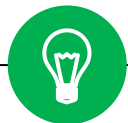


Applied Deep Learning



Beyond Supervised Learning



May 24th, 2021 <http://adl.miulab.tw>



**National
Taiwan
University**
國立臺灣大學

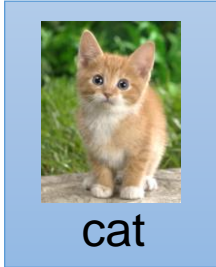
Introduction

- Big data \neq Big annotated data
- Machine learning techniques include:
 - Supervised learning (if we have labelled data)
 - Reinforcement learning (if we have an environment for reward)
 - **Unsupervised learning (if we do not have labelled data)**

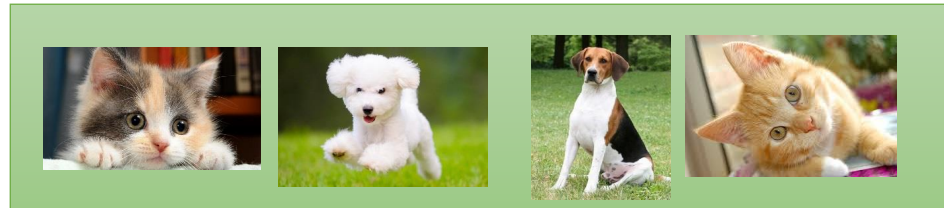
What can we do if there is no sufficient training data?

3 Semi-Supervised Learning

Labelled Data



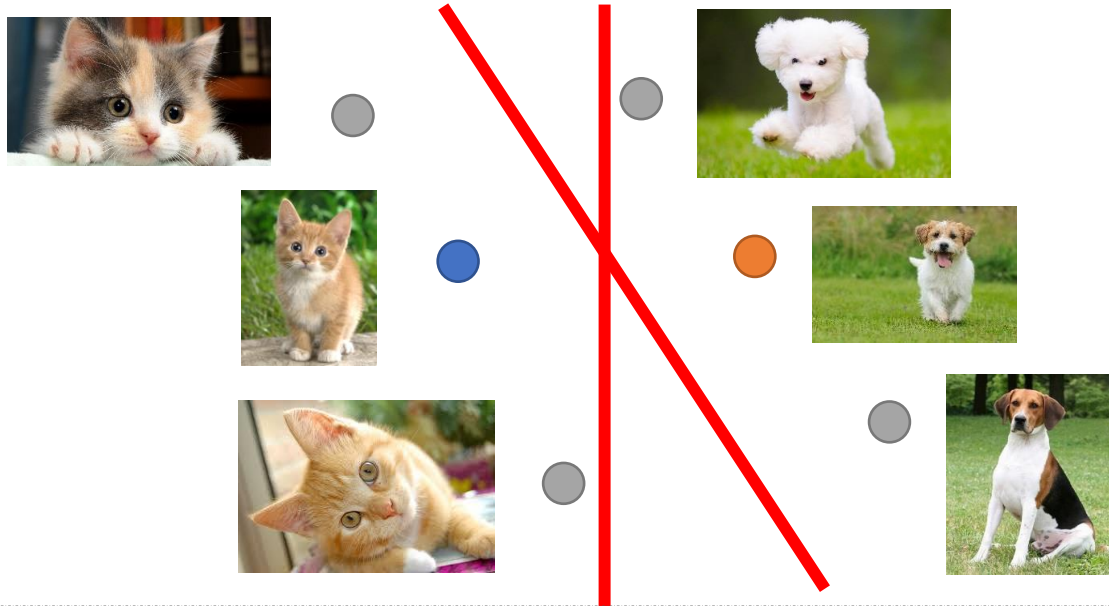
Unlabeled Data



(Image of cats and dogs without labeling)

4 Semi-Supervised Learning

- Why semi-supervised learning helps?



The distribution of the unlabeled data provides some cues

Transfer Learning

Source Data



cat



dog

Target Data



elephant



elephant



tiger

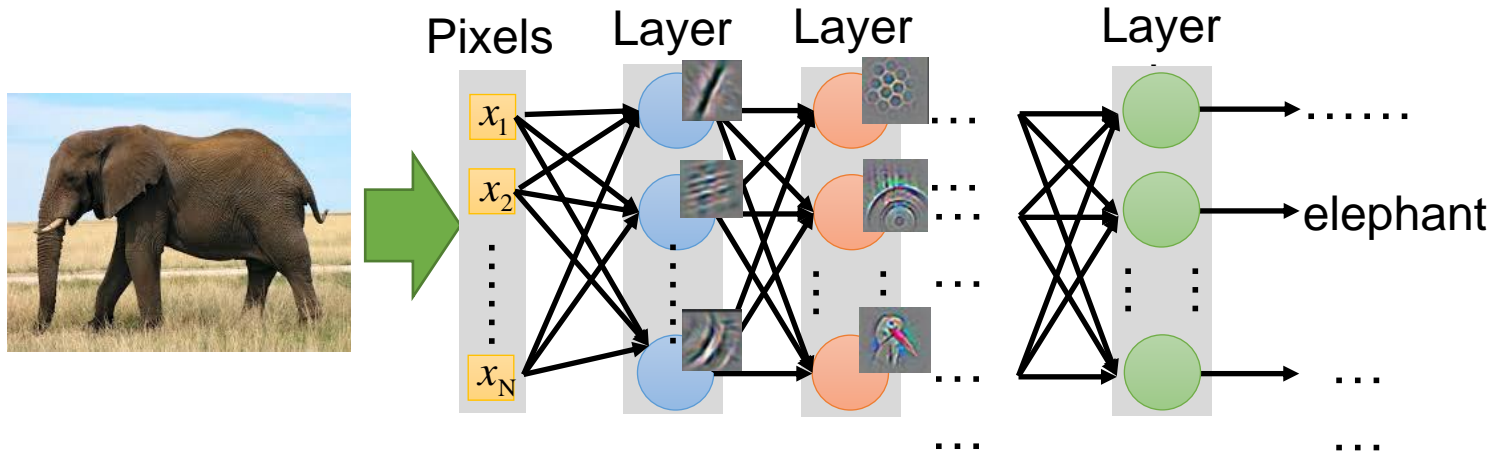


tiger

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



7 Transfer Learning Example

研究生 online

漫畫家 online

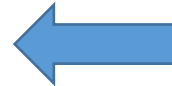
研究生
生存守則

研究生

指導教授

跑實驗

投稿期刊



漫畫家

責編

畫分鏡

投稿 jump



爆漫王

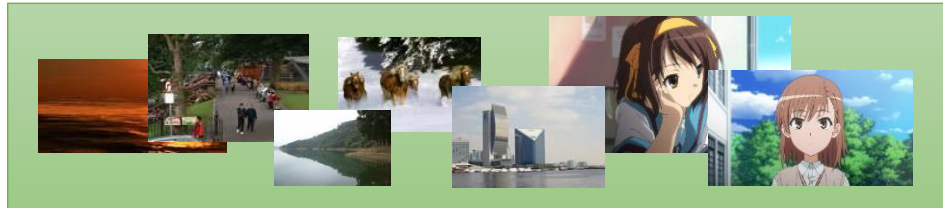
Self-Taught Learning

- ⦿ The unlabeled data sometimes is not related to the task

Labelled Data




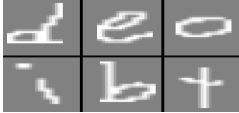

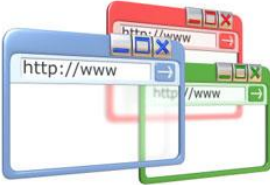
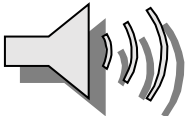

Unlabeled Data



(Just crawl millions of images from the Internet)

