# FOUNDATIONS OF NATURAL LANGUAGE PROCESSING
# 自然語言處理的基礎

張家瑋 博士
國立臺中科技大學資訊工程系助理教授

# NATURAL LANGUAGE PROCESSING
# 自然語言處理的原理與應用

# 自然語言處理的主要範疇

- 機器翻譯 (Machine Translation)
- 自然語言理解/語意分析 (Natural Language Understanding / Semantic Analysis)
  1. 問答系統 (Question Answering)
  2. 萃取式摘要 (Extractive Summarization)
  3. 文件分類 (Text Categorization)
- 自然語言生成 (Natural Language Generation)
  1. 進階問答系統 (Advanced Question Answering)
  2. 抽象式摘要 (Abstractive Summarization)
  3. 聊天機器人 (Chatbot)

- 語法分析 (Syntactic Parsing)
  1. 中文斷詞 (Chinese word segmentation)
  2. 詞性標註 (Part-of-speech Tagging)
  3. 實體辨識 (Named Entity Recognition)
  4. 詞彙依存 (Typed Dependencies)
  5. 文法樹 (Parse Tree)
- 語音辨識 (Speech Recognition)
- 文字轉語音 (Text to Speech)
- 語音轉文字 (Speech to Text)

# DATA PREPARATION

- Data preprocessing and cleaning
  - Preprocess data in order to reduce noise and handle missing values
  - 斷字, 斷詞
- Relevance analysis (feature selection)
  - Remove the irrelevant or redundant attributes
  - 移除stop words, 擷取有用資訊(TF-IDF)
- Data transformation
  - Generalize and/or normalize data
  - 轉成向量(Vector representation)

# SEGMENTATION

- Segment by word/ sentence
- Segment in English
  - In English, we can directly segment the word by space " "
  - Ex: I love machine learning. → [I, love, machine, learning]
- Segment in Chinese
  - In Chinese, we segment the word by meaningful word rather than directly segment by characters.
  - Ex: 我喜歡機器學習 → [我, 喜歡, 機器, 學習]
    
    rather than [我, 喜, 歡, 機, 器, 學, 習]

根據字詞結構將一句話斷字

Dear 小明,
這是目前公司的最新技術，利用
apples 和 pens 的特性可以讓產能最
佳化............

Dear, 小明, 這是, 目前, 公司, 的, 最新,
技術, 利用, apples, 和, pens, 的, 特性,
可以, 讓, 產能, 最佳化, ............

# REMOVING STOP WORDS

- Remove the word which is meaningless.

- Usually do after segment.

- Remove stop words in Chinese

  - Example of stop words: 的, 了, 且, 個, 是

  - Ex: 今天的空氣品質不好 → [今天, 空氣, 品質, 不好 ]

- Remove stop words in English

  - Example of stop words: is, the, an, and, a

  - Ex: Today‘s air quality is not good → [Today's, air, quality, not, good]

Dear, 小明, 這是, 目前, 公司, 的, 最新, 技術, 利用, apples, 和, pens, 的, 特性, 可以, 讓, 產能, 最佳化, ............

移除stop-word

小明, 目前, 最新, 技術, 利用, apples, pens, 特性, 產能, 最佳化, ............

# STEMMING

- Stemming is to transform the word into its original type by removing word endings such as -s , -ed and -ing.
  - "bikes"   is replaced with   "bike"   ,
  - "raining"   is replaced with   "rain"
  - "tried"   is replaced with   "try"

小明, 目前, 最新, 技術, 利用, apples, pens, 特性, 產能, 最佳化, …………

stemming

小明, 目前, 最新, 技術, 利用, apple, pen, 特性, 產能, 最佳化, …………

# REPRESENTATION

- Select features from the data
- Transform data into vector model

- Ex)
  - WordNet
  - TF-IDF (Term Frequency - Inverse Document Frequency)
  - Word2Vec

# WORD-SENSE DISAMBIGUATION(1/2)

- Ambiguity: a word or phrase with multiple meanings.
  1. "procure"   (I will get the drinks)
  2. "become"   (she got scared)
  3. "have"   (I have got three dollars)
  4. "understand"   (I get it)

# WORDNