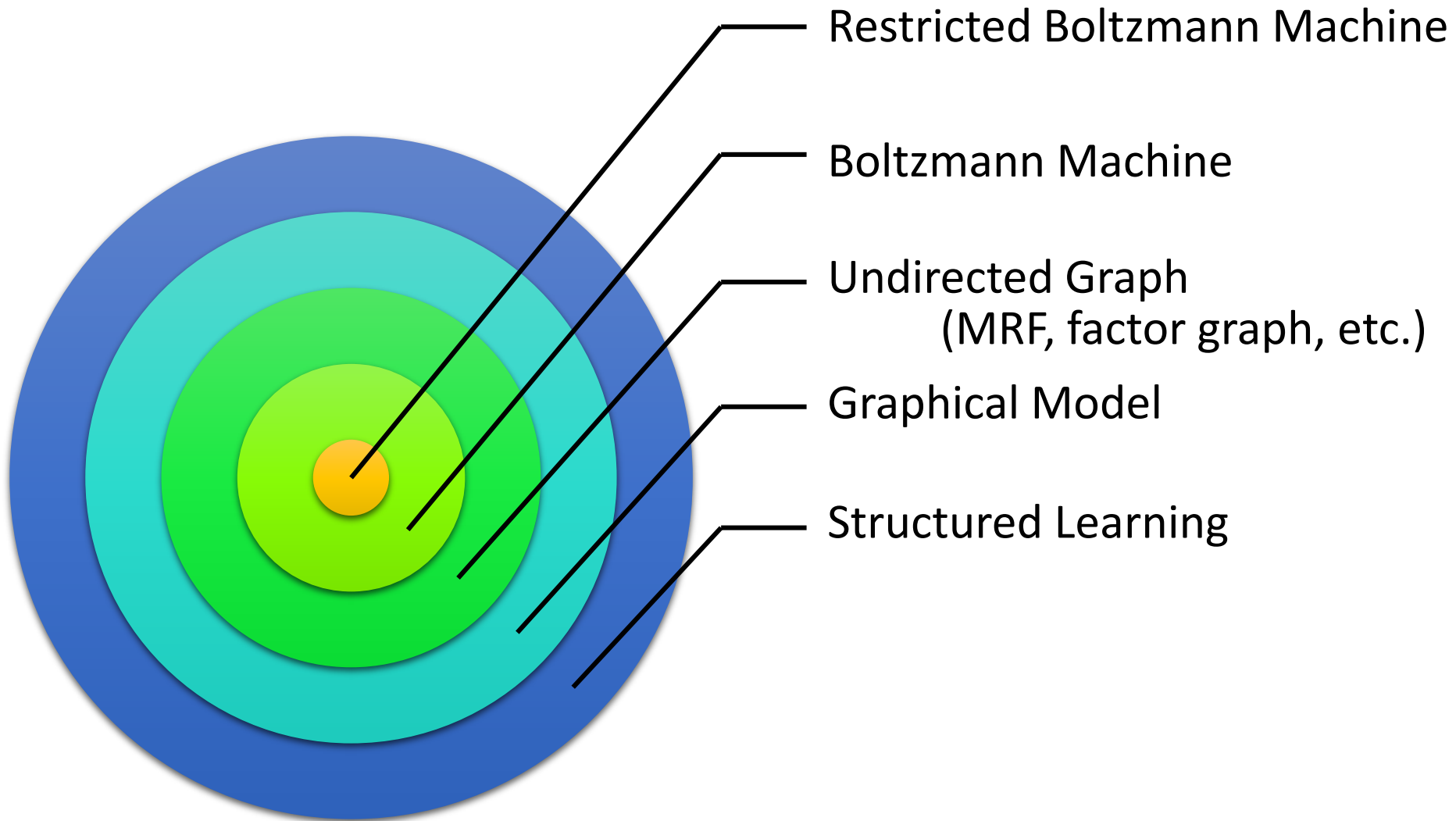
Represented by  
28 X 28 = 784 pixels

[1 0 1 0 1 0 .....]  
(simpler representation)

# Restricted Boltzmann Machine (RBM)



# Where are we?



# Boltzmann Machine

- There is a set of variables  $S = \{s_1, s_2, \dots, s_i, \dots, s_K\}$
- The values of each variable  $s_i \in \{0,1\}$ 
  - will be generalized later

Evaluation function

$$E(S) = \sum_i a_i s_i + \sum_{i < j} w_{ij} s_i s_j$$

When  $s_i = 1$ , the evaluation function gains  $a_i$

When  $s_i = 1$  and  $s_j = 1$ , the evaluation function gains  $w_{ij}$

$a_i$  and  $w_{ij}$  are learned from data.

# Boltzmann Machine - Example

$$s = \{s_1, s_2, s_3, s_4\}$$

$$E(s_1 = 1, s_2 = 0, s_3 = 0, s_4 = 1) = a_1 + a_4 + w_{14}$$

$$E(s_1 = 0, s_2 = 0, s_3 = 0, s_4 = 0) = 0$$

$$E(s_1 = 1, s_2 = 1, s_3 = 1, s_4 = 0) = a_1 + a_2 + a_3 \\ + w_{12} + w_{13} + w_{23}$$

If  $a_i > 0$ ,  $s_i$  is likely to be 1

If  $w_{ij} > 0$ ,  $s_i$  and  $s_j$  are likely to be 1 together

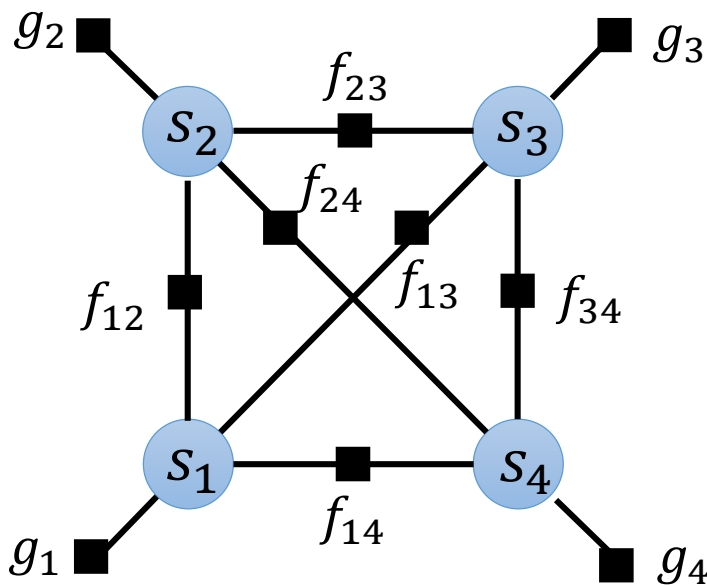
## Probability Point of View

$$P(S) = \frac{e^{E(S)}}{\sum_{S'} e^{E(S)}}$$

# Boltzmann Machine

## - Factor Graph

- There are factors for each node and factors between each node pair



Factor for a node:

$$g_i(s_i) = \begin{cases} a_i & s_i = 1 \\ 0 & \text{else} \end{cases}$$

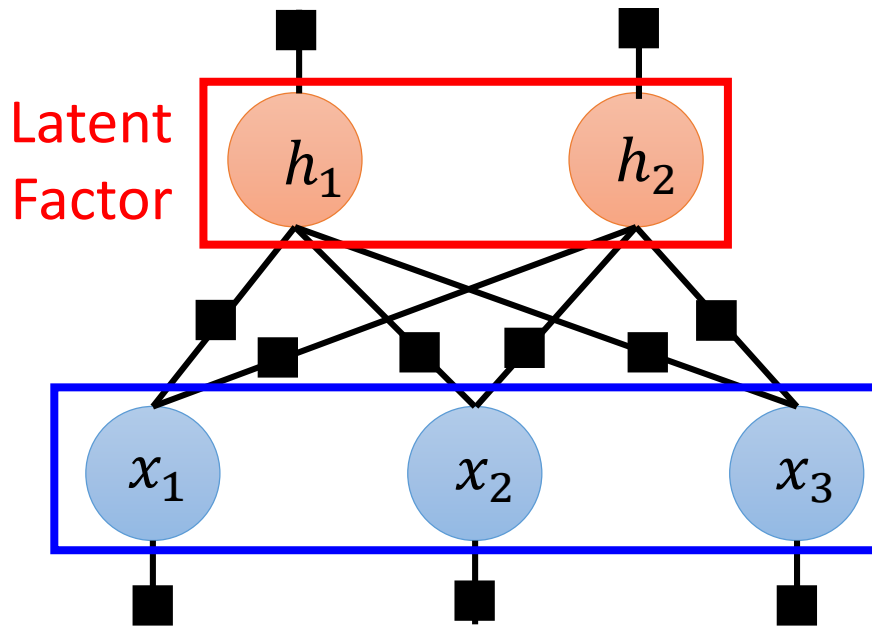
Factor between nodes:

$$f_{ij}(s_i, s_j) = \begin{cases} w_{ij} & s_i = 1, s_j = 1 \\ 0 & \text{else} \end{cases}$$

$$E(S) = \sum_i g_i(s_i) + \sum_{i < j} f_{ij}(s_i, s_j) = \sum_i a_i s_i + \sum_{i < j} w_{ij} s_i s_j$$

# Restricted Boltzmann Machine (RBM)

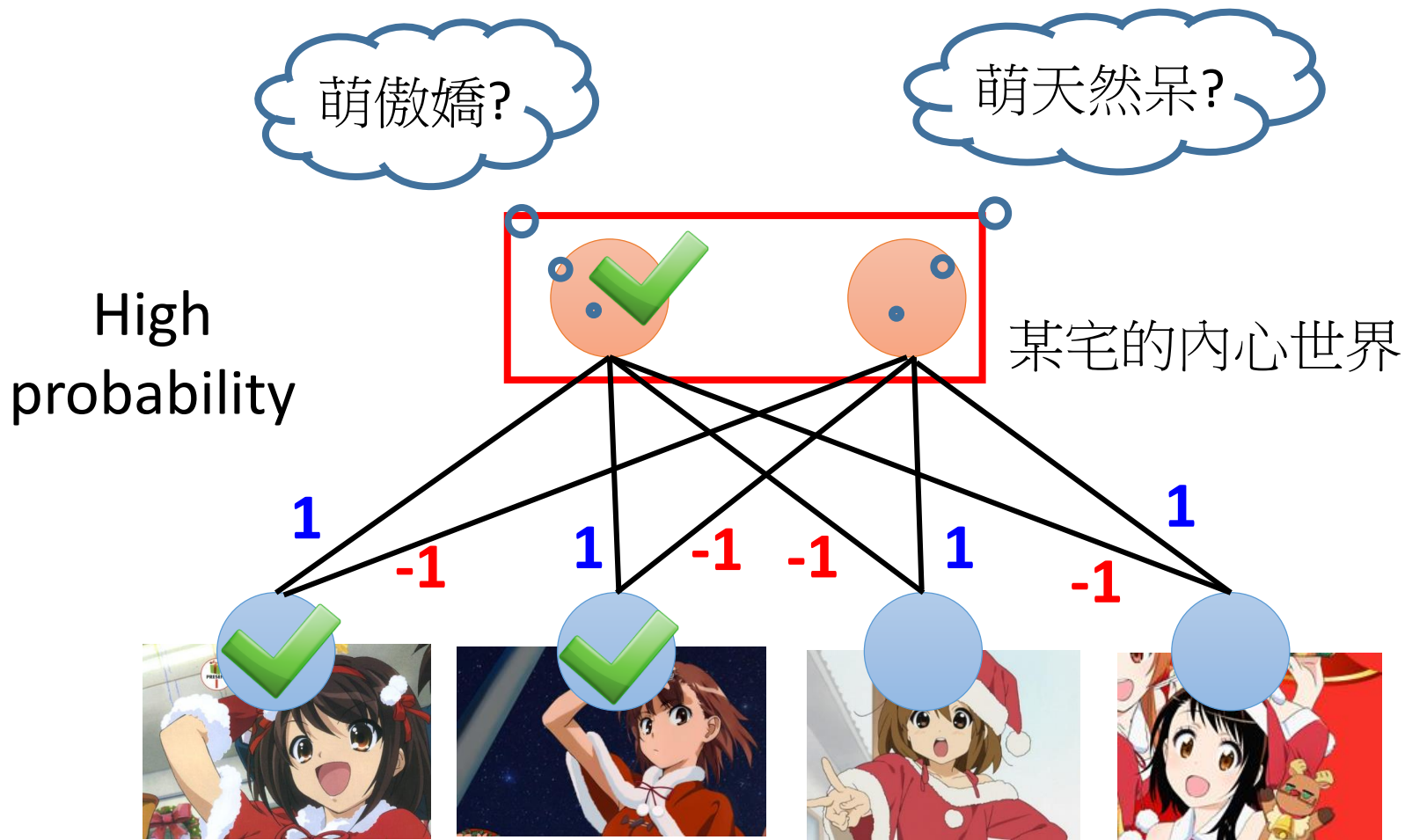
- The variables are separated into two sets
- The variables in the same sets are not connected



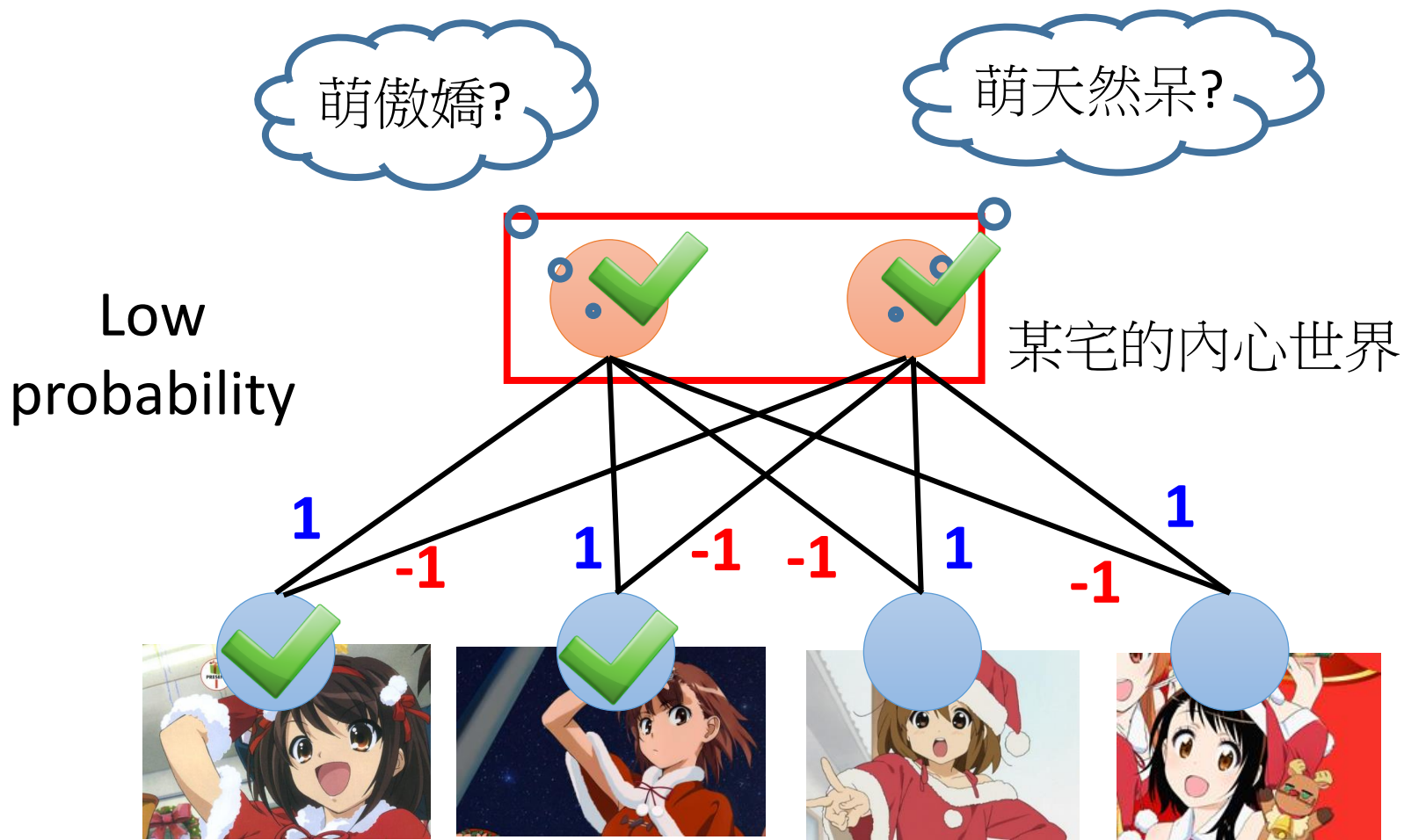
What we can observed.