Self-Taught Learning

- The unlabeled data sometimes is not related to the task

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

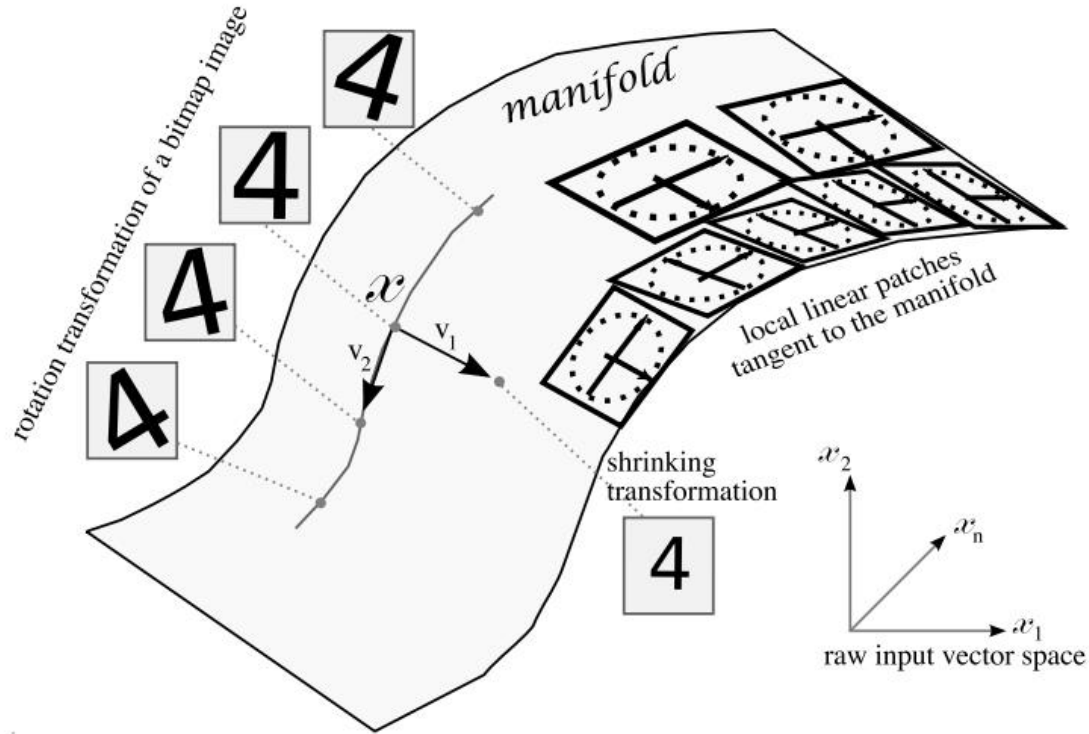
Why can we use unlabeled and unrelated data to help our tasks?

Self-Taught Learning

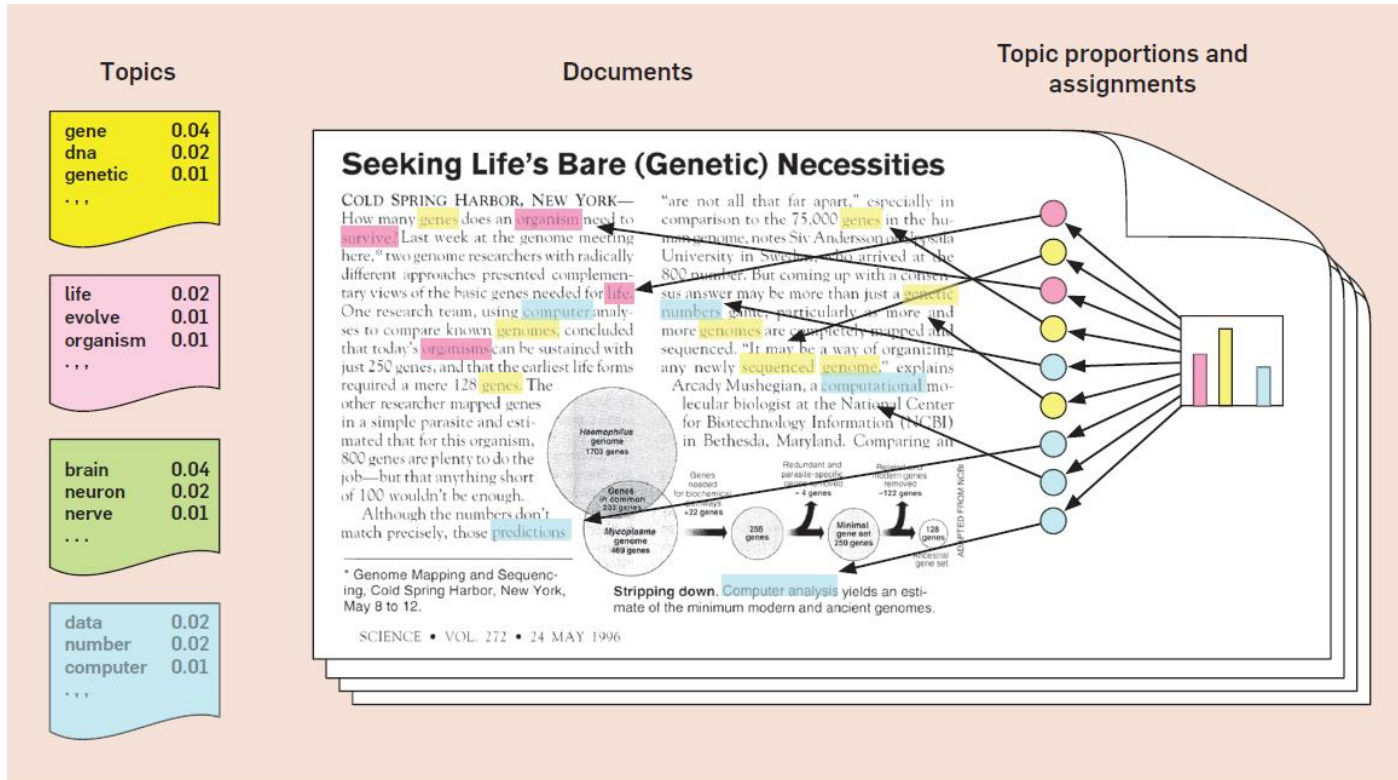
- ① How does self-taught learning work?
- ① Why does unlabeled and unrelated data help the tasks?

Finding latent factors that control the observations

Latent Factors for Handwritten Digits






Latent Factors for Documents



Latent Factors for Recommendation System

單純呆

傲嬌

A	 ✓	 ✓	 ✓	 ✓
B			✓	✓
C	✓	?	✓	

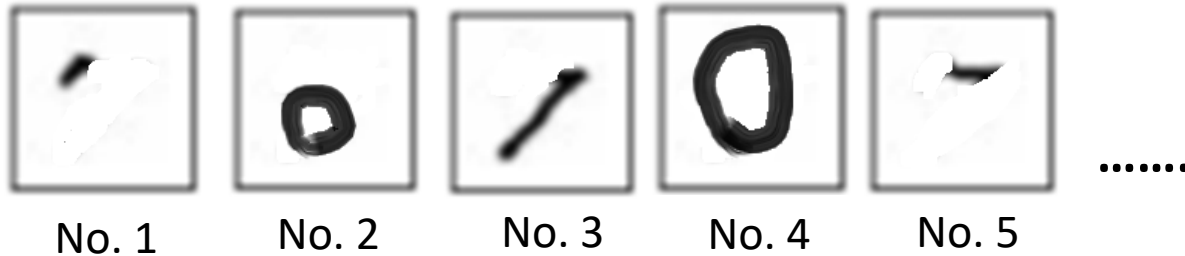
Latent Factor Exploitation

Handwritten digits



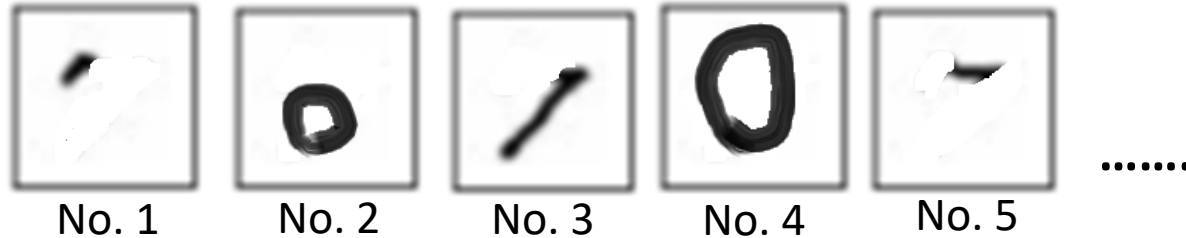
The handwritten images are composed of **strokes**

Strokes (Latent Factors)



Latent Factor Exploitation

Strokes (Latent Factors)



A 28x28 grayscale image of the digit 7 is shown on the left. To its right is an equals sign, followed by three 28x28 grayscale images of the latent factors No. 1, No. 3, and No. 5, separated by plus signs. Below the digit 7 is the text "Represented by 28 X 28 = 784 pixels". Below the latent factors is the text "[1 0 1 0 1 0]" and "(simpler representation)".

$$\begin{matrix} 28 \\ 28 \end{matrix} \begin{matrix} 28 \\ \text{7} \end{matrix} = \begin{matrix} \text{No. 1} \\ \text{No. 3} \\ \text{No. 5} \end{matrix} + \begin{matrix} \text{No. 1} \\ \text{No. 3} \\ \text{No. 5} \end{matrix} + \begin{matrix} \text{No. 1} \\ \text{No. 3} \\ \text{No. 5} \end{matrix}$$

