*Applied Deep Learning*

# Unsupervised Learning

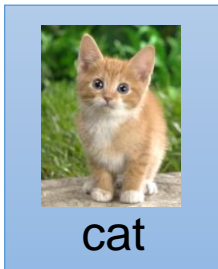**May 25th, 2020**  **http://adl.miulab.tw**

National Taiwan University

# **Introduction**

◉ Big data ≠ Big annotated data

◉ Machine learning techniques include:
  ○ Supervised learning (if we have labelled data)
  ○ Reinforcement learning (if we have an environment for reward)
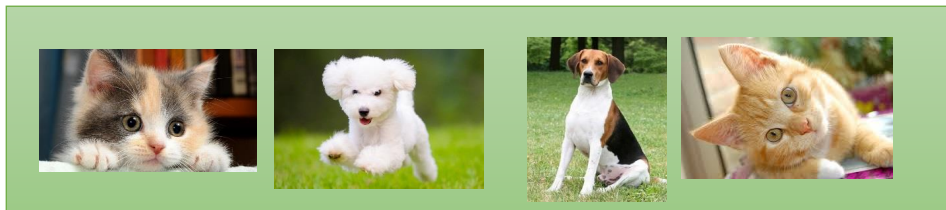  ○ <span style="color:red">Unsupervised learning (if we do not have labelled data)</span>

What can we do if there is no sufficient training data?
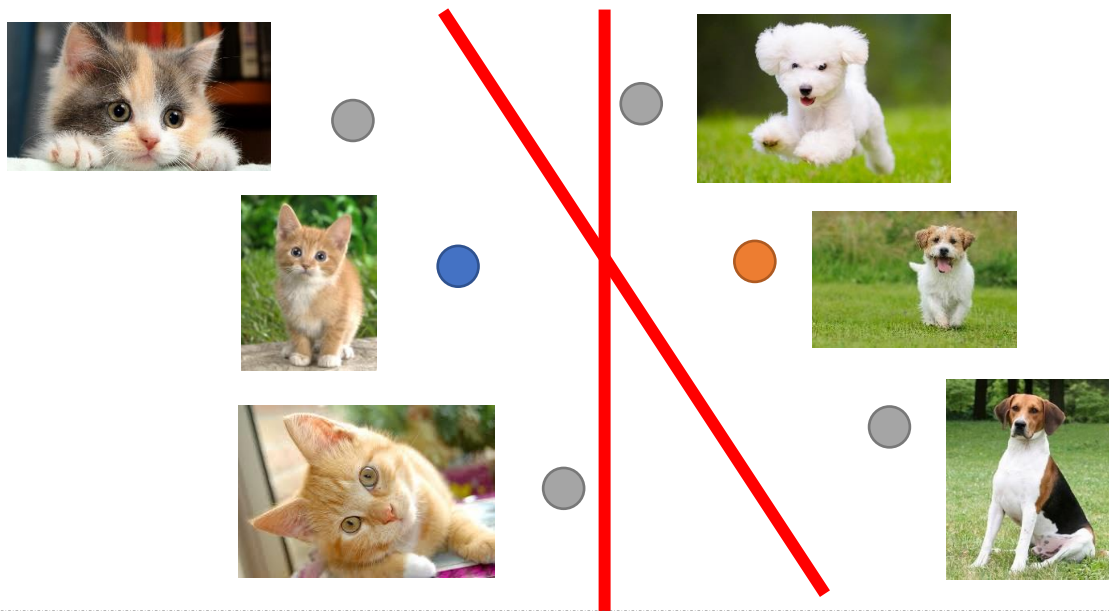
# Semi-Supervised Learning

Labelled Data


cat


dog

Unlabeled Data



(Image of cats and dogs without labeling)
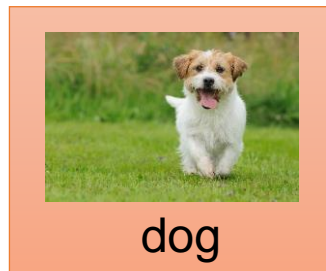
# Semi-Supervised Learning

◎ Why semi-supervised learning helps?