$$E(x, h) = \sum_{h_i} b_i h_i + \sum_{x_j} c_j x_j + \sum_{h_i, x_j} w_{ij} h_i x_j$$

# RBM – Example

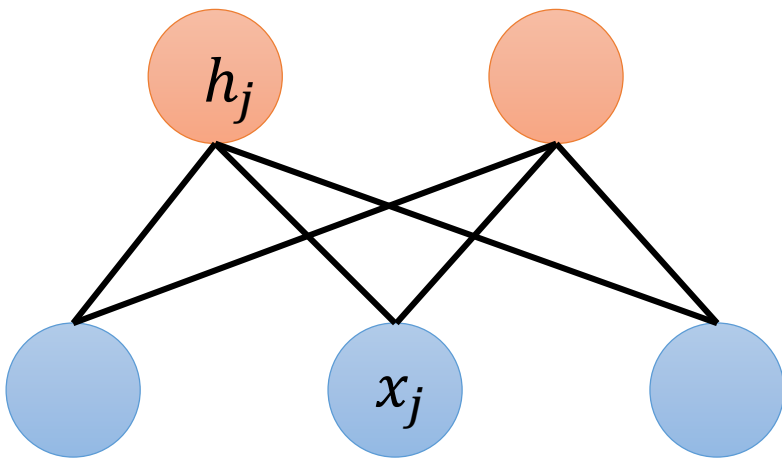


# RBM – Example



# RBM - Inference

- Given the parameters, compute  $P(h_j = 1|x)$  is simple



Given  $x$ , all  $h_j$  are independent

Given  $h$ , all  $x_i$  are independent

(see the reference)



# RBM

$$E(x, h) = \sum_{h_i} b_i h_i + \sum_{x_j} c_j x_j + \sum_{h_i, x_j} w_{ij} h_i x_j$$

- Given the parameters, compute  $P(h_j = 1|x)$  is simple

$$P(h_i = 1|x) = P(h_i = 1|x, h_{-i}) \quad P(x, h) = \frac{e^{E(x, h)}}{\sum_{x', h'} e^{E(x', h')}} \\ = \frac{P(x, h_{-i}, h_i = 1)}{P(x, h_{-i}, h_i = 1) + P(x, h_{-i}, h_i = 0)}$$

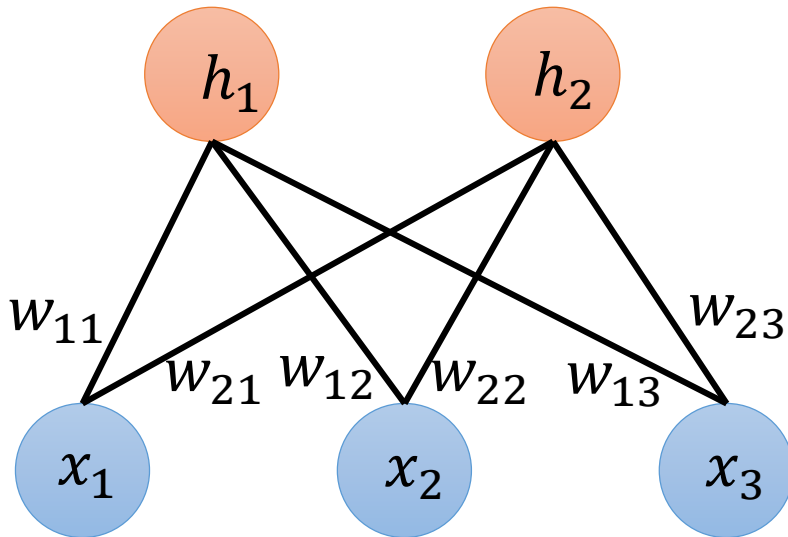
$$= \frac{e^{E(x, h_{-i}, h_i = 1)}}{e^{E(x, h_{-i}, h_i = 1)} + e^{E(x, h_{-i}, h_i = 0)}}$$

Only the terms related to  $h_i$  are different

$$= \frac{e^{b_i + \sum_{x_j} w_{ij} x_j}}{e^{b_i + \sum_{x_j} w_{ij} x_j} + 1} = \frac{1}{1 + e^{-(b_i + \sum_{x_j} w_{ij} x_j)}} = \text{sig} \left( b_i + \sum_{x_j} w_{ij} x_j \right)$$

# RBM - Inference

- Given the parameters, compute  $P(h_j = 1|x)$  is simple



Given  $x_1, x_2$  and  $x_3$

$$P(h_1 = 1|x)$$

$$= \sigma(b_1 + w_{11}x_1 + w_{12}x_2 + w_{13}x_3)$$

$$P(h_2 = 1|x)$$

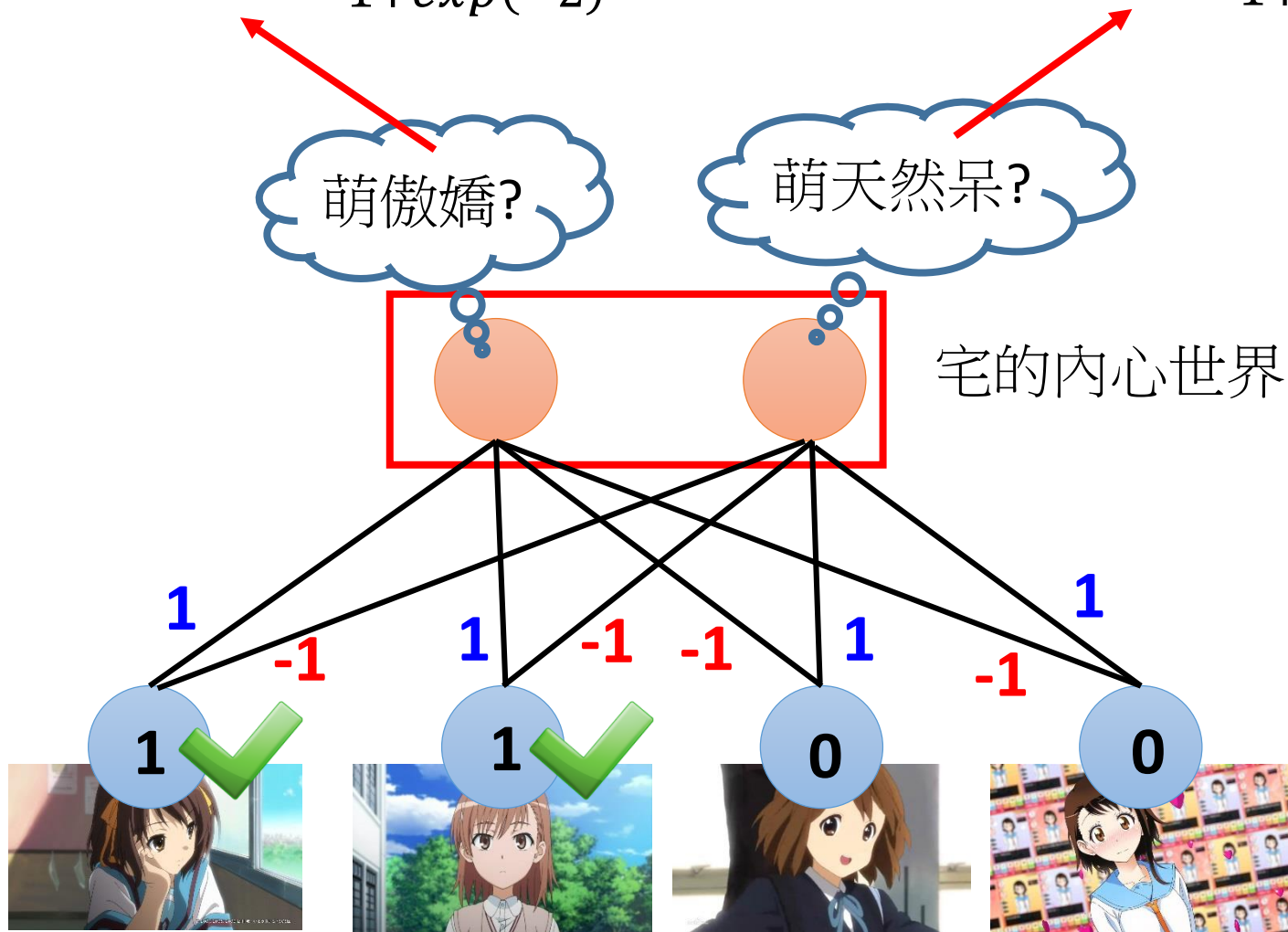
$$= \sigma(b_2 + w_{21}x_1 + w_{22}x_2 + w_{23}x_3)$$

Given  $h_1, h_2$  .....

Neural network with sigmoid function as activation function (誤)

$$P(\text{萌傲娇}=1 | \mathbf{x}) = \frac{1}{1 + \exp(-2)}$$

$$P(\text{萌天然呆}=1 | \mathbf{x}) = \frac{1}{1 + \exp(2)}$$



# RBM – Training without Labeling?

			
$x^1$ ✓	✓		
		✓	✓
$x^2$			
✓	✓	✓	✓
$x^3$			

Training data:  $\{x^1, x^2 \dots x^R\}$

Find the parameters which

maximize  $\prod_{r=1}^R P(x^r)$

Maximizing the  
likelihood of  
observed data

# RBM – Training without Labeling?

- Maximizing  $P(x)$

$$P(x) = \sum_{h''} P(x, h'') = \sum_{h''} \frac{e^{E(x, h'')}}{\sum_{x', h'} e^{E(x', h')}} = \frac{\sum_{h''} e^{E(x, h'')}}{\sum_{x'} \sum_{h'} e^{E(x', h')}}$$

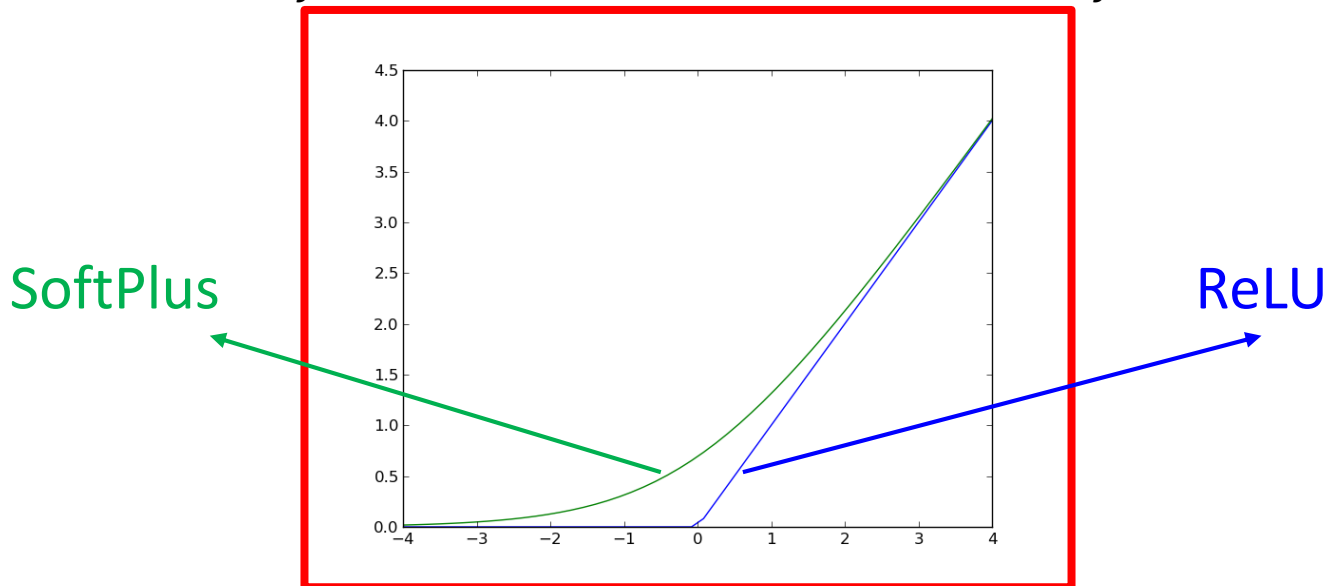
$$F(x) = \sum_{h''} e^{E(x, h'')} \quad \longrightarrow \quad P(x) = \frac{F(x)}{\sum_{x'} F(x')}$$

$$F(x) = \exp \left( \sum_{x_j} c_j x_j + \sum_{h_i} \log \left( 1 + \exp \left( b_i + \sum_{x_j} w_{ij} x_j \right) \right) \right)$$

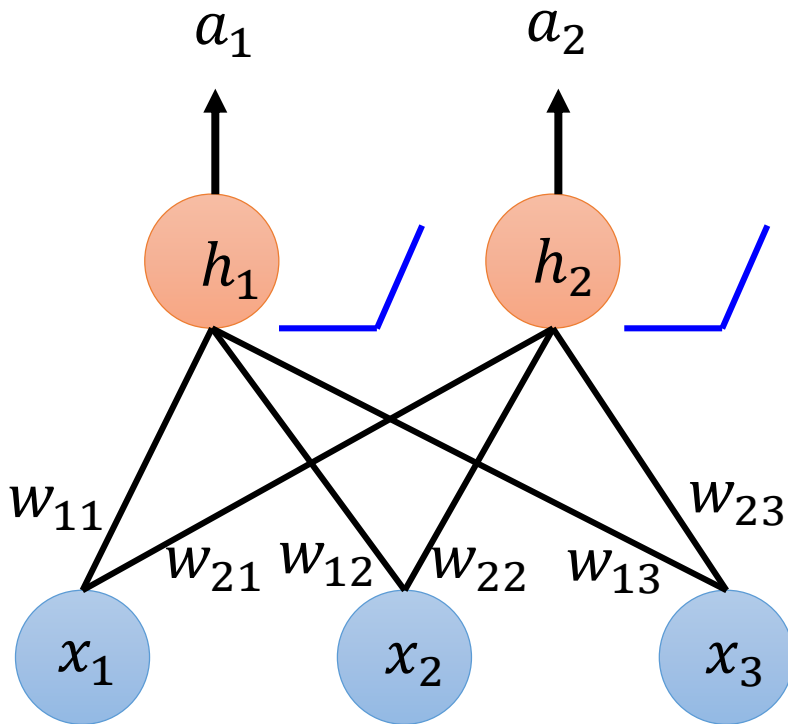
(see the reference)

# RBM – Training without Labeling?

$$F(x) = \exp \left( \sum_{x_j} c_j x_j + \sum_{h_i} \log \left( 1 + \exp \left( b_i + \sum_{x_j} w_{ij} x_j \right) \right) \right)$$
$$= \exp \left( \sum_{x_j} c_j x_j + \sum_{h_i} \underline{\text{softplus}} \left( b_i + \sum_{x_j} w_{ij} x_j \right) \right)$$



# RBM – Training without Labeling?



Neural network with  
Softplus as activation  
function (誤)

$$F(x) = \exp(a_1 + a_2)$$

The degree that the  
hidden layers are  
activated.

$$P(x) = \frac{F(x)}{\sum_{x'} F(x')}$$

→ Increase for x in data  
→ Decrease for any x

# RBM

## – Training by Gradient Ascent

- Given data  $\{x^1, x^2 \dots x^R\}$

$$\text{maximize } \prod_{r=1}^R P(x^r) \quad \longrightarrow \quad \text{maximize } \sum_{r=1}^R \log P(x^r)$$

$$P(x) = \sum_{h''} P(x, h'') = \sum_{h''} \frac{e^{E(x, h'')}}{\sum_{x', h'} e^{E(x', h')}} = \frac{\sum_{h''} e^{E(x, h'')}}{\sum_{x', h'} e^{E(x', h')}}$$

$$\log P(x) = \log \sum_{h''} e^{E(x, h'')} - \log \sum_{x', h'} e^{E(x', h')}$$

$$w_{ij} \leftarrow w_{ij} + \eta \frac{\partial \log P(x)}{\partial w_{ij}} \quad \text{compute } \frac{\partial \log P(x)}{\partial w_{ij}} = ?$$