Represented by
28 X 28 = 784 pixels

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

16

Autoencoder

Representation Learning

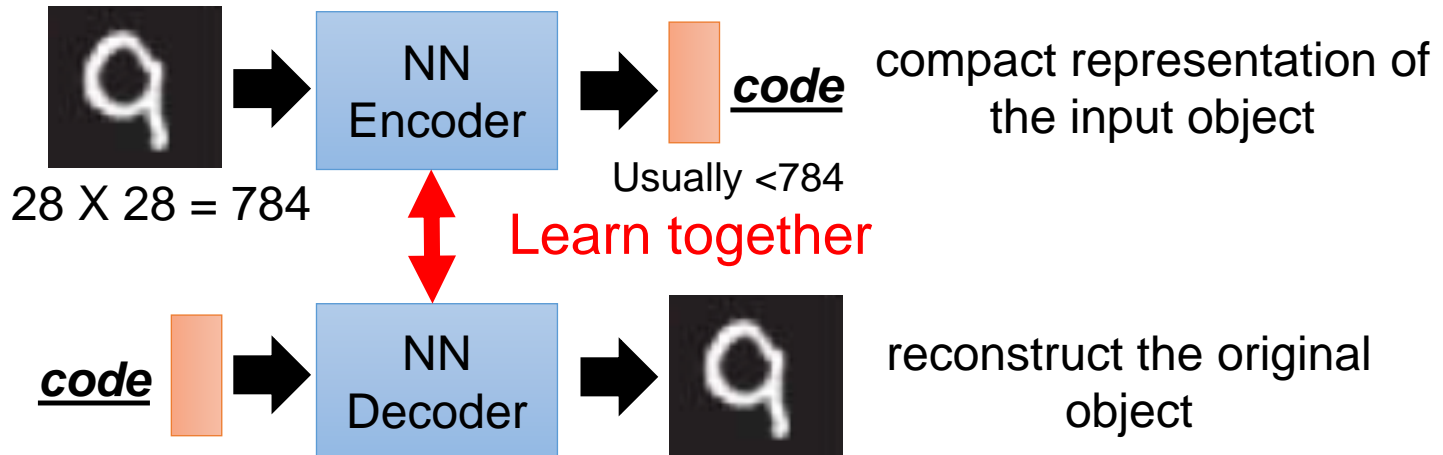
17

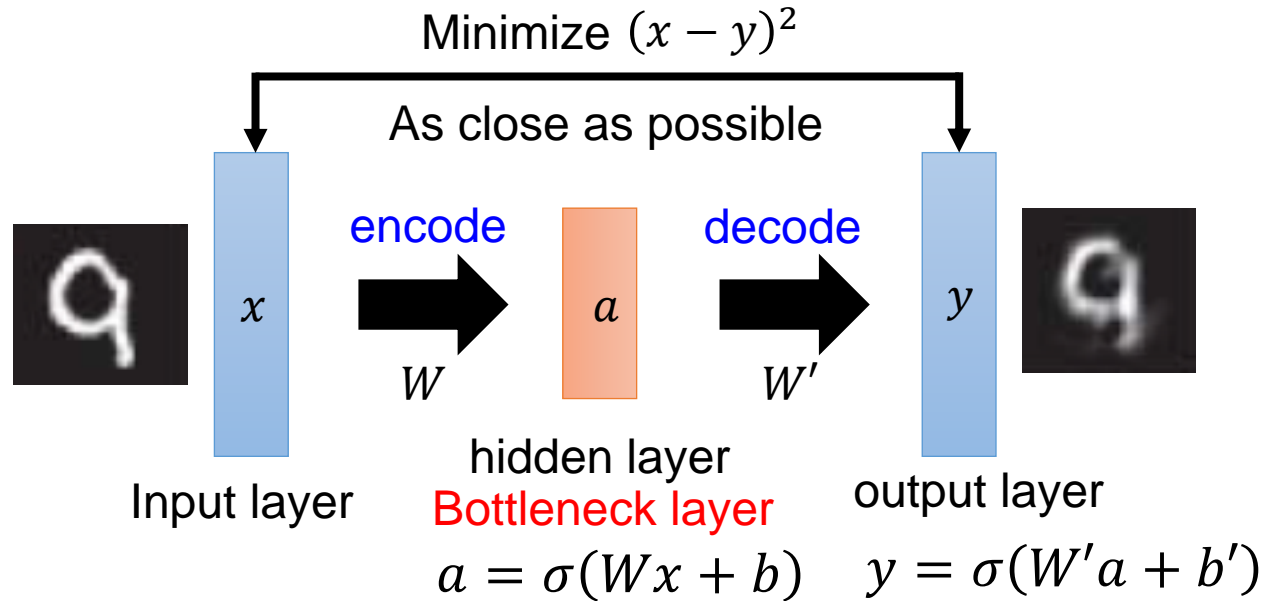
Autoencoder



- Represent a digit using 28 X 28 dimensions
- Not all 28 X 28 images are digits

Idea: represent the images of digits in a more compact way

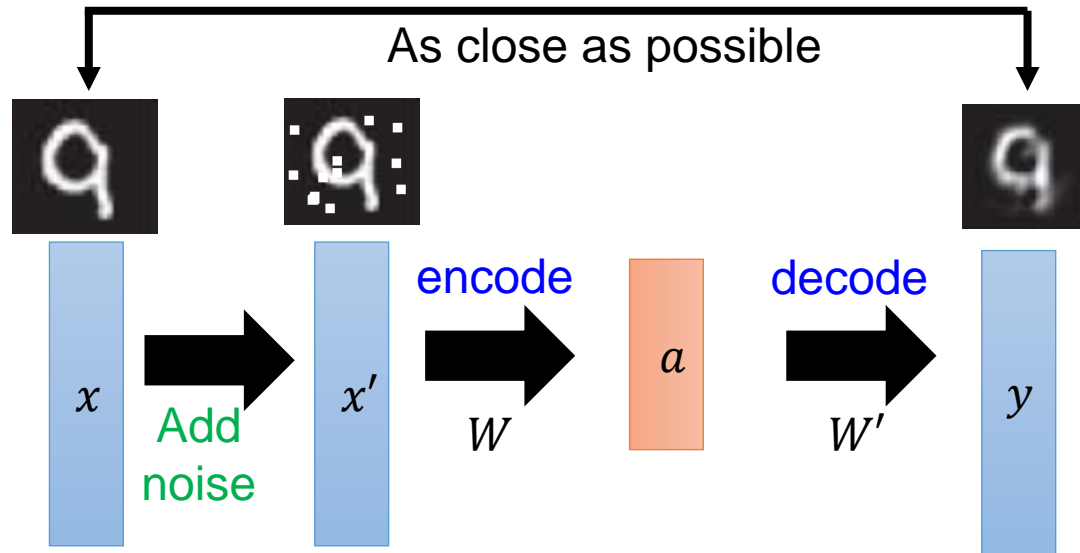