The distribution of the unlabeled data provides some cues

# Transfer Learning

Source Data

cat

dog

Target Data

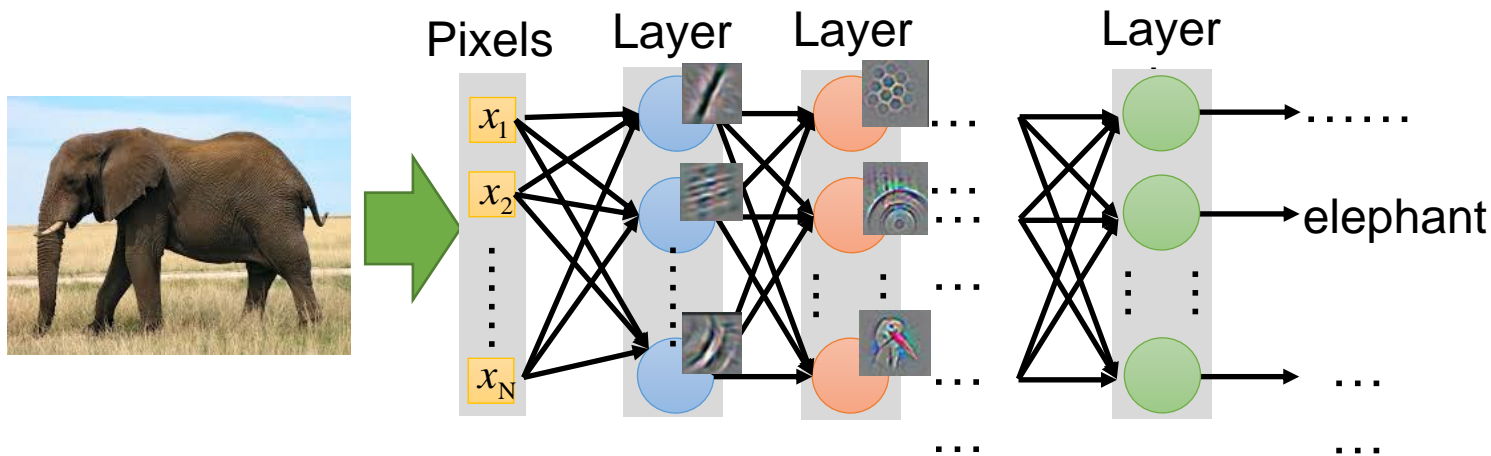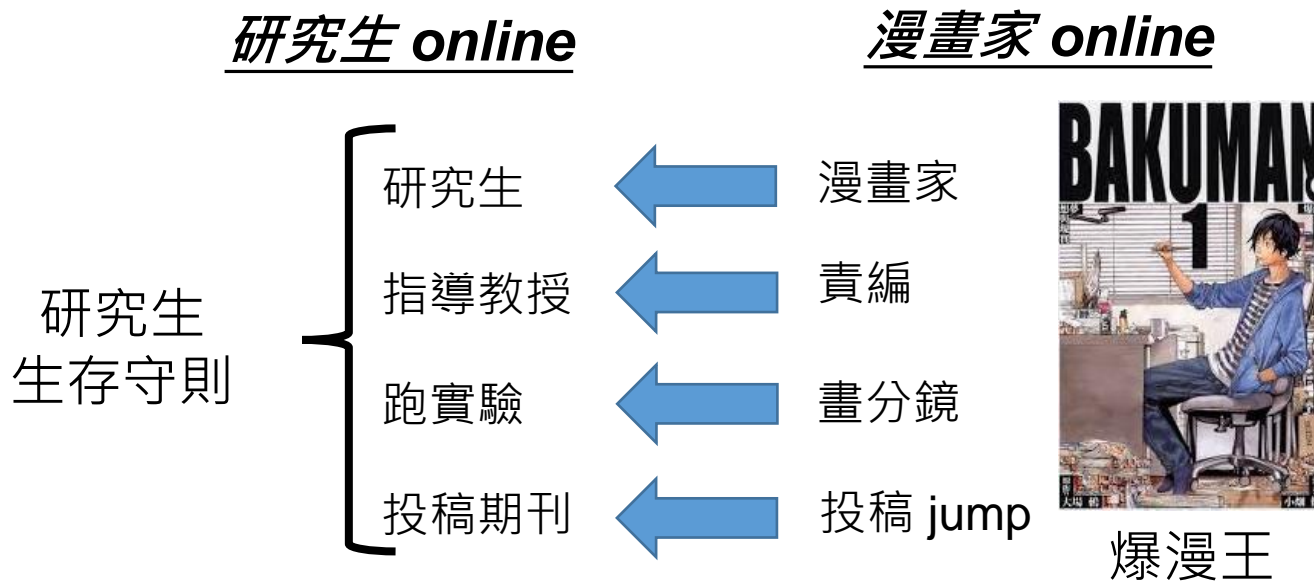elephant          elephant          tiger          tiger