Warning of Math

$$\log P(x) = \underbrace{\log \sum_{h''} e^{E(x, h'')}}_A - \underbrace{\log \sum_{x', h'} e^{E(x', h')}}_B \quad \frac{\partial \log P(x)}{\partial w_{ij}} = ?$$

$$E(x, h) = \sum_{h_i} b_i h_i + \sum_{x_j} c_j x_j + \sum_{h_i, x_j} w_{ij} h_i x_j \quad \frac{\partial E(x, h)}{\partial w_{ij}} = h_i x_j$$

$$\begin{aligned} \frac{\partial A}{\partial w_{ij}} &= \frac{1}{\sum_{h''} e^{E(x, h'')}} \sum_{h''} \frac{\partial e^{E(x, h'')}}{\partial w_{ij}} \\ &= \frac{1}{\sum_{h''} e^{E(x, h'')}} \sum_{h''} e^{E(x, h'')} \frac{\partial E(x, h'')}{\partial w_{ij}} \\ &= \frac{1}{\sum_{h''} e^{E(x, h'')}} \sum_{h''} e^{E(x, h'')} h''_i x_j = \sum_{h''} \frac{e^{E(x, h'')}}{\sum_{h''} e^{E(x, h'')}} h''_i x_j \\ &= \sum_{h''} P(h'' | x) h''_i x_j \end{aligned}$$

$P(h'' | x)$

$$\log P(x) = \underbrace{\log \sum_{h''} e^{E(x, h'')}}_A - \underbrace{\log \sum_{x', h'} e^{E(x', h')}}_B \quad \frac{\partial \log P(x)}{\partial w_{ij}} = ?$$

$$\begin{aligned} \frac{\partial B}{\partial w_{ij}} &= \frac{1}{\sum_{x', h'} e^{E(x', h')}} \sum_{x', h'} \frac{\partial e^{E(x', h')}}{\partial w_{ij}} \\ &= \frac{1}{\sum_{x', h'} e^{E(x', h')}} \sum_{x', h'} e^{E(x', h')} \frac{\partial E(x', h')}{\partial w_{ij}} \\ &= \frac{1}{\sum_{x', h'} e^{E(x', h')}} \sum_{x', h'} e^{E(x', h')} h'_i x'_j \\ &= \sum_{x', h'} \frac{e^{E(x', h')}}{\sum_{x', h'} e^{E(x', h')}} h'_i x'_j = \sum_{x', h'} P(x', h') h'_i x'_j \end{aligned}$$

End of Warning

$$\log P(x) = \underbrace{\log \sum_{h''} e^{E(x, h'')}}_A - \underbrace{\log \sum_{x', h'} e^{E(x', h')}}_B \quad \frac{\partial \log P(x)}{\partial w_{ij}} = ?$$

$$\frac{\partial A}{\partial w_{ij}} = \sum_{h''} P(h'' | x) h''_i x_j = \underline{E_{P(h'' | x)}[h''_i x_j]}$$

The expected value of  $h''_i x_j$  based on  $P(h'' | x)$  given the current parameters

$$\log P(x) = \underbrace{\log \sum_{h''} e^{E(x, h'')}}_A - \underbrace{\log \sum_{x', h'} e^{E(x', h')}}_B \quad \frac{\partial \log P(x)}{\partial w_{ij}} = ?$$

$$\frac{\partial A}{\partial w_{ij}} = \sum_{h''} P(h'' | x) h''_i x_j = \underline{E_{P(h'' | x)}[h''_i x_j]}$$

The expected value of  $h''_i x_j$  based on  $P(h'' | x)$  given the current parameters

$$\frac{\partial B}{\partial w_{ij}} = \sum_{x', h'} P(x', h') h'_i x'_j = \underline{E_{P(x', h')}[h'_i x'_j]}$$

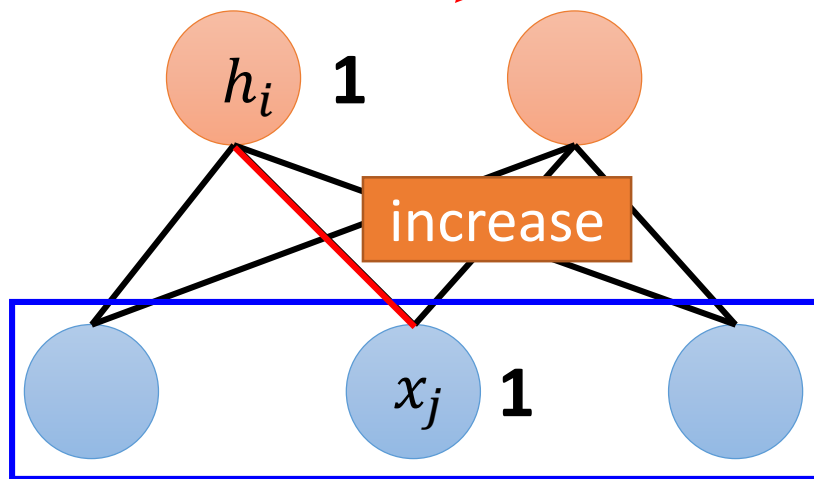
The expected value of  $h'_i x'_j$  based on  $P(x', h')$  given the current parameters

# RBM

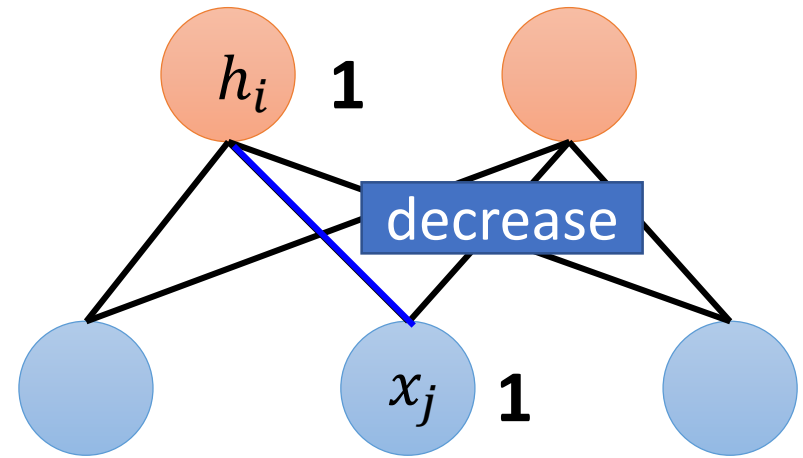
## – Training by Gradient Ascent

$$w_{ij} \leftarrow w_{ij} + \eta \frac{\partial \log P(x)}{\partial w_{ij}}$$

$$\frac{\partial \log P(x)}{\partial w_{ij}} = \underbrace{E_{P(h''|x)}[h''_i x_j]}_{\text{red arrow}} - \underbrace{E_{P(x',h')}[h'_i x'_j]}_{\text{blue arrow}}$$



Given  $x$

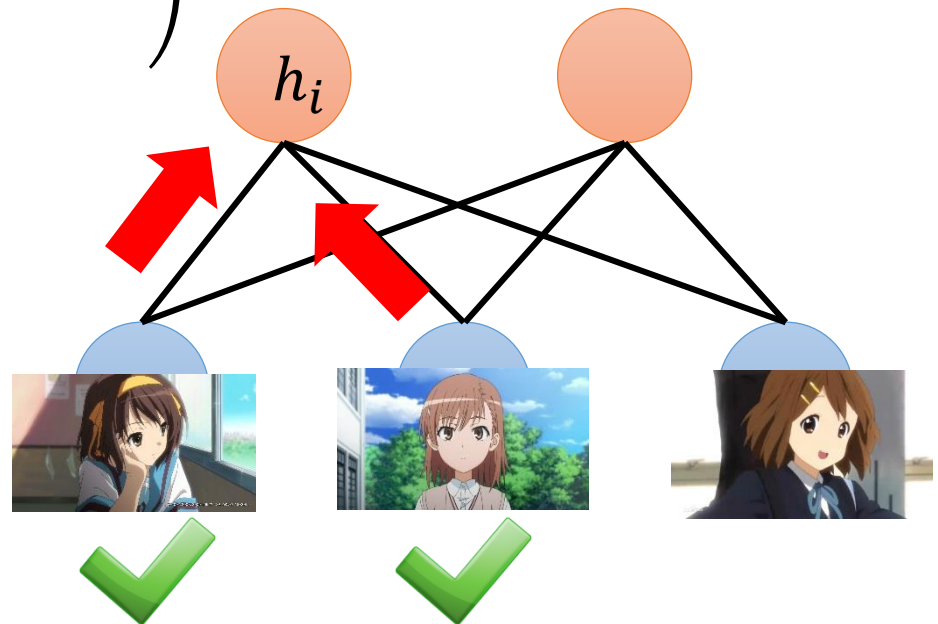


$$\log P(x) = \underbrace{\log \sum_{h''} e^{E(x, h'')}}_A - \underbrace{\log \sum_{x', h'} e^{E(x', h')}}_B \quad \frac{\partial \log P(x)}{\partial w_{ij}} = ?$$

$$\frac{\partial A}{\partial w_{ij}} = \sum_{h''} P(h'' | x) h''_i x_j = \underbrace{E_{P(h'' | x)} [h''_i x_j]}_{\text{Just a sigmoid}}$$

$$= \underline{P(h''_i = 1 | x)} x_j \quad \text{Just a sigmoid}$$

$$P(h''_i = 1 | x) = \text{sig} \left( b_i + \sum_{x_j} w_{ij} x_j \right)$$





$$\log P(x) = \underbrace{\log \sum_{h''} e^{E(x, h'')}}_A - \underbrace{\log \sum_{x', h'} e^{E(x', h')}}_B \quad \frac{\partial \log P(x)}{\partial w_{ij}} = ?$$

$$\frac{\partial A}{\partial w_{ij}} = \sum_{h''} P(h''|x) h''_i x_j = \underbrace{E_{P(h''|x)}[h''_i x_j]}_{= P(h''_i = 1|x) x_j} \quad \text{Just a sigmoid}$$

$$\frac{\partial B}{\partial w_{ij}} = \sum_{x', h'} P(x', h') h'_i x'_j = \underbrace{E_{P(x', h')}[h'_i x'_j]}_{\text{Exact computing is not tractable}}$$

Sample by  
Gibbs  
sampling

$$x^1, h^1 \sim P(x', h')$$

$$x^2, h^2 \sim P(x', h')$$

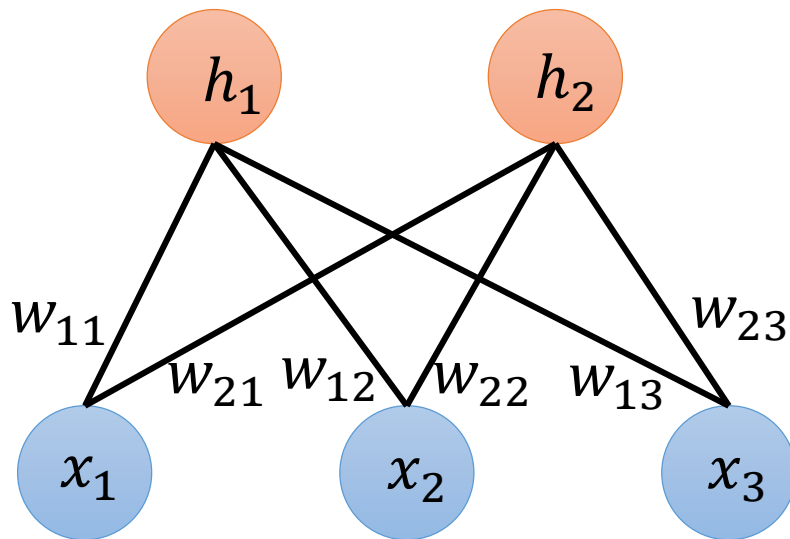
⋮

$$x^N, h^N \sim P(x', h')$$

$$\frac{1}{N} \sum_{n=1}^N h_i^n x_j^n$$

# RBM –Gibbs Sampling

Use Gibbs sampling  
to sample from  $P(x,h)$



Random initialize  $x^0, h^0$

For  $n = 1$  to  $N$

$$x_1^n \sim P(x_1 | \cancel{x_2^{n-1}}, \cancel{x_3^{n-1}}, h_1^{n-1}, h_2^{n-1})$$

$$x_2^n \sim P(x_2 | \cancel{x_1^n}, \cancel{x_3^{n-1}}, h_1^{n-1}, h_2^{n-1})$$

$$x_3^n \sim P(x_3 | \cancel{x_1^n}, \cancel{x_2^n}, h_1^{n-1}, h_2^{n-1})$$

$$h_1^n \sim P(h_1 | x_1^n, x_2^n, x_3^n, \cancel{h_2^{n-1}})$$

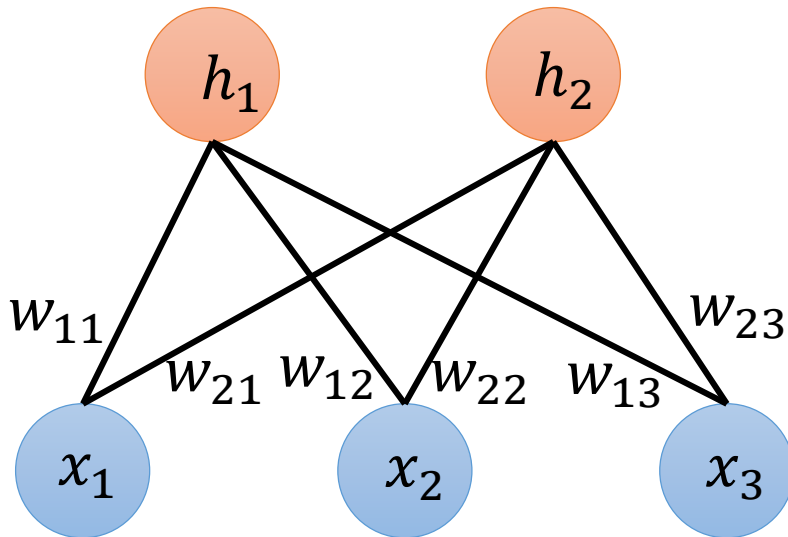
$$h_2^n \sim P(h_2 | x_1^n, x_2^n, x_3^n, \cancel{h_1^n})$$

Obtain one sample  $x^n, h^n$

( as sample from  $P(x,h)$  )

# RBM –Gibbs Sampling

Use Gibbs sampling  
to sample from  $P(x,h)$



Random initialize  $x^0, h^0$

For  $n = 1$  to  $N$

$$x_1^n \sim P(x_1 | h_1^{n-1}, h_2^{n-1})$$

$$x_2^n \sim P(x_2 | h_1^{n-1}, h_2^{n-1})$$

$$x_3^n \sim P(x_3 | h_1^{n-1}, h_2^{n-1}) \quad x^n \sim P(x | h^{n-1})$$

$$h_1^n \sim P(h_1 | x_1^n, x_2^n, x_3^n)$$

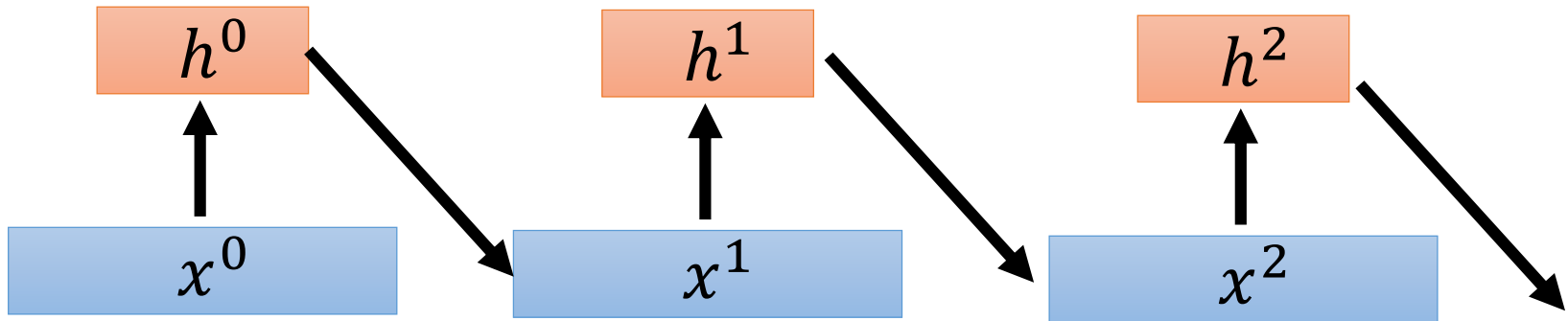
$$h_2^n \sim P(h_2 | x_1^n, x_2^n, x_3^n) \quad h^n \sim P(h | x^n)$$

Obtain one sample  $x^n, h^n$

( as sample from  $P(x,h)$  )

# RBM –Gibbs Sampling

$$x^n \sim P(x|h^{n-1})$$
$$h^n \sim P(h|x^n)$$

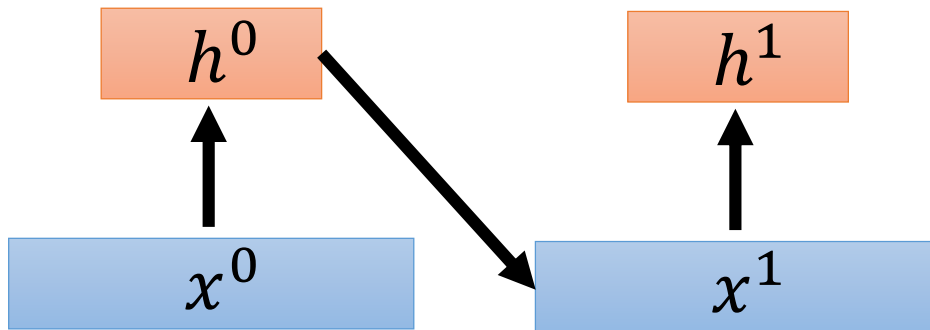


$x^0$  is data  $x$

$$E_{P(x',h')} [h'_i x'_j] \approx \frac{1}{N} \sum_{n=1}^N h_i^n x_j^n$$

Each time we update parameters, we should do Gibbs sampling?!