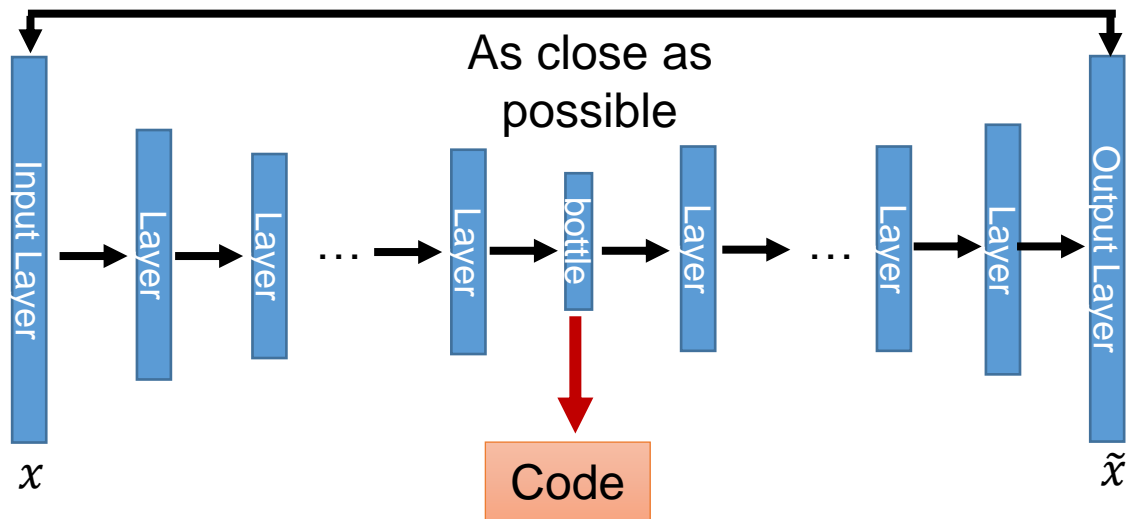


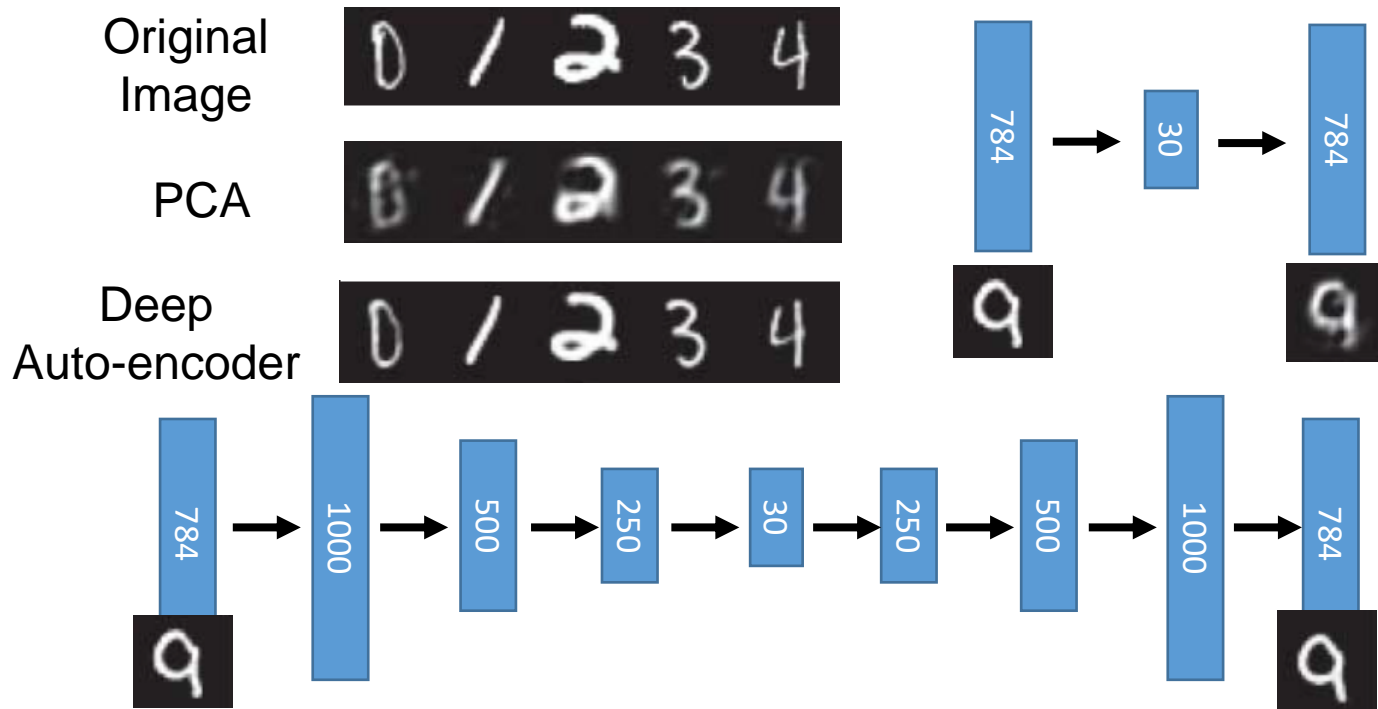
Output of the hidden layer is the code

Autoencoder

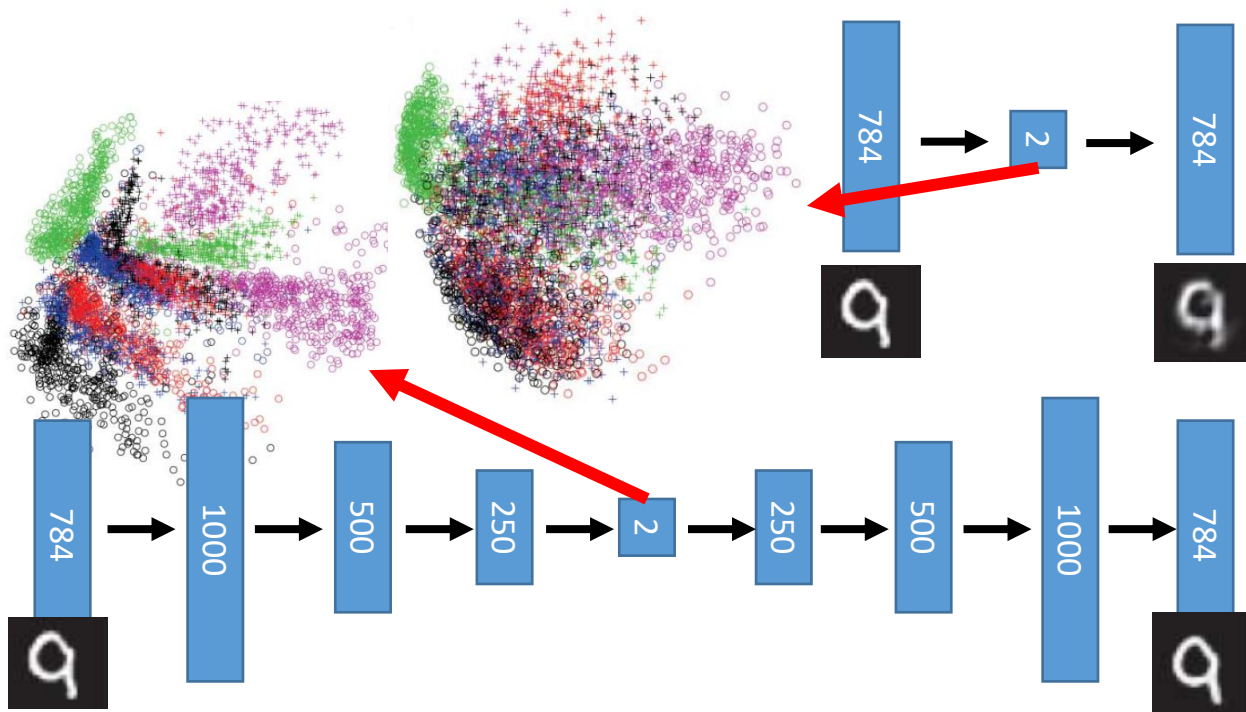
De-noising auto-encoder





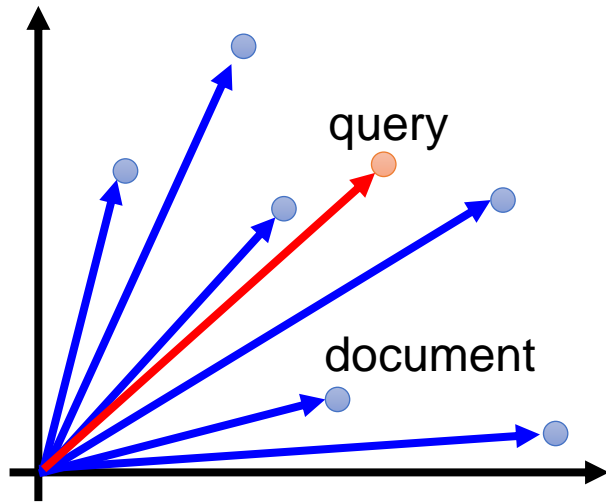


Feature Representation



Auto-encoder – Text Retrieval

Vector Space Model



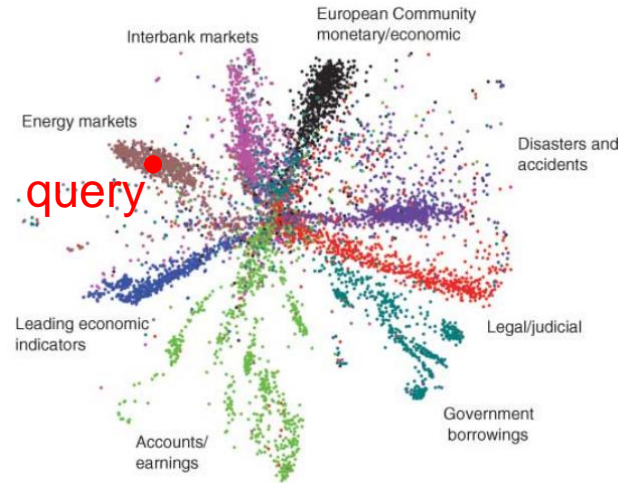
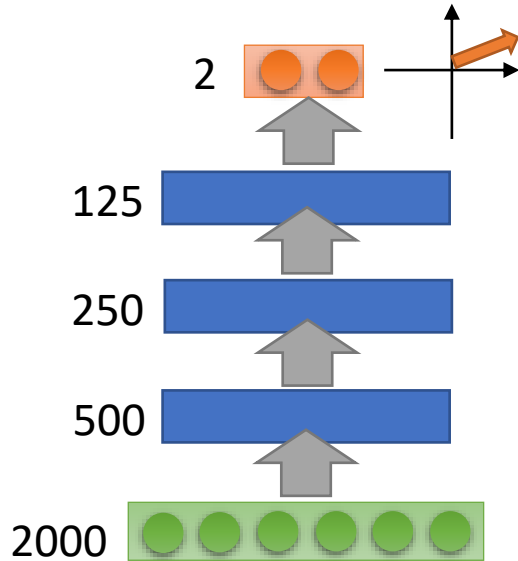
Bag-of-words

word string:
"This is an apple"