Not related to the task considered

# Transfer Learning

◉ Widely used on image processing
  ○ Using sufficient labeled data to learn a CNN
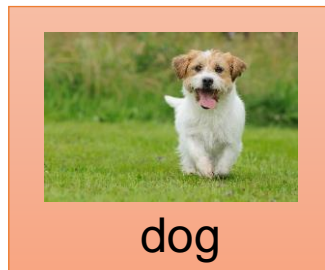  ○ Using this CNN as feature extractor
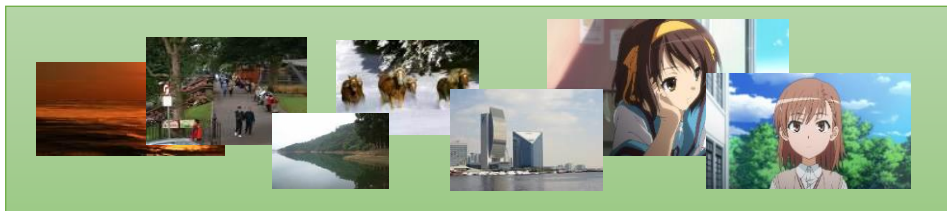
# Transfer Learning Example



研究生 online　　　　　漫畫家 online

研究生
生存守則

研究生 ← 漫畫家

指導教授 ← 責編

跑實驗 ← 畫分鏡

投稿期刊 ← 投稿 jump

爆漫王

# Self-Taught Learning

◉ The unlabeled data sometimes is not related to the task

Labelled Data


cat


dog

Unlabeled Data



(Just crawl millions of images from the Internet)

# **Self-Taught Learning**

◉ The unlabeled data sometimes is not related to the task

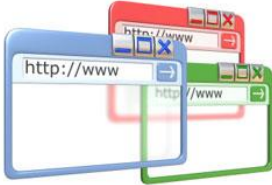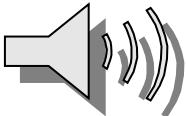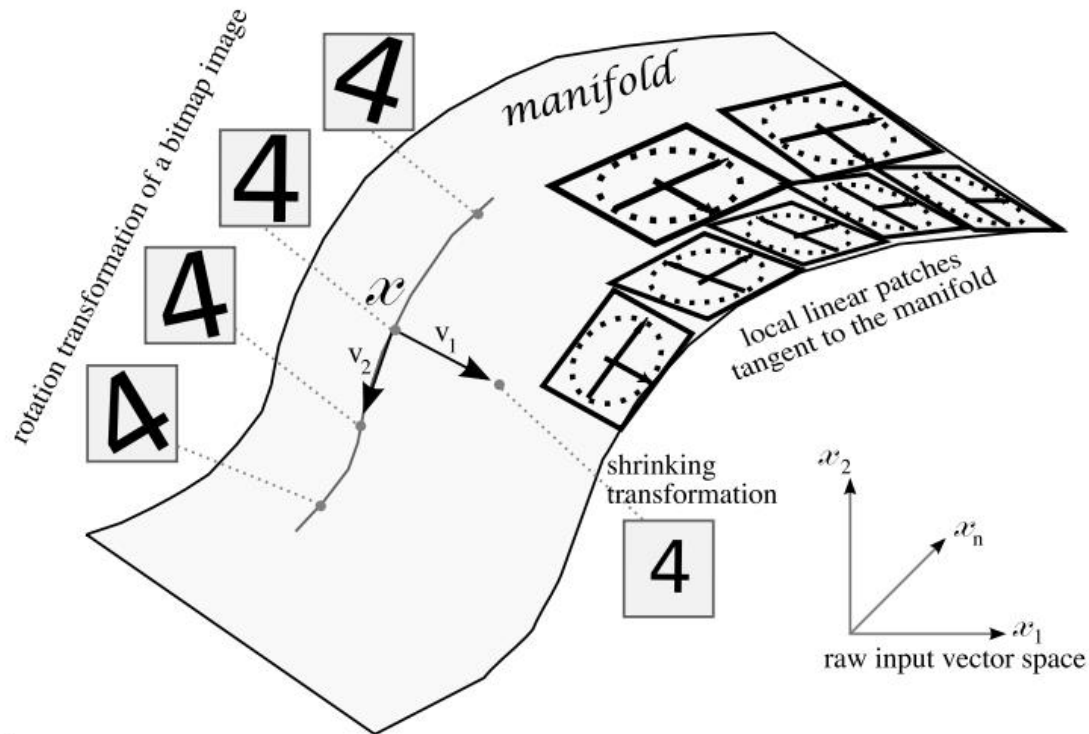|  | Labelled Data | Unlabeled Data |
|---|---|---|
| Digit Recognition |  Digits |  character |
| Document Classification |  News |  Webpages |
| Speech Recognition |  Taiwanese |  English Chinese …… |