# RBM – Contrastive Divergence (CD)



$x^0$  is data  $x$

Stop!

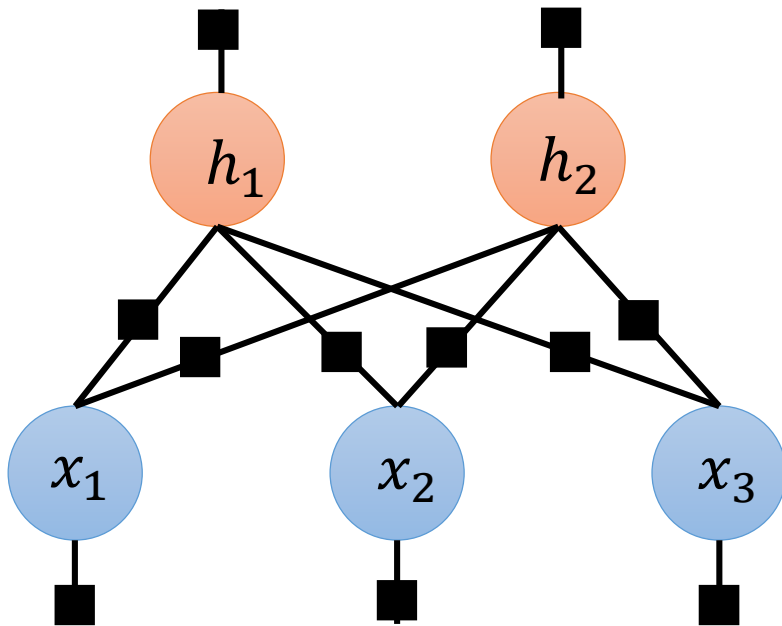
$$E_{P(x', h')} [h'_i x'_j] \approx h_i^1 x_j^1$$

It works in reality!

## Persistent CD

Ref: Tieleman, Tijmen. "Training restricted Boltzmann machines using approximations to the likelihood gradient." *Proceedings of the 25th international conference on Machine learning*. ACM, 2008.

# RBM - Generalization



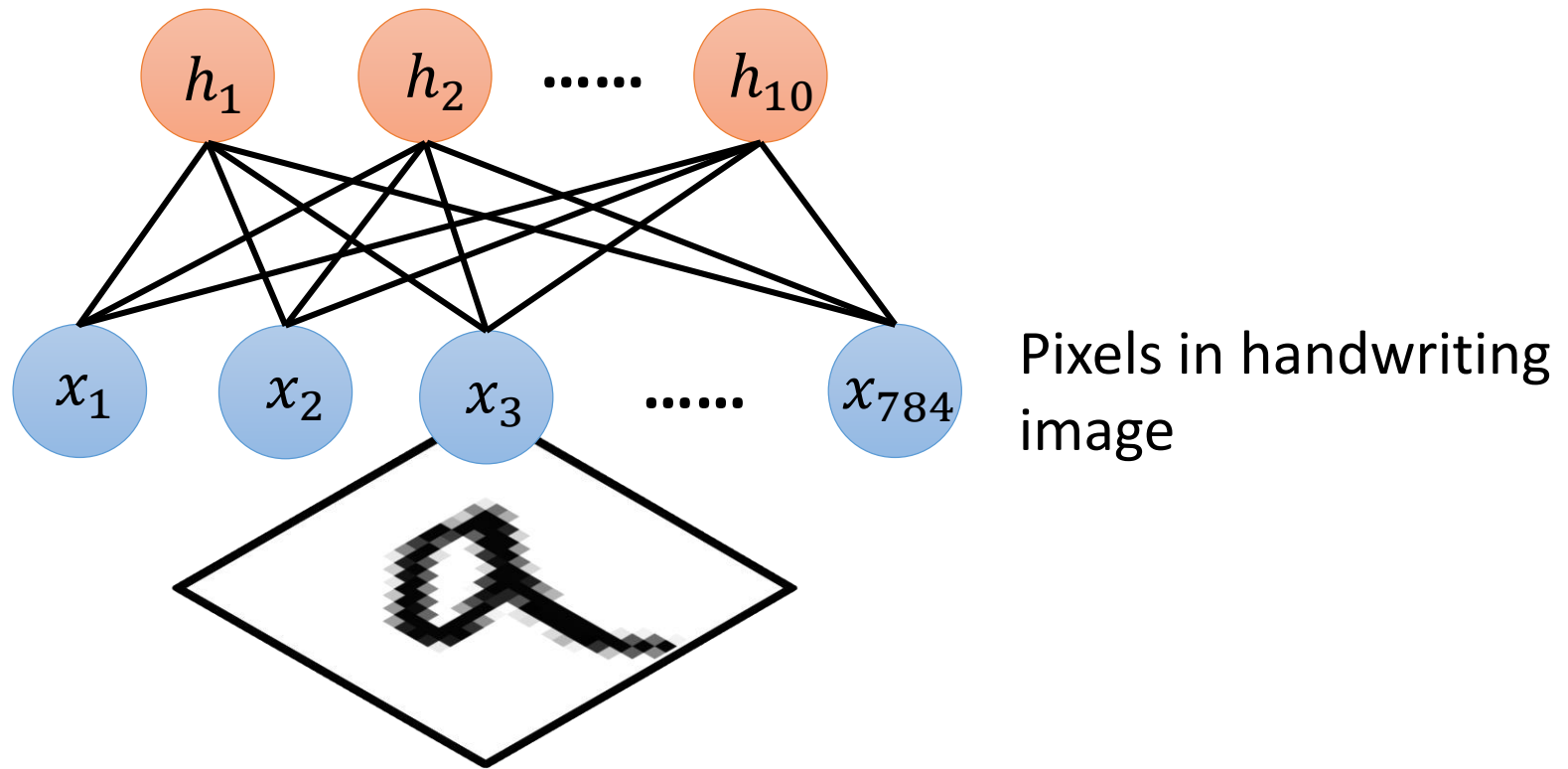
$$g_j(x_j) = \begin{cases} c_j & x_j = 1 \\ 0 & \text{else} \end{cases}$$

If  $x$  are real numbers

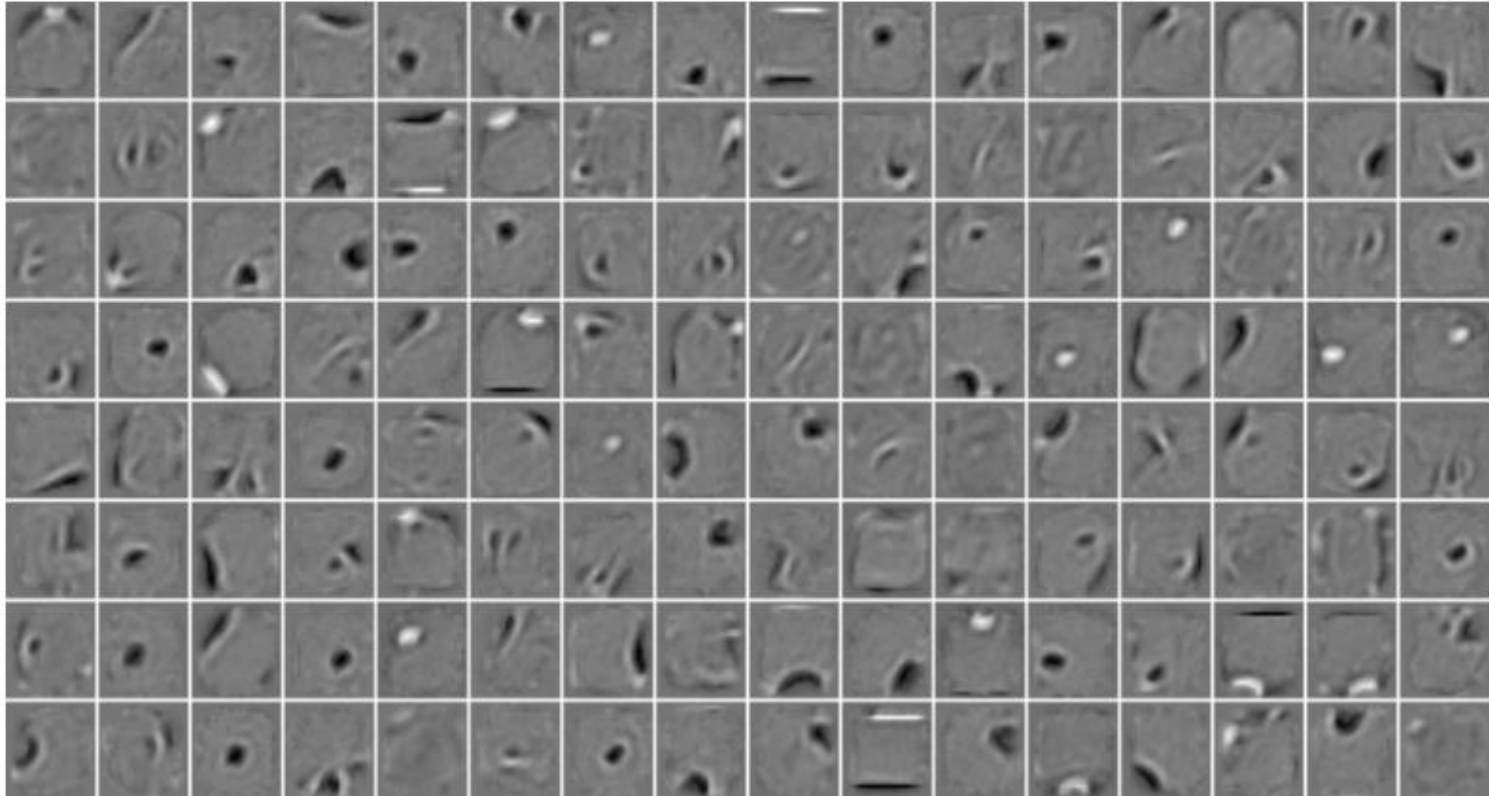
$$g_j(x_j) = -(x_j - c_j)^2$$

$$E(x, h) = \sum_{h_i} b_i h_i + \sum_{h_i, x_j} w_{ij} h_i x_j - \sum_{x_j} (x_j - c_j)^2$$

# RBM - Handwritten Digits



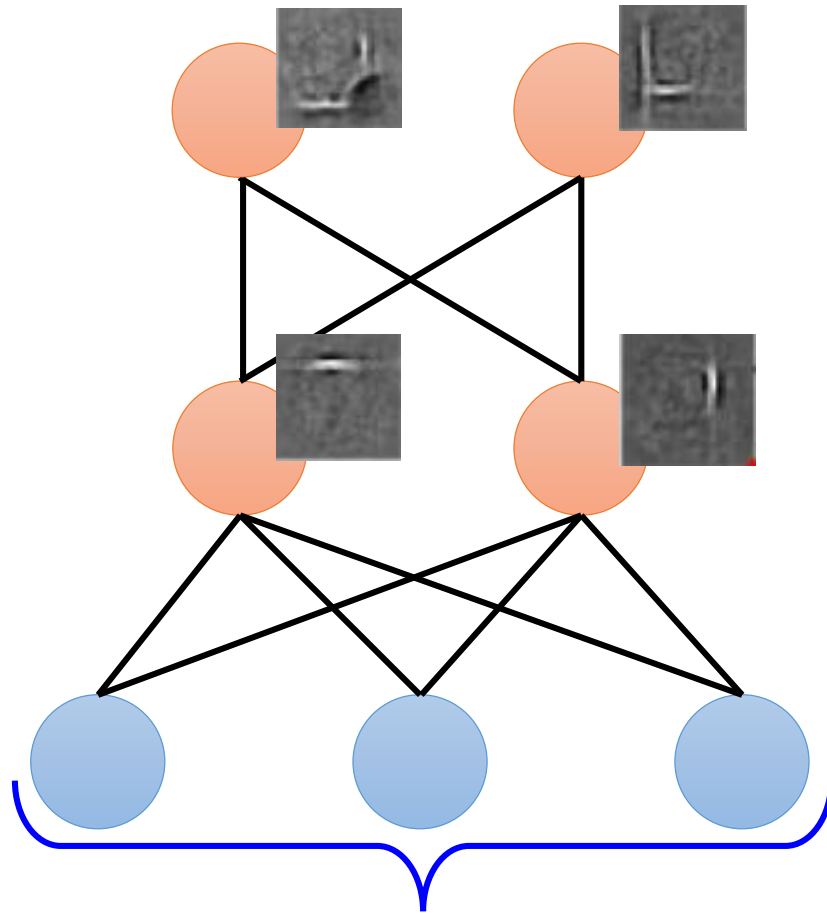
# RBM – Handwritten Digits



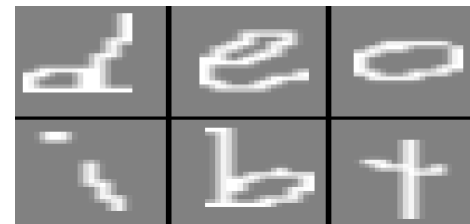
Source of image: Larochelle, H., Bengio, Y., Louradour, J., & Lamblin, P. (2009). Exploring strategies for training deep neural networks. *The Journal of Machine Learning Research*, 10, 1-40.



# Deep Boltzmann Machines



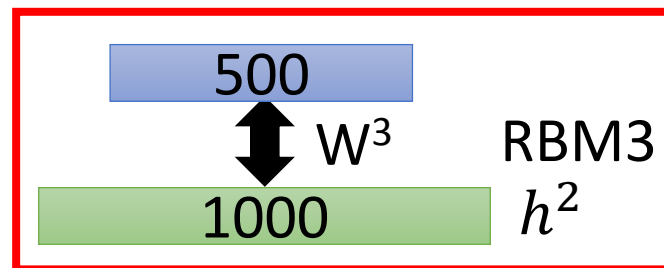
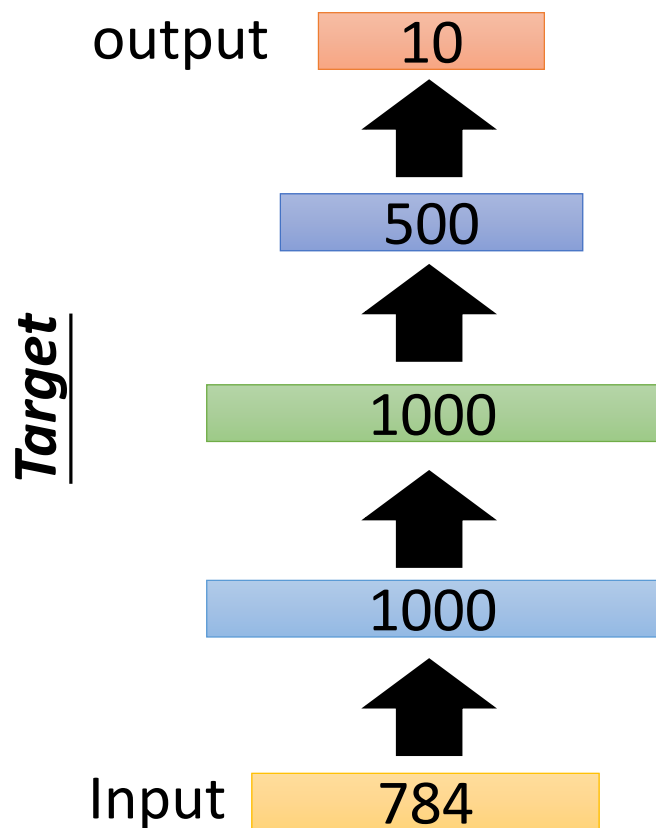
Ref: Salakhutdinov, Ruslan, and Geoffrey E. Hinton. "Deep Boltzmann machines." *International Conference on Artificial Intelligence and Statistics*. 2009.