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

Semantics are not considered

Autoencoder – Text Retrieval



Bag-of-words (document or query)

The documents talking about the same thing will have close code

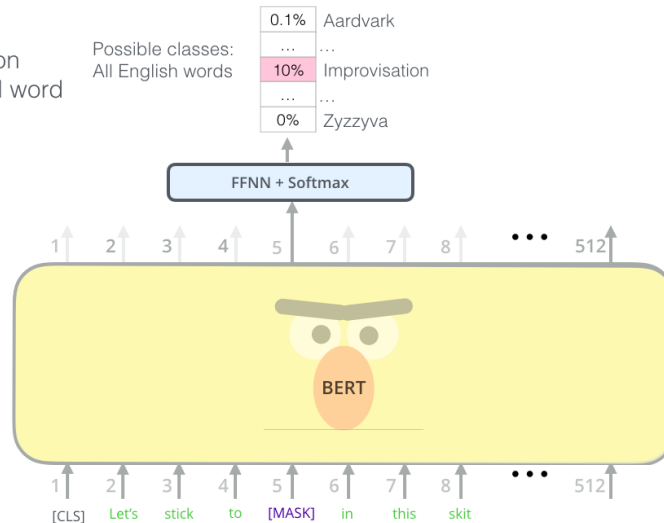
Auto-Encoding (AE)

- Objective: reconstructing \bar{x} from \hat{x}

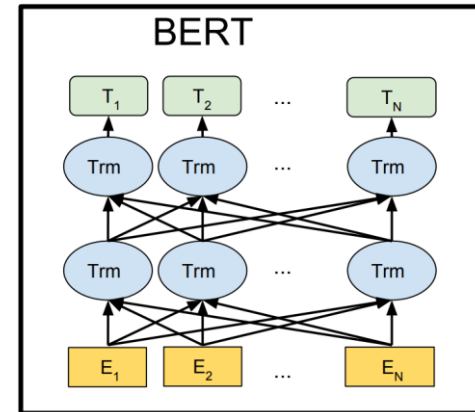
$$\max_{\theta} \log p_{\theta}(\bar{x} | \hat{x}) \approx \sum_{t=1}^T m_t \log p_{\theta}(x_t | \hat{x}) = \sum_{t=1}^T m_t \log \frac{\exp(H_{\theta}(\hat{x})_t^{\top} e(x_t))}{\sum_{x'} \exp(H_{\theta}(\hat{x})_t^{\top} e(x'))}$$

- dimension reduction or denoising (masked LM)