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

# Self-Taught Learning

◉ How does self-taught learning work?

◉ 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

# Latent Factor Exploitation

⊙ Handwritten digits



The handwritten images are composed of **strokes**

## *Strokes (Latent Factors)*



| No. 1 | No. 2 | No. 3 | No. 4 | No. 5 | ....... |

# Latent Factor Exploitation

**_Strokes (Latent Factors)_**

No. 1    No. 2    No. 3    No. 4    No. 5    .......

28

28    =    No. 1    +    No. 3    +    No. 5

Represented by
28 X 28 = 784 pixels

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

# Autoencoder

Representation Learning

# 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

28 X 28 = 784 → NN Encoder → *__code__*  compact representation of the input object

Usually <784

Learn together

*__code__* → NN Decoder → reconstruct the original object

# Autoencoder

Minimize $(x - y)^2$

As close as possible

encode

decode

$x$

$a$

$y$

$W$

$W'$

Input layer

hidden layer
Bottleneck layer

output layer

$$a = \sigma(Wx + b) \qquad y = \sigma(W'a + b')$$

Output of the hidden layer is the code

# Autoencoder

◉ De-noising auto-encoder



Rifai, et al. "Contractive auto-encoders: Explicit invariance during feature extraction," in *ICML*, 2011.

# Deep Autoencoder



As close as possible

Input Layer — Layer — Layer — ... — Layer — bottle — Layer — ... — Layer — Layer — Output Layer

$x$

Code

$\tilde{x}$

# Deep Autoencoder

# Feature Representation

# Auto-encoder – Text Retrieval

## *Vector Space Model*



query

document

## *Bag-of-word*

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



2

125

250

500

2000

Bag-of-word (document or query)

query

Interbank markets

European Community monetary/economic

Energy markets

Disasters and accidents

Leading economic indicators

Legal/judicial

Accounts/ earnings

Government borrowings

The documents talking about the same thing will have close code

# Auto-Encoding (AE)

◉ Objective: reconstructing $\bar{x}$ from $\hat{x}$

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

○ dimension reduction or denoising (masked LM)

# Autoencoder – Similar Image Retrieval

◉ Retrieved using Euclidean distance in pixel intensity space



dist: 0.0

Krizhevsky et al. "Using very deep autoencoders for content-based image retrieval," in *ESANN*, 2011.

# Autoencoder – Similar Image Retrieval



32x32 → 8192 → 4096 → 2048 → 1024 → 512 → 256 (code)

(crawl millions of images from the Internet)