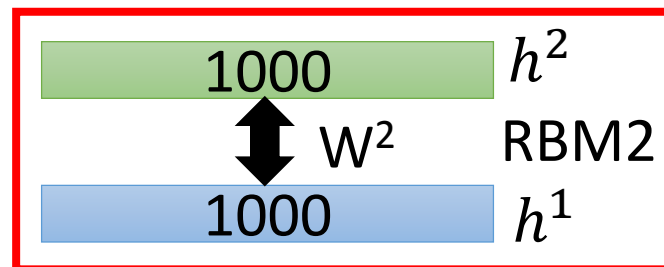
Pixels

# RBM - Pre-training DNN

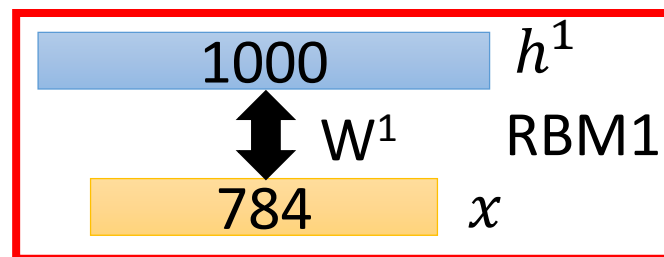
## Greedy Layer-wise Pre-training



Sample from  $P_{RBM2}(h^2|h^1)$

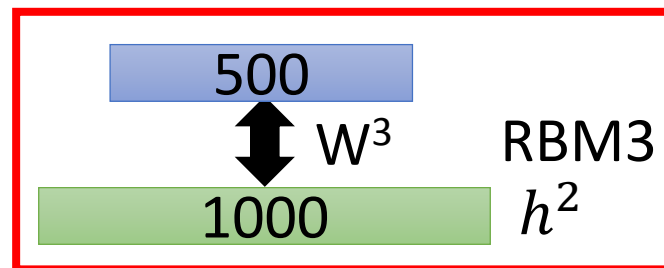
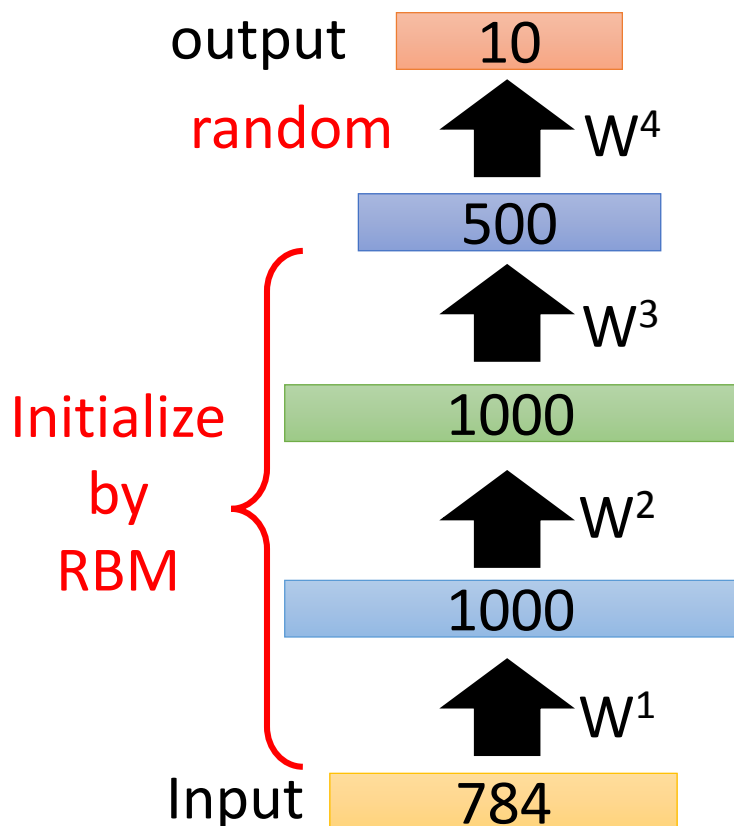


Sample from  $P_{RBM1}(h^1|x)$

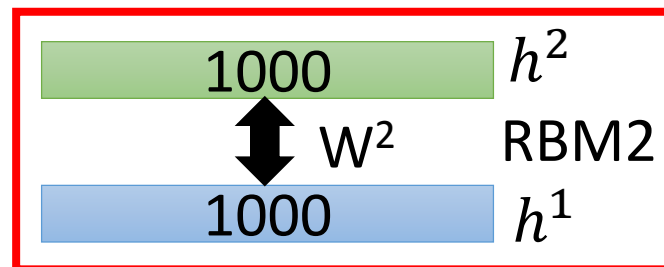


# RBM - Pre-training DNN

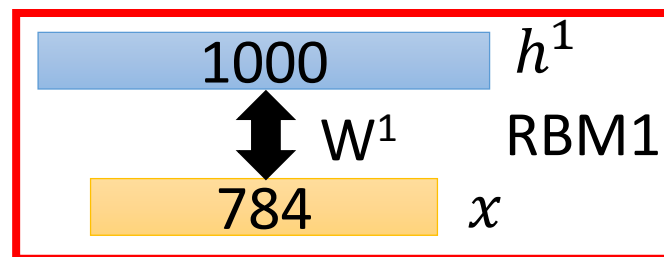
## Greedy Layer-wise Pre-training



Sample from  $P_{RBM2}(h^2|h^1)$



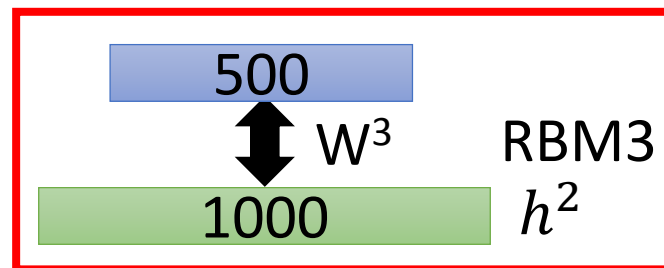
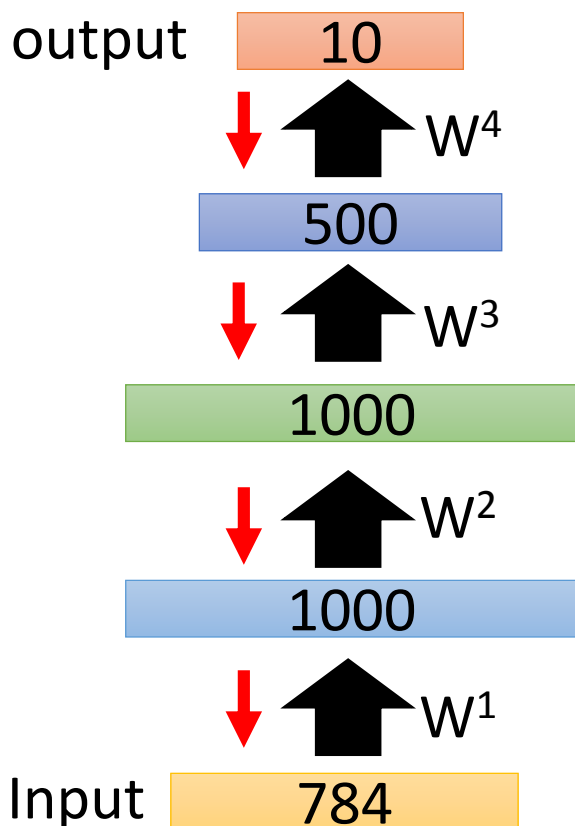
Sample from  $P_{RBM1}(h^1|x)$



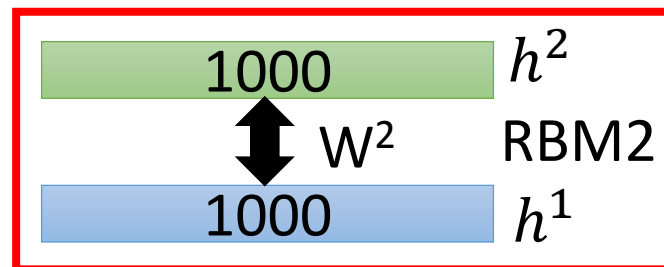
# RBM - Pre-training DNN

Then do back propagation

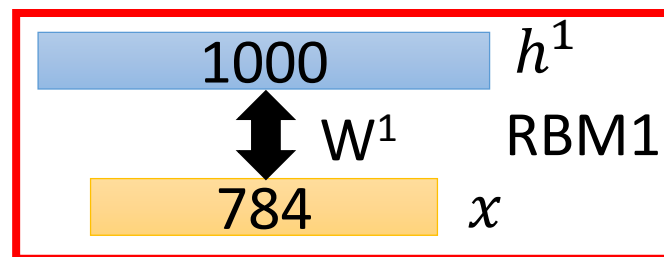
➔ **Fine tuning**



Sample from  $P_{RBM2}(h^2|h^1)$



Sample from  $P_{RBM1}(h^1|x)$

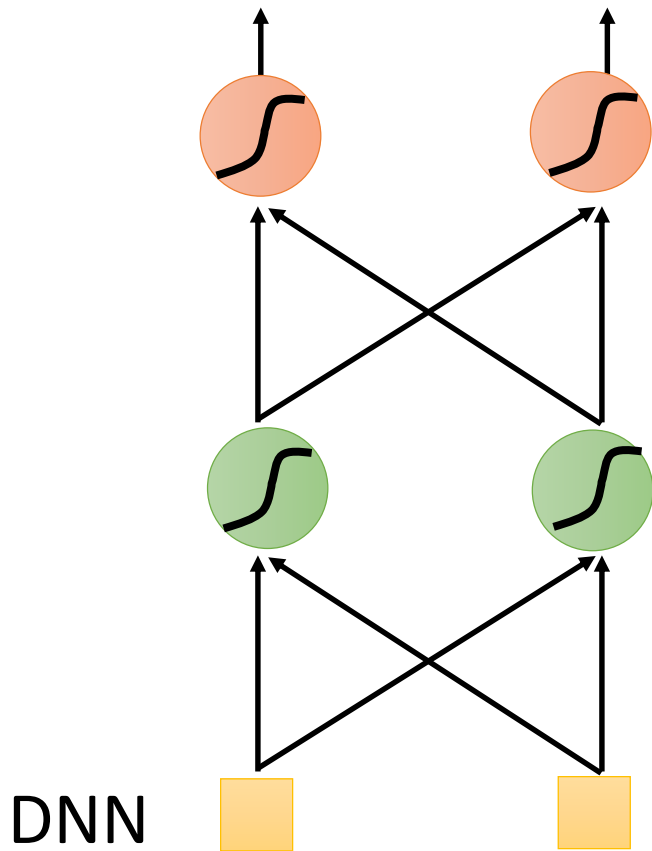


# Reference for RBM

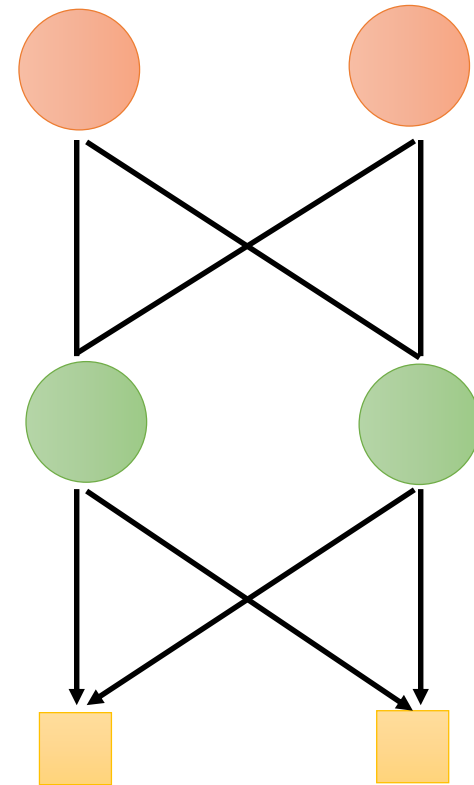
- Independent
  - Neural networks [5.2] : Restricted Boltzmann machine – inference
  - [https://www.youtube.com/watch?v=lekCh\\_i32iE&list=PL6Xpj9I5qXYEcOhn7TqghAJ6NAPrNmUBH&index=37](https://www.youtube.com/watch?v=lekCh_i32iE&list=PL6Xpj9I5qXYEcOhn7TqghAJ6NAPrNmUBH&index=37)
- Intuition for maximizing likelihood
  - Neural networks [5.3] : Restricted Boltzmann machine - free energy
  - [https://www.youtube.com/watch?v=e0Ts\\_7Y6hZU&list=PL6Xpj9I5qXYEcOhn7TqghAJ6NAPrNmUBH&index=38](https://www.youtube.com/watch?v=e0Ts_7Y6hZU&list=PL6Xpj9I5qXYEcOhn7TqghAJ6NAPrNmUBH&index=38)

# Deep Belief Network (DBN)

- DBN  $\neq$  DNN



DBN  
(Graphical Model)



# Reference for DBN

- Neural networks [7.7] : Deep learning - deep belief network
  - <https://www.youtube.com/watch?v=vkb6AWYXZ5I&list=PL6Xpj9I5qXYEcOhn7TqghAJ6NAPrNmUBH&index=57>
- Neural networks [7.8] : Deep learning - variational bound
  - <https://www.youtube.com/watch?v=pStDscJh2Wo&list=PL6Xpj9I5qXYEcOhn7TqghAJ6NAPrNmUBH&index=58>
- Neural networks [7.9] : Deep learning - DBN pre-training
  - <https://www.youtube.com/watch?v=35MUIYCColk&list=PL6Xpj9I5qXYEcOhn7TqghAJ6NAPrNmUBH&index=59>