Use the output of the masked word's position to predict the masked word



Randomly mask
15% of tokens

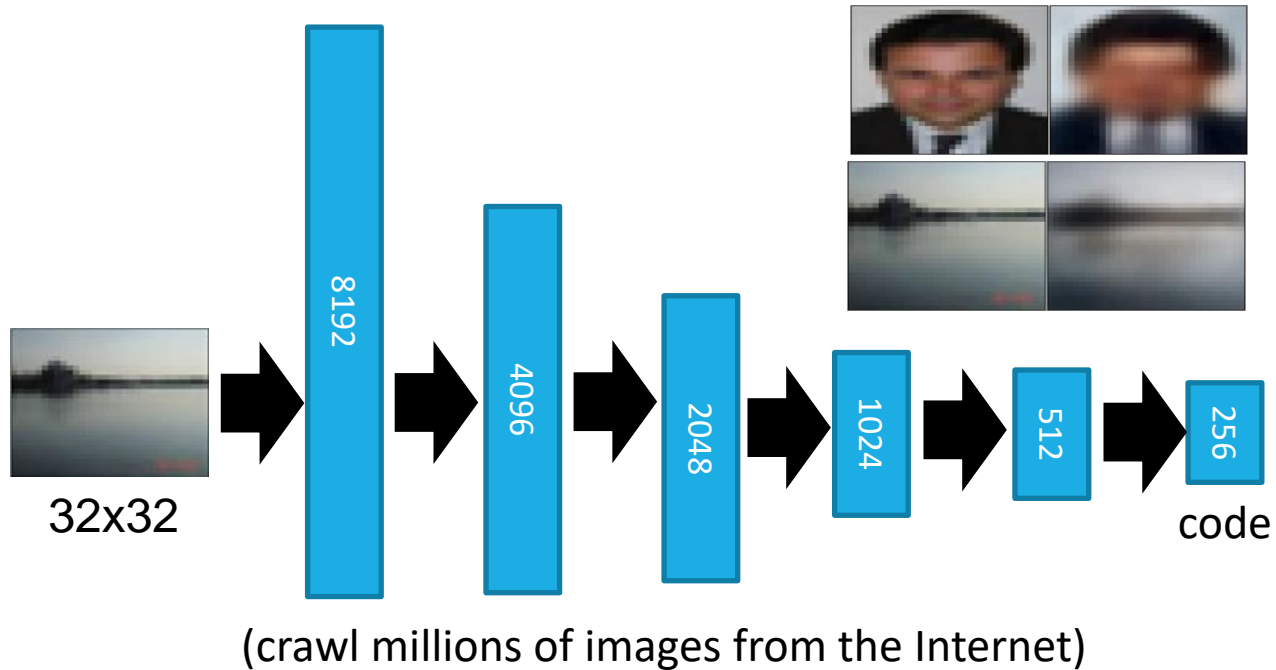


Autoencoder – Similar Image Retrieval

- Retrieved using Euclidean distance in pixel intensity space

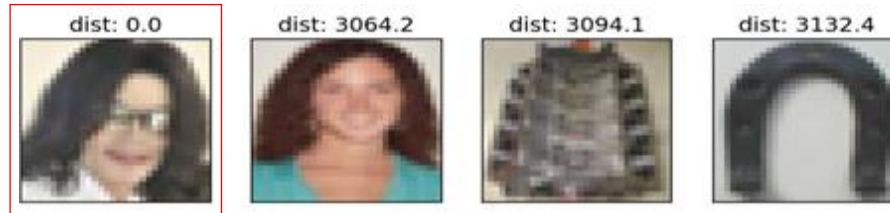


Autoencoder – Similar Image Retrieval



Autoencoder – Similar Image Retrieval

- Images retrieved using Euclidean distance in pixel intensity space



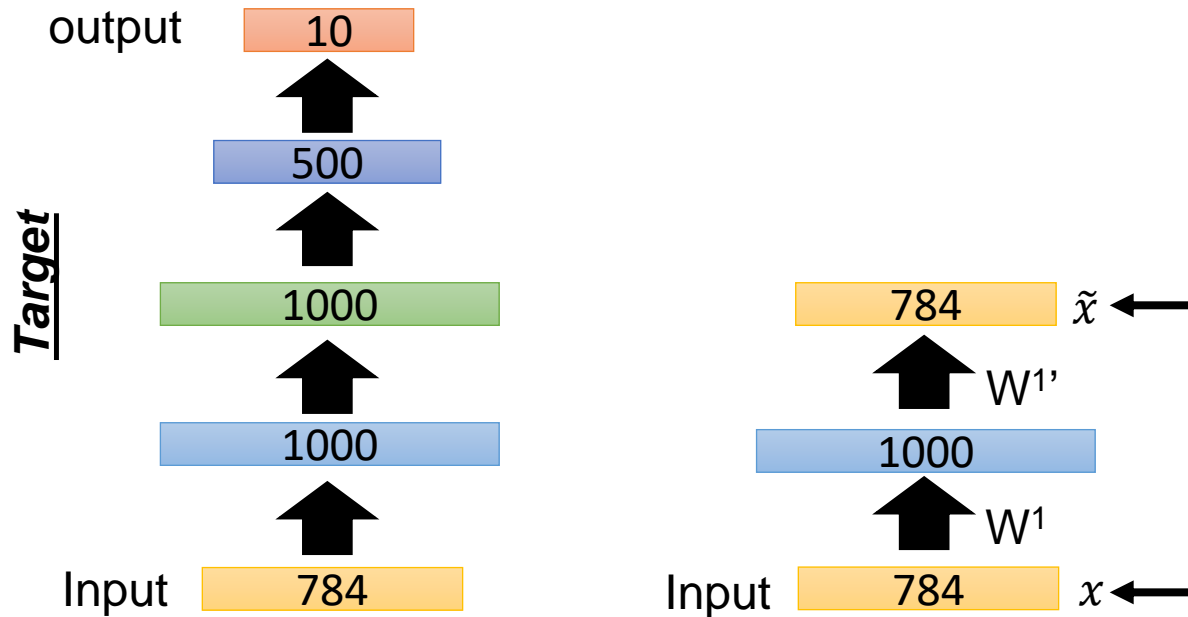
- Images retrieved using 256 codes



Learning the useful latent factors

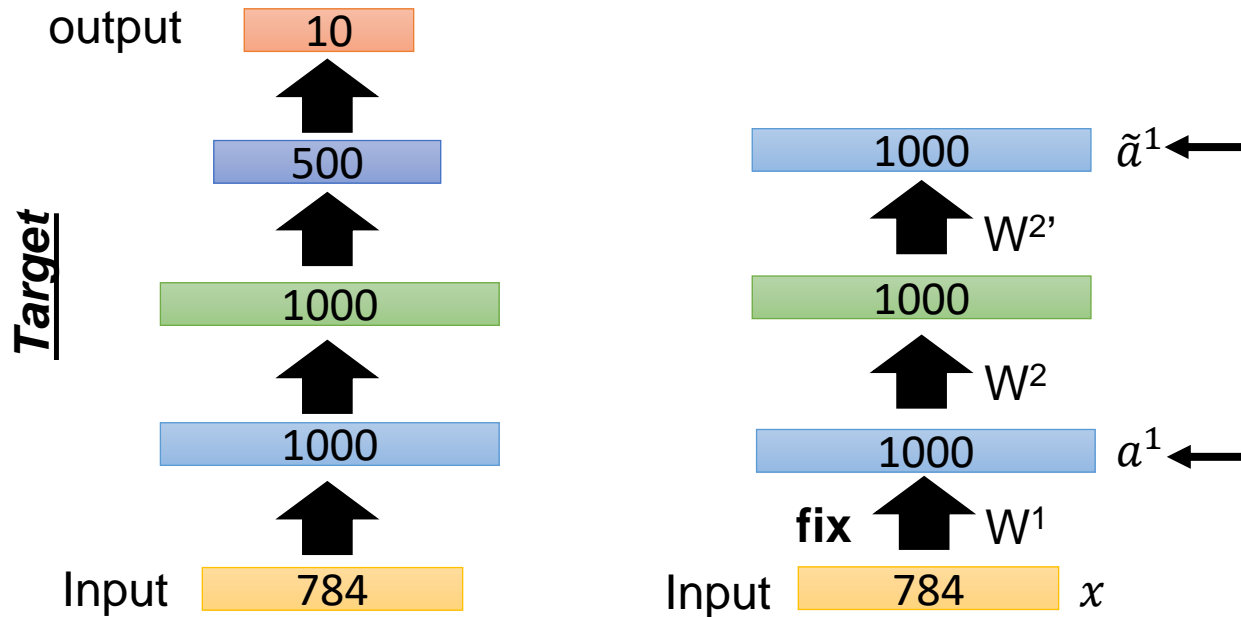
Autoencoder for DNN Pre-Training

- Greedy layer-wise pre-training *again*



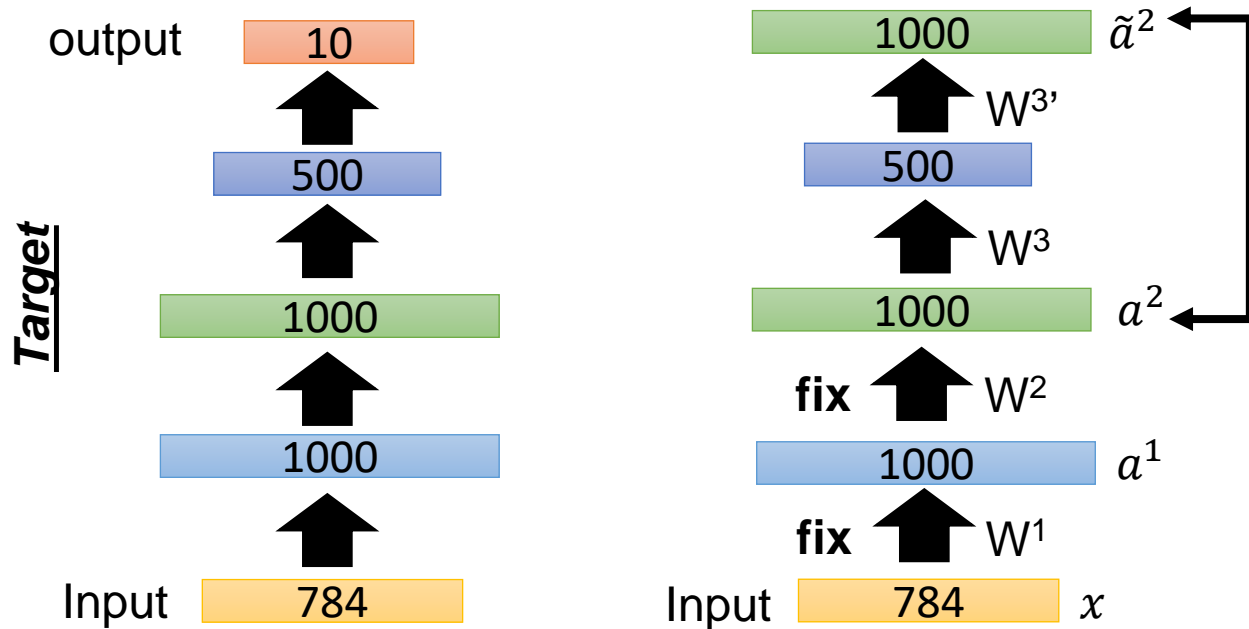
Autoencoder for DNN Pre-Training

- Greedy layer-wise pre-training *again*



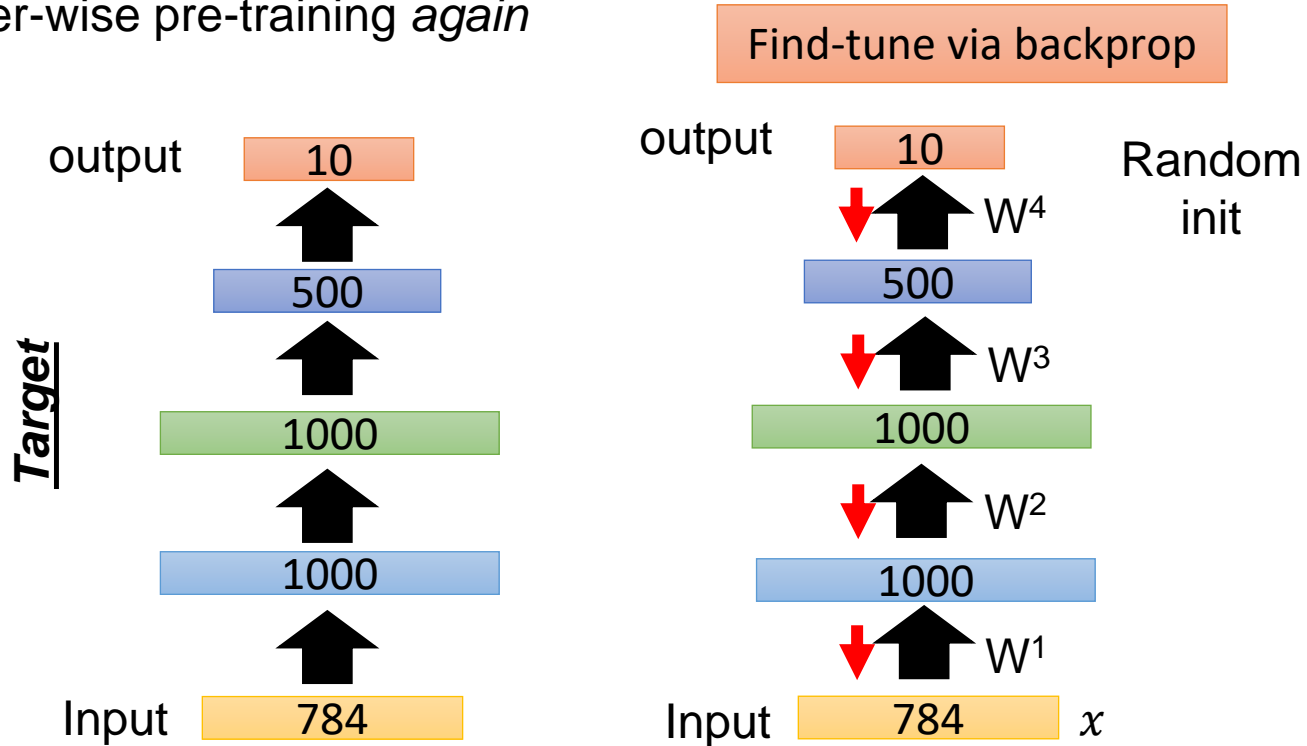
Autoencoder for DNN Pre-Training

- Greedy layer-wise pre-training *again*



32 Autoencoder for DNN Pre-Training

- Greedy layer-wise pre-training *again*

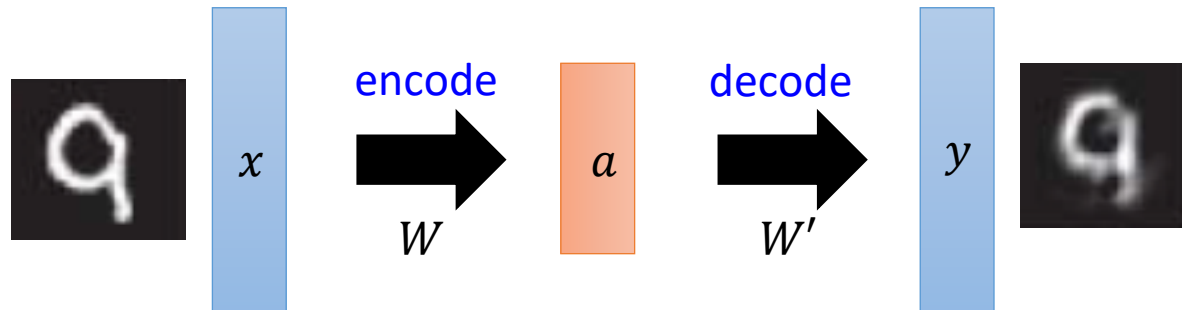


33

Variational Autoencoder