# 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*

output  | 10 |

| 500 |

***Target***

| 1000 |

| 1000 |

Input | 784 |

| 1000 | $\tilde{a}^1$

$W^{2'}$

| 1000 |

$W^2$

| 1000 | $a^1$

**fix** $W^1$

Input | 784 | $x$

# Autoencoder for DNN Pre-Training

- Greedy layer-wise pre-training *again*



output: 10

500

**Target**

1000

1000

Input: 784

$\tilde{a}^2$: 1000

$W^{3'}$

500

$W^3$

$a^2$: 1000

**fix** $W^2$

$a^1$: 1000

**fix** $W^1$

Input: 784 $x$

# Autoencoder for DNN Pre-Training

● Greedy layer-wise pre-training *again*

Find-tune via backprop

# Variational Autoencoder

Representation Learning and Generation

# Generation from Latent Codes



encode $W$

decode $W'$

$x$   $a$   $y$

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

# Reconstruction

AE

VAE

# 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: $l_m$

Max Pooling Operation

Convolutional Layer: $l_c$

Convolution Matrix: $W_c$

Word Hashing Layer: $l_h$

Word Hashing Matrix: $W_h$

Word Sequence: $x$



$P(D_1 \mid Q)$ $P(D_2 \mid Q)$ $P(D_n \mid Q)$

$\text{CosSim}(Q, D_i)$

$Q$ Query

$D_1$ $D_2$ $\dots$ $D_n$

Documents

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

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

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

maximizes the likelihood of clicked documents given queries

Huang et al.,  "Learning deep structured semantic models for web search using clickthrough data," in *Proc. of CIKM*, 2013.
Shen et al., "Learning semantic representations using ´ convolutional neural networks for web search," in *Proc. of WWW*, 2014.

# 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



Slot Tagging Model — STM — SHARED — LM

**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).
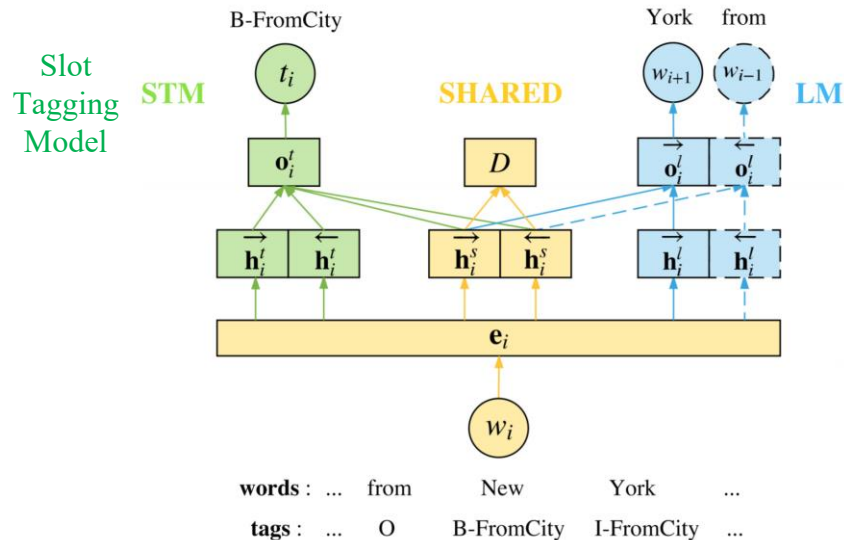6    Train the task discriminator and the shared model by Eq.(6) or Eq.(7) as LM task ($y = 0$).
7  **until** *convergence*;

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 $1, 2, ..., T$ **do**

 | Pack the dataset $t$ into mini-batch: $D_t$.

**end**

**for** $epoch$ in $1, 2, ..., epoch_{max}$ **do**

 1. Merge all the datasets:

  $D = D_1 \cup D_2... \cup D_T$

 2. Shuffle $D$

 **for** $b_t$ in $D$ **do**

  $//b_t\ is\ a\ mini\text{-}batch\ of\ task\ t.$

  3. Compute loss : $L(\Theta)$

   $L(\Theta) = $ Eq. 6 for classification

   $L(\Theta) = $ Eq. 7 for regression

   $L(\Theta) = $ Eq. 8 for ranking

  4. Compute gradient: $\nabla(\Theta)$

  5. Update model: $\Theta = \Theta - \epsilon\nabla(\Theta)$

 **end**

**end**

https://github.com/namisan/mt-dnn

# **Concluding Remarks**

◉ Labeling data is expensive, but we have large unlabeled data

◉ Autoencoder
  ○ exploits unlabeled data to learn latent factors as representations
  ○ learned representations can be transfer to other tasks

◉ 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