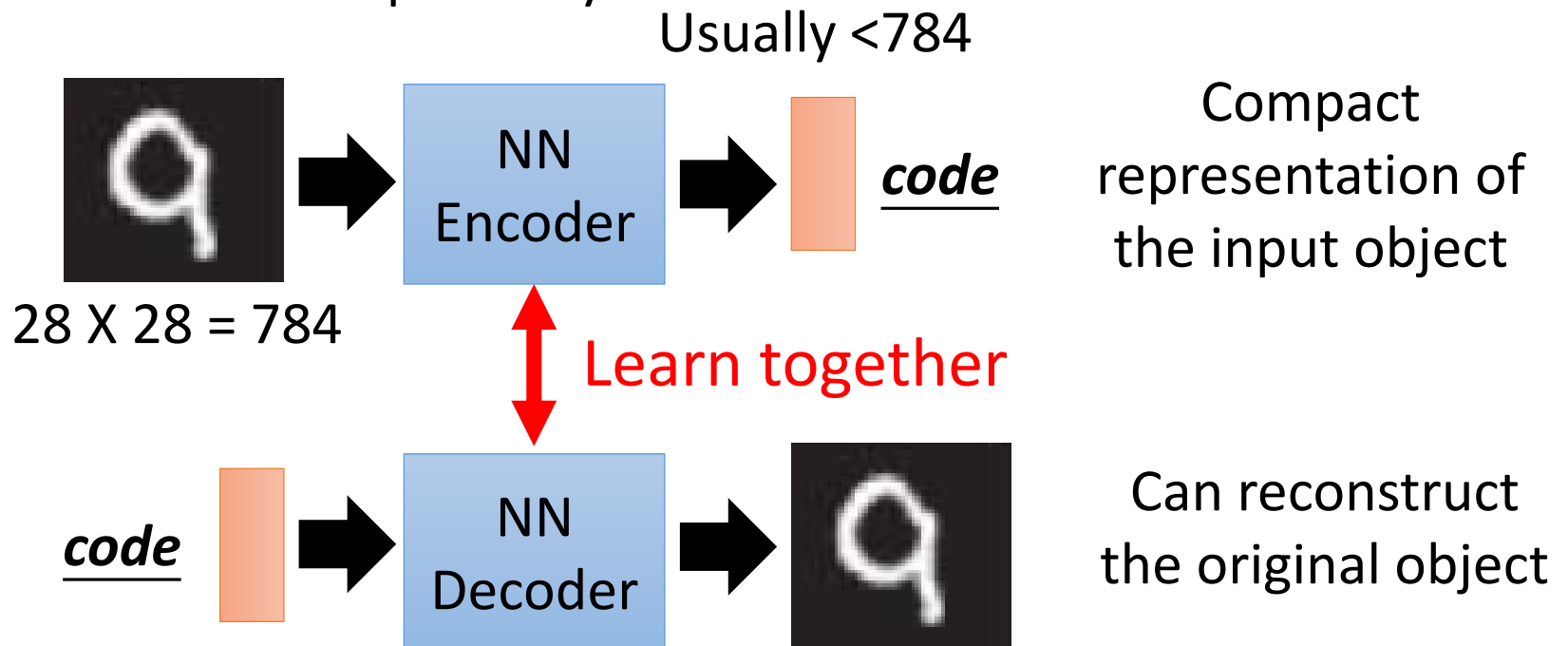
Auto-encoder



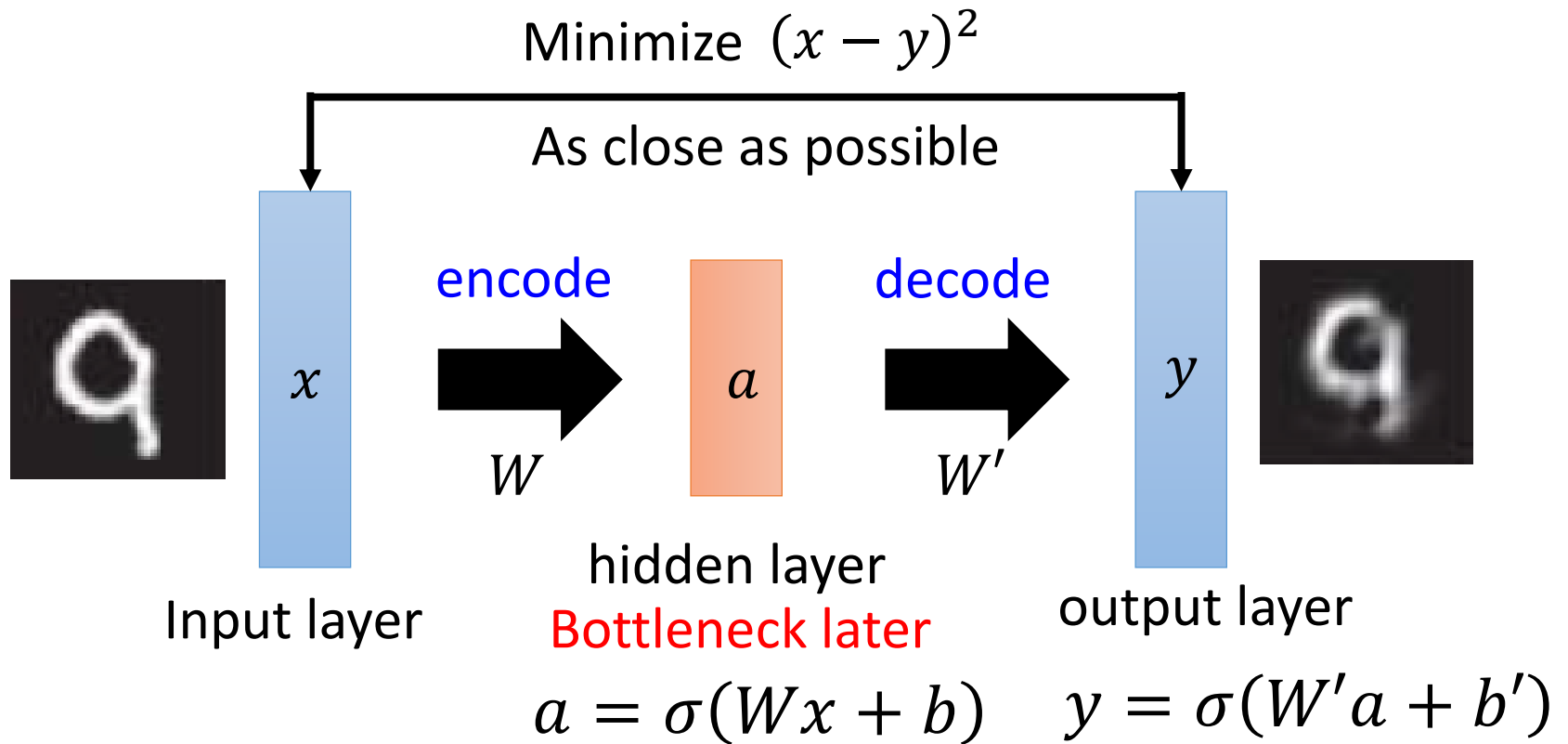
# Auto-encoder



- We use 28 X 28 d to represent a digit
- Not all 28 X 28 images are digit
- It is possible to represent the images of digits in a more compact way



# Auto-encoder



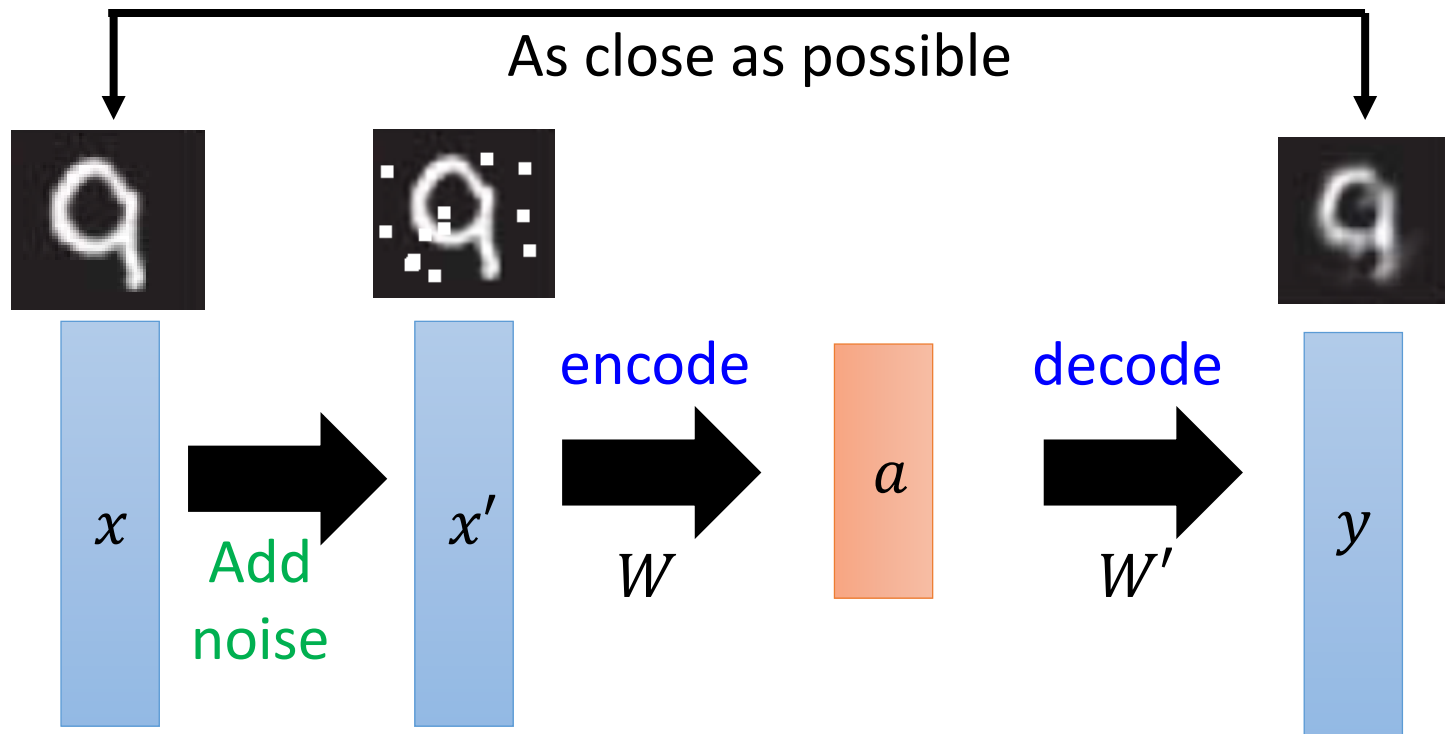
Output of the hidden layer is the code

## More: Contractive auto-encoder

# Auto-encoder

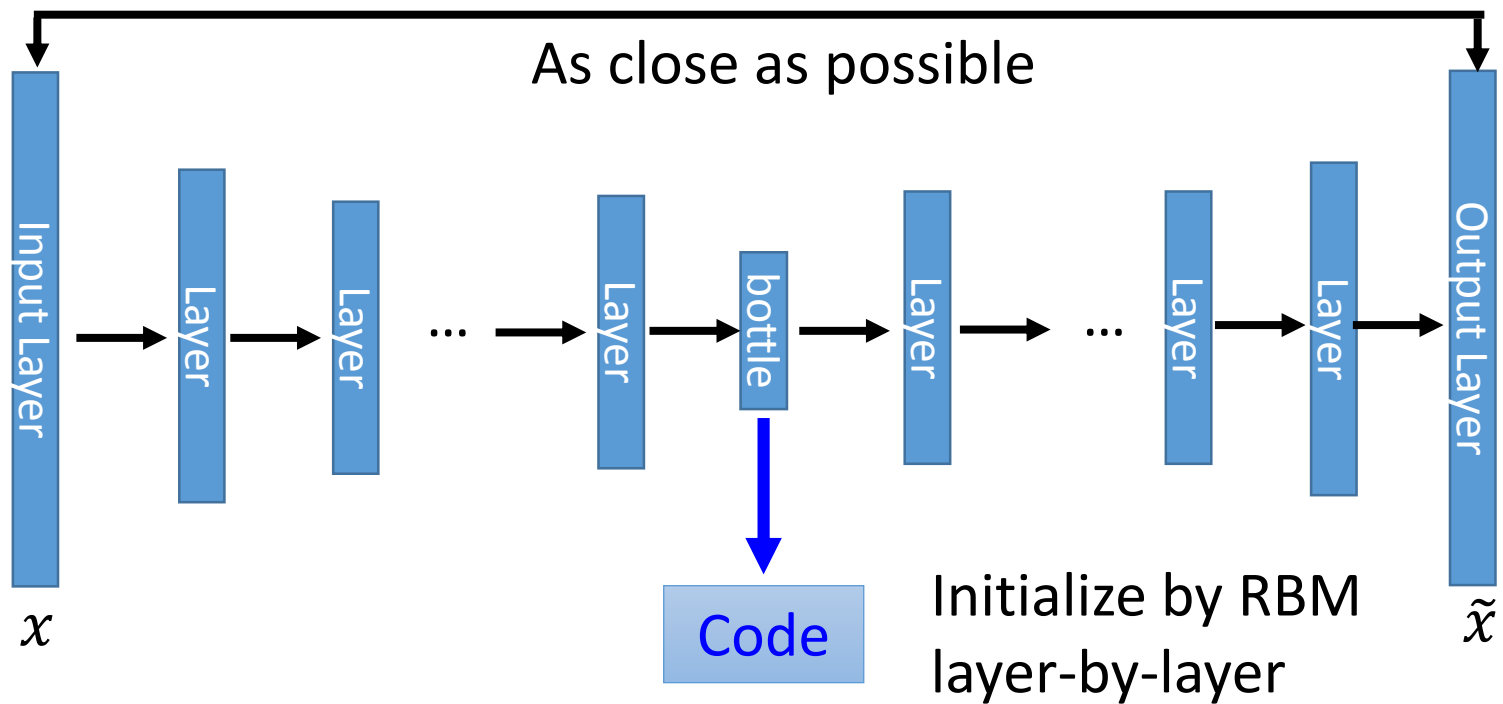
Ref: Rifai, Salah, et al. "Contractive auto-encoders: Explicit invariance during feature extraction." *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*. 2011.

- De-noising auto-encoder



# Deep Auto-encoder

- Of course, the auto-encoder can be deep



Reference: Hinton, Geoffrey E., and Ruslan R. Salakhutdinov. "Reducing the dimensionality of data with neural networks." *Science* 313.5786 (2006): 504-507

# Deep Auto-encoder

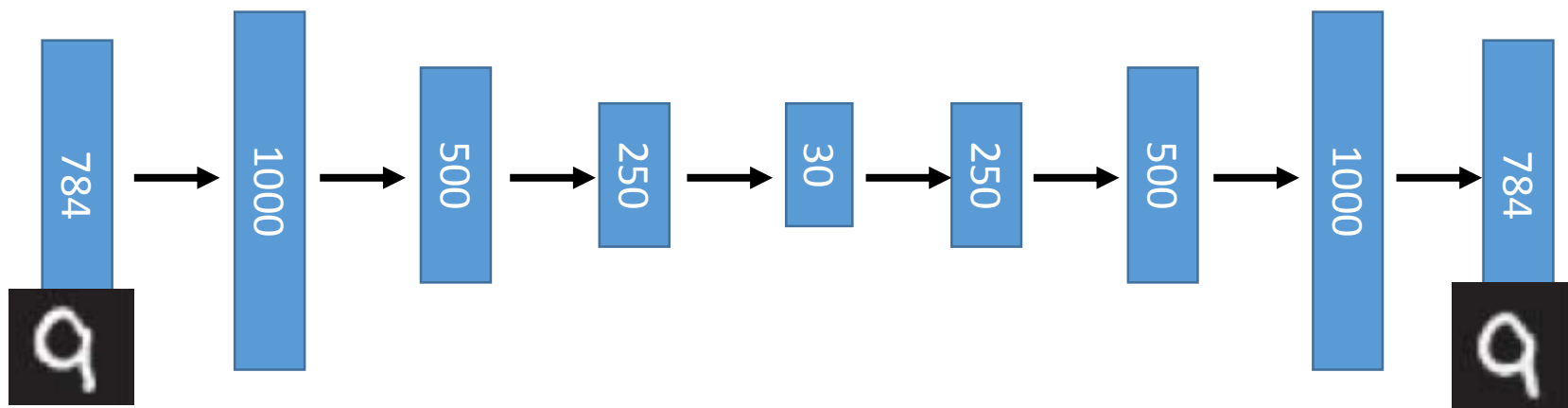
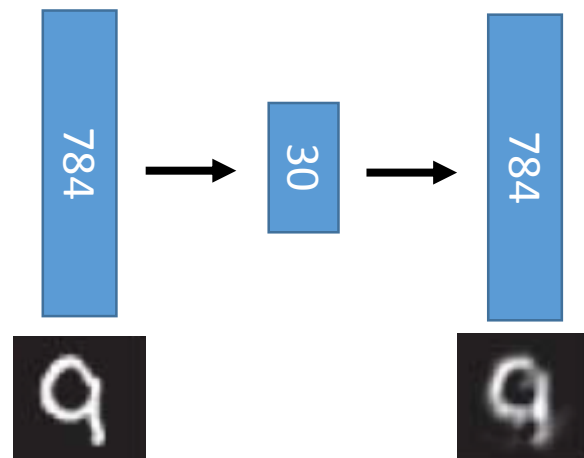
Original Image