Representation Learning and Generation

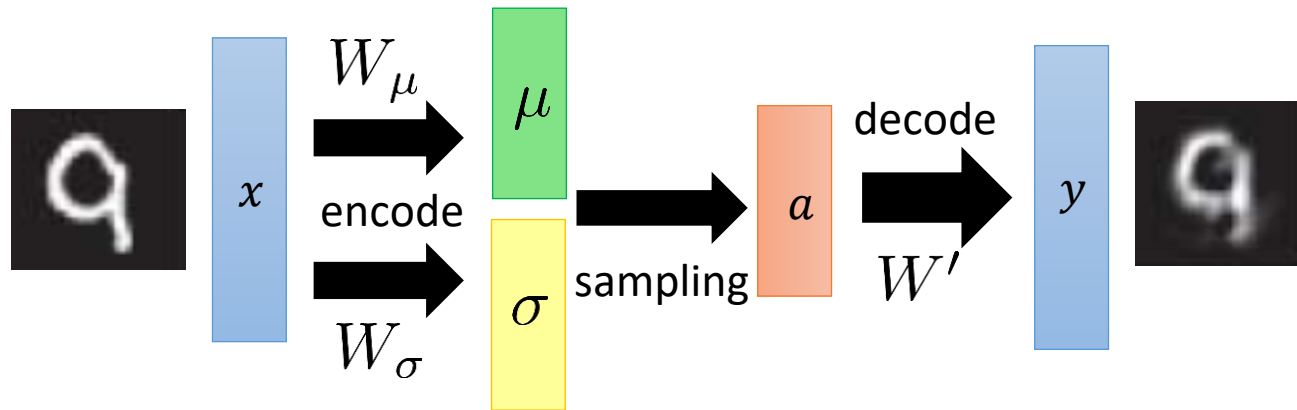
Generation from Latent Codes



How can we set a latent code for generation?

Latent Code Distribution Constraints

- ⦿ Constrain the data distribution for learned latent codes
- ⦿ Generate the latent code via a prior distribution



AE



VAE



37

Distant Supervision

Representation Learning by Weak Labels

Convolutional Deep Structured Semantic Models (CDSSM/DSSM)

Semantic Layer: y

Semantic Projection Matrix: W_s

Max Pooling Layer: I_m

Max Pooling Operation

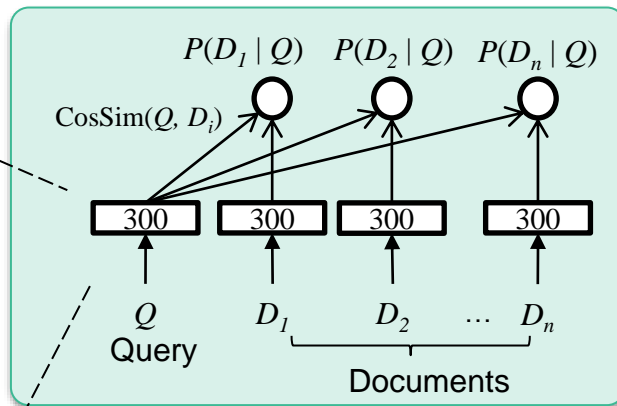
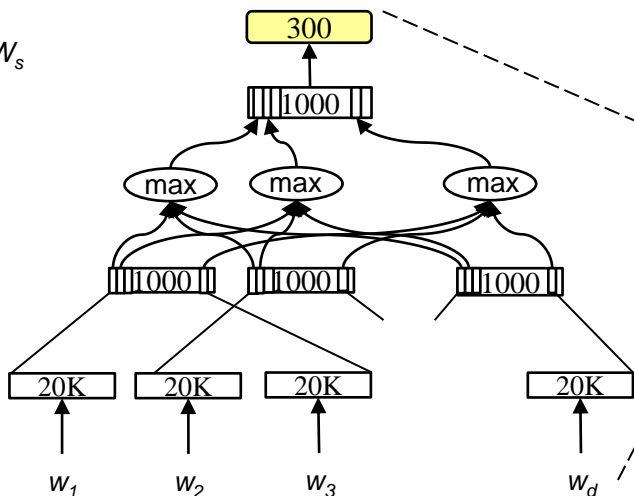
Convolutional Layer: I_c

Convolution Matrix: W_c

Word Hashing Layer: I_h

Word Hashing Matrix: W_h

Word Sequence: x



$$P(D | Q) = \frac{\exp(\text{CosSim}(Q, D))}{\sum_{D'} \exp(\text{CosSim}(Q, D'))}$$

$$\Lambda(\theta) = \log \prod_{(Q, D^+)} P(D^+ | Q)$$

maximizes the likelihood of clicked documents given queries

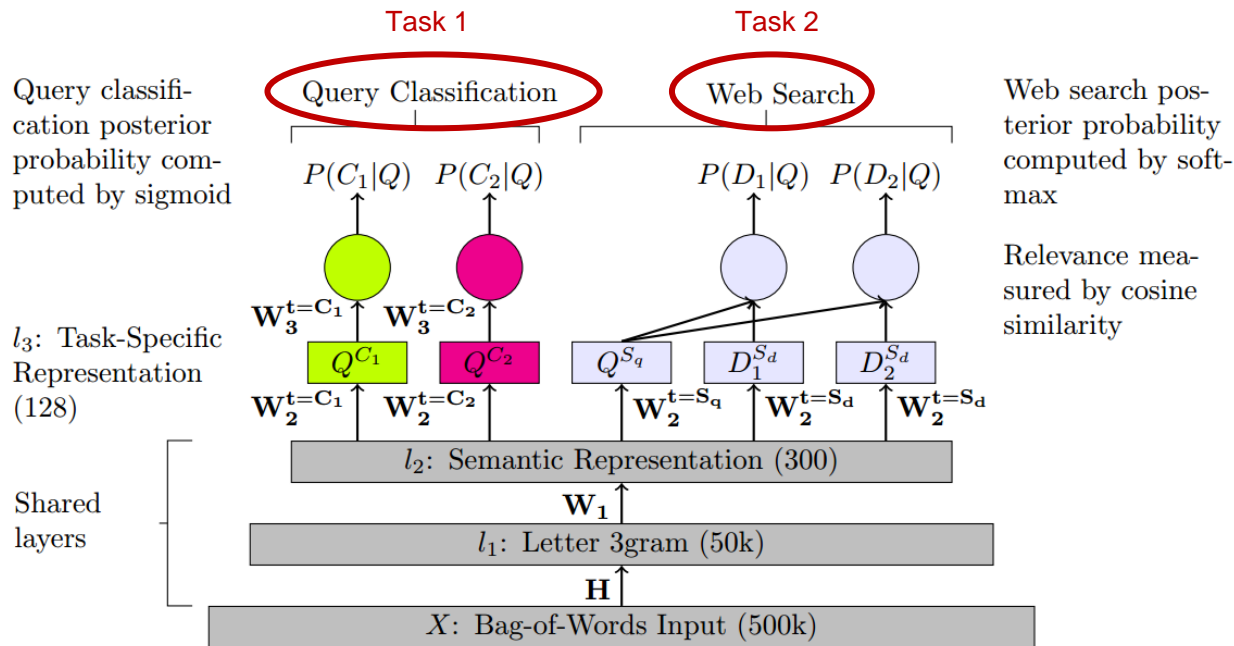
Semantically related documents are close to the query in the encoded space

39

Multi-Task Learning

Representation Learning by Different Tasks

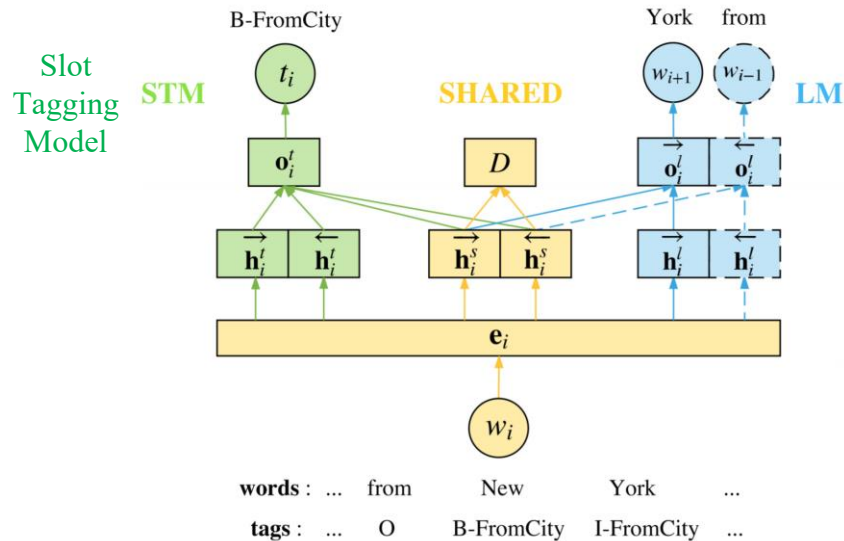
Task-Shared Representation



The latent factors can be learned by different tasks

Semi-Supervised Multi-Task SLU (Lan et al., 2018)

- Idea: language understanding objective can enhance other tasks



Algorithm 1: Adversarial Multi-task Learning for SLU

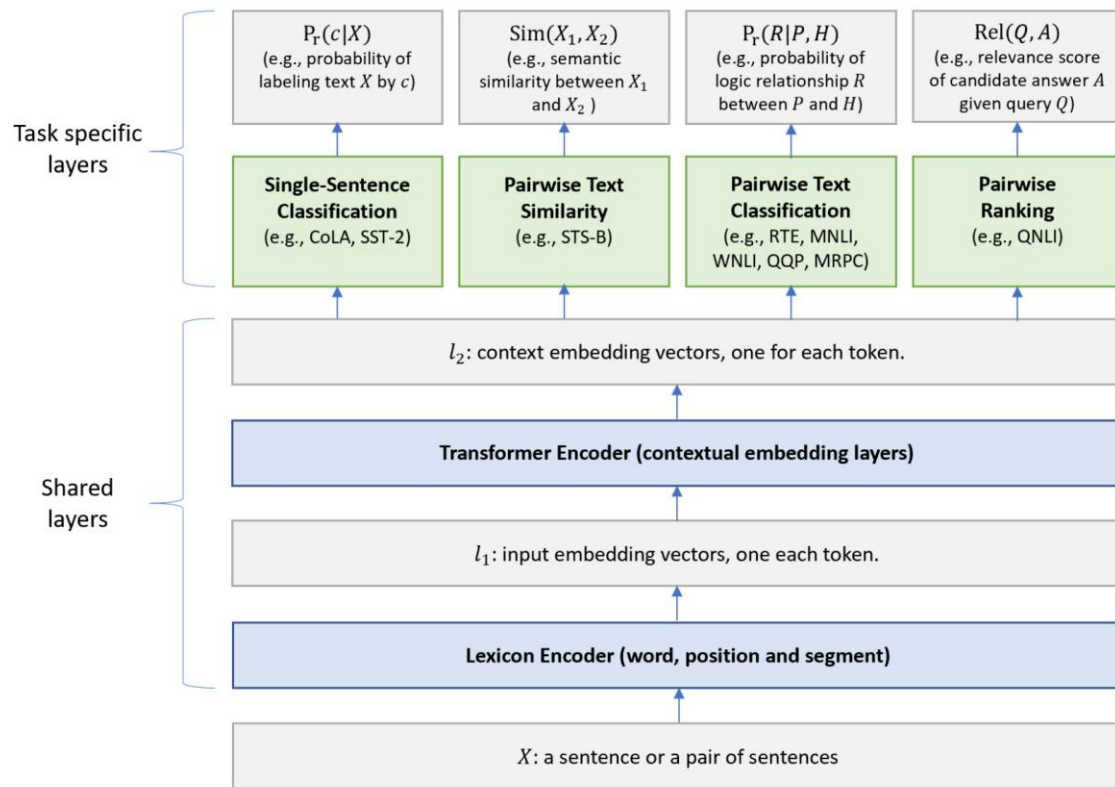
Input : Labeled training data $\{(\mathbf{w}^l, \mathbf{t}^l)\}$
Unlabeled data $\{\mathbf{w}^u\}$

Output: Adversarially enhanced slot tagging model

- 1 Initialize parameters $\{\theta^s, \theta^t, \theta^l, \theta^d\}$ randomly.
 - 2 **repeat**
 - /* Sample from $\{(\mathbf{w}^l, \mathbf{t}^l)\}$ */
 - 3 Train the STM and shared model by Eq.(8).
 - 4 Train the task discriminator and the shared model by Eq.(6) or Eq.(7) as slot tagging task ($y = 1$).
 - /* Sample from $\{\mathbf{w}^l\}$ and $\{\mathbf{w}^u\}$ */
 - 5 Train the LM and shared models by Eq.(9) (and Eq.(10) for BLM).
 - 6 Train the task discriminator and the shared model by Eq.(6) or Eq.(7) as LM task ($y = 0$).
 - 7 **until** convergence;
-

BLM exploits the *unsupervised knowledge*, the *shared-private framework* and *adversarial training* make the slot tagging model more generalized

MT-DNN (Liu et al., 2019)



Algorithm 1: Training a MT-DNN model.

```

Initialize model parameters  $\Theta$  randomly.
Pre-train the shared layers (i.e., the lexicon
encoder and the transformer encoder).
Set the max number of epoch:  $epoch_{max}$ .
//Prepare the data for  $T$  tasks.
for  $t$  in  $1, 2, \dots, T$  do
    | Pack the dataset  $t$  into mini-batch:  $D_t$ .
end
for  $epoch$  in  $1, 2, \dots, epoch_{max}$  do
    1. Merge all the datasets:
        $D = D_1 \cup D_2 \dots \cup D_T$ 
    2. Shuffle  $D$ 
    for  $b_t$  in  $D$  do
        // $b_t$  is a mini-batch of task  $t$ .
    3. Compute loss :  $L(\Theta)$ 
        $L(\Theta) = \text{Eq. 6}$  for classification
        $L(\Theta) = \text{Eq. 7}$  for regression
        $L(\Theta) = \text{Eq. 8}$  for ranking
    4. Compute gradient:  $\nabla(\Theta)$ 
    5. Update model:  $\Theta = \Theta - \epsilon \nabla(\Theta)$ 
    end
end
  
```

Concluding Remarks

- ① Labeling data is expensive, but we have large unlabeled data
- ① Autoencoder
 - exploits unlabeled data to learn latent factors as representations
 - learned representations can be transfer to other tasks
- ① Distant Labels / Labels from Other Tasks
 - learn the representations that are useful for other tasks
 - learned representations may be also useful for the target task