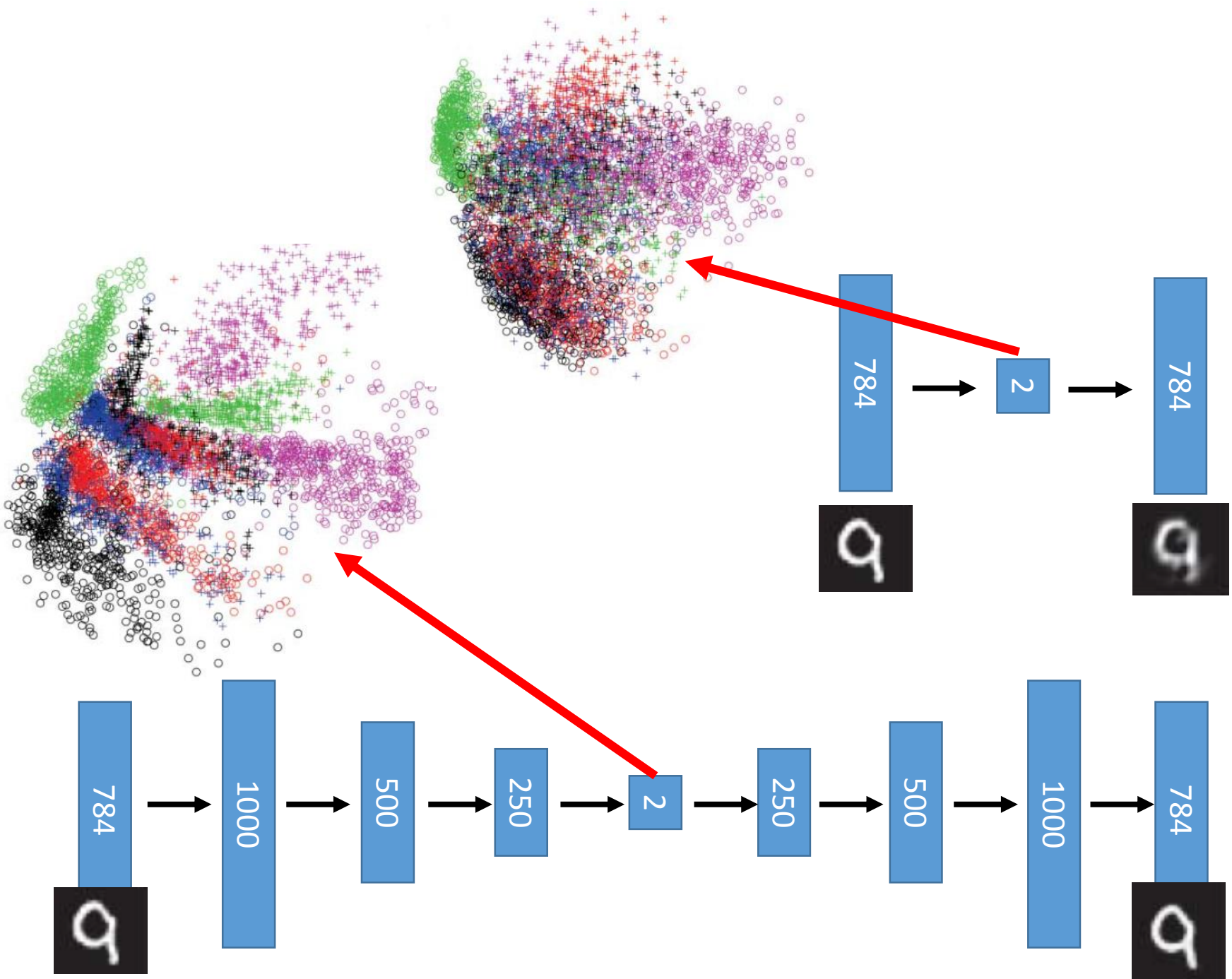


PCA



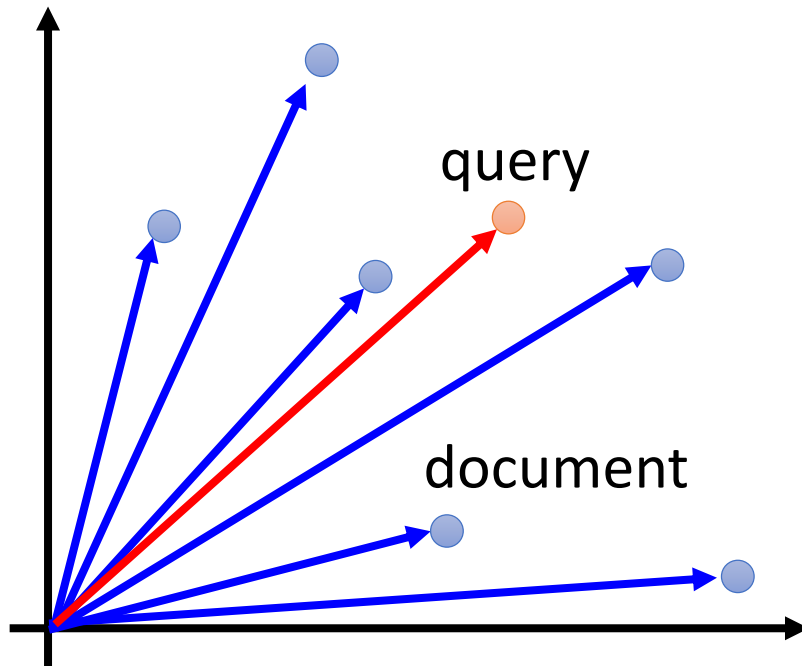
Deep Auto-encoder





# Auto-encoder – Text Retrieval

## Vector Space Model



## Bag-of-words

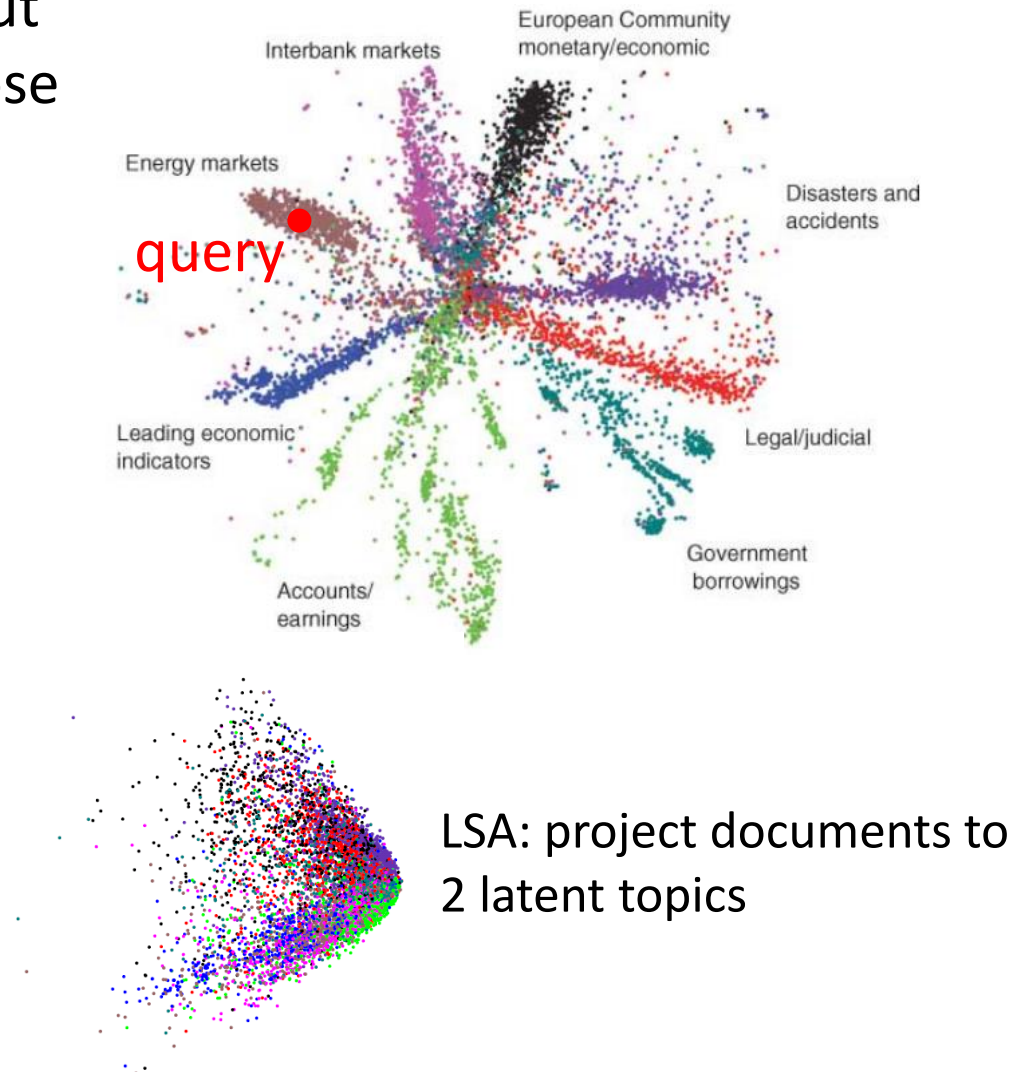
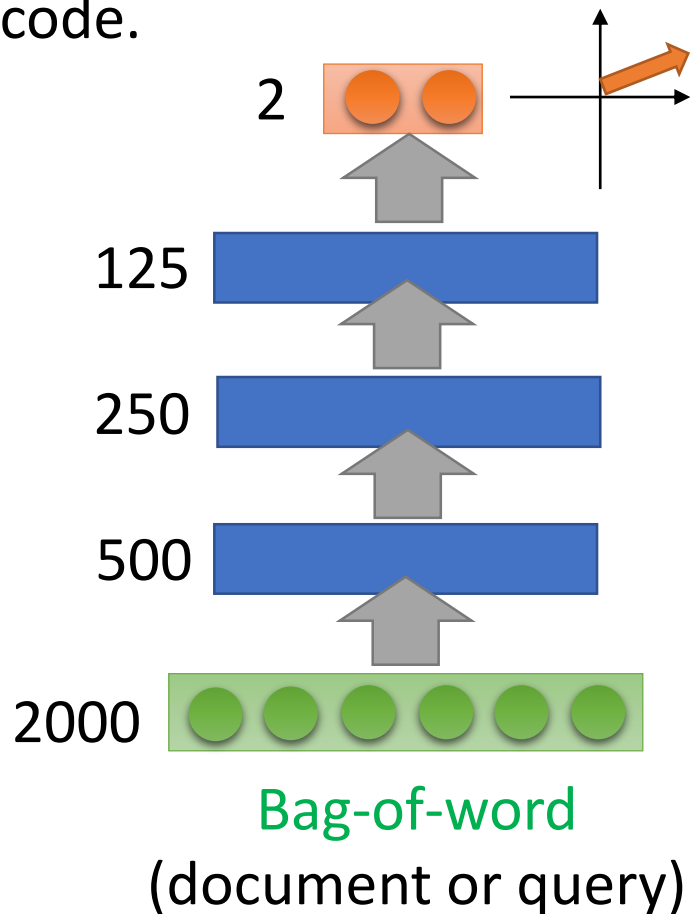
word string:  
"This is an apple"

this	●	1
is	●	1
a	●	0
an	●	1
apple	●	1
pen	●	0
⋮	●	

Semantics are not considered.

# Auto-encoder – Text Retrieval

The documents talking about the same thing will have close code.





# Auto-encoder – Similar Image Search

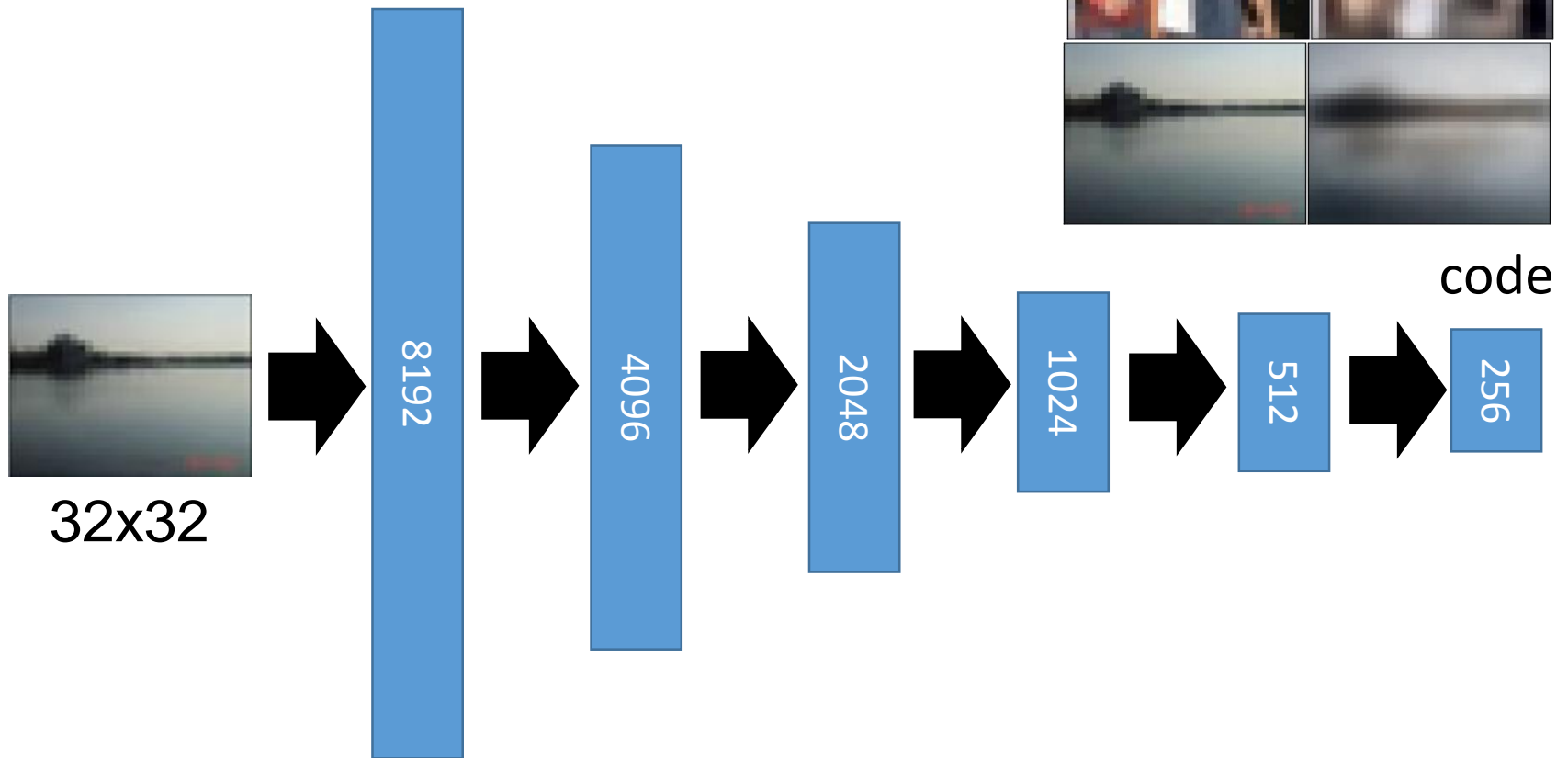
Retrieved using Euclidean distance in pixel intensity space



(Images from Hinton's slides on Coursera)

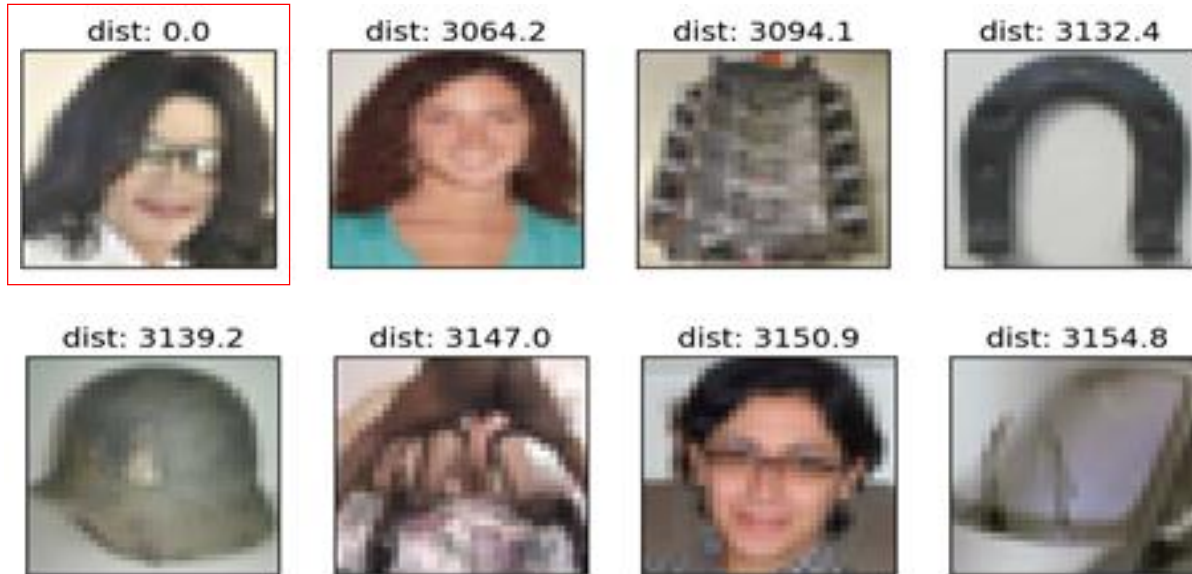
Reference: Krizhevsky, Alex, and Geoffrey E. Hinton. "Using very deep autoencoders for content-based image retrieval." *ESANN*. 2011.

# Auto-encoder – Similar Image Search



(crawl millions of images from the Internet)

# Retrieved using Euclidean distance in pixel intensity space

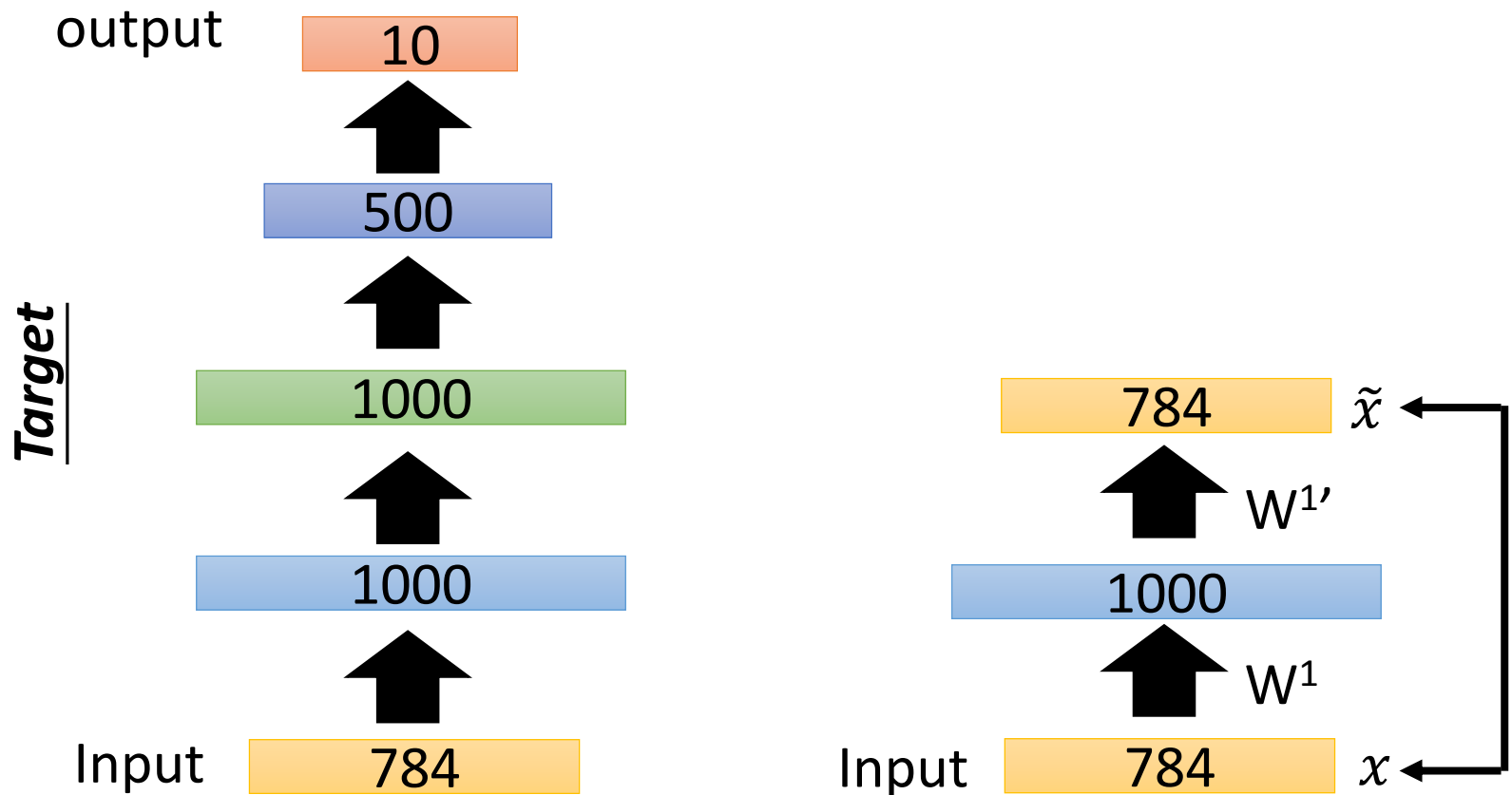


retrieved using 256 codes



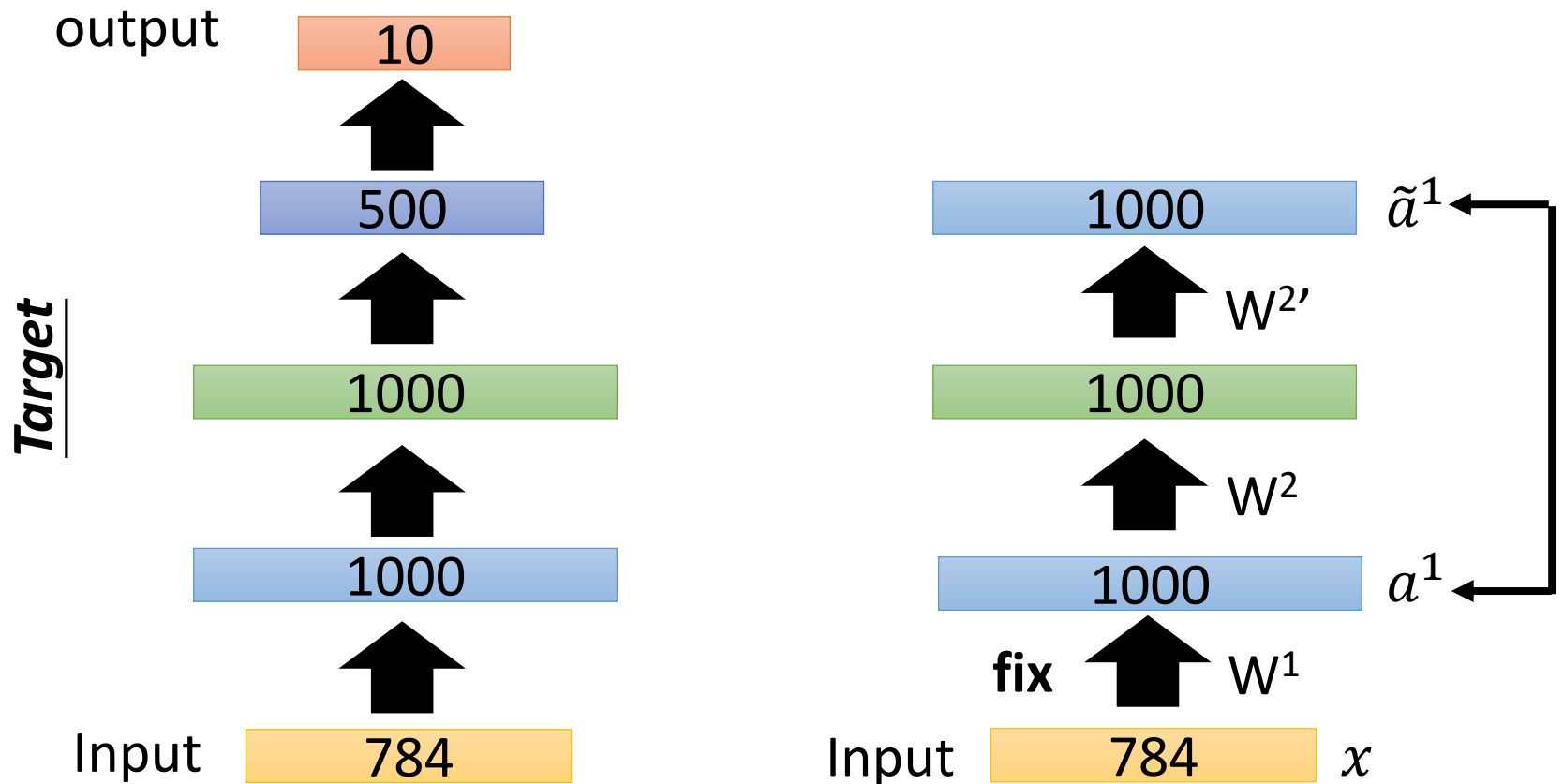
# Auto-encoder – Pre-training DNN

- Greedy Layer-wise Pre-training *again*



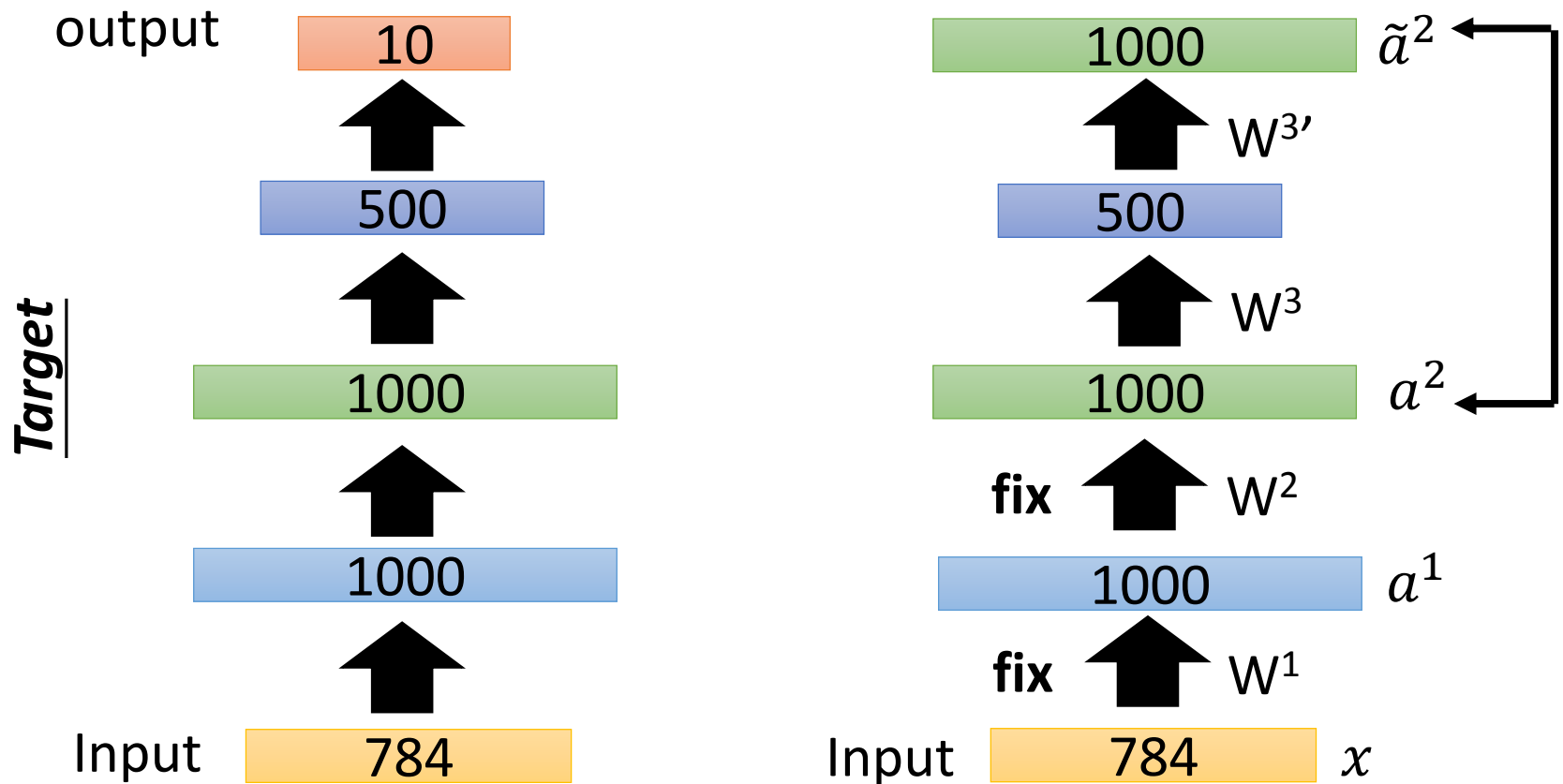
# Auto-encoder – Pre-training DNN

- Greedy Layer-wise Pre-training *again*



# Auto-encoder – Pre-training DNN

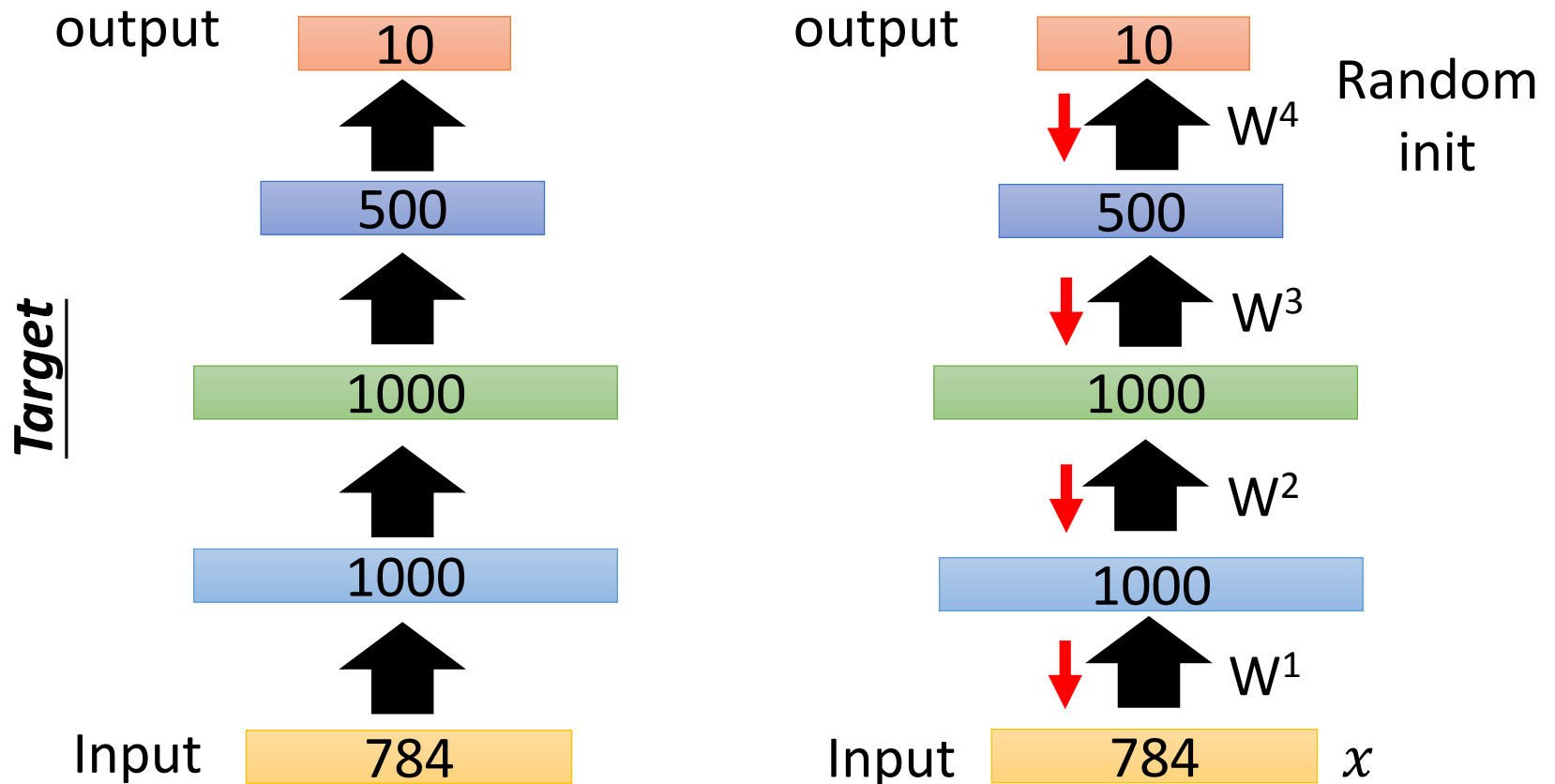
- Greedy Layer-wise Pre-training *again*



# Auto-encoder – Pre-training DNN

- Greedy Layer-wise Pre-training *again*

Find-tune by  
backpropagation



# Concluding Remarks



# Concluding Remarks

- Labeling data is expensive, but it is relatively easy to collect lots of unlabeled data.
- RBM and auto-encoder exploit the unlabeled data
- RBM and auto-encoder had been popular for pre-training DNN before.
- With sufficient labelled data and ReLU, pre-training is not that important.
  - However, it is still useful when you have lots of unlabeled data but little labelled data

# Plan

- 1/1 (五): 元旦放假
- 1/8 (五):
  - 2:30 ~ 3:30: TensorFlow: next generation of deep learning in Google (資工系 Seminar)
    - <https://www.csie.ntu.edu.tw/app/news.php?Sn=10358>
  - 4:00~: Attention-based Model
  - 23:59: presentation team decided
- 1/13 (三) 23:59: Presentation slides deadline
- 1/15 (五)
  - 14:20~: Presentation, 返鄉投票
- 1/16 (六): 投票
- 1/20 (三) 23:59: Report deadline

Merry Christmas

&

Happy New Year

# Acknowledgement

- 感謝 呂慶輝 同學指出投影片上的錯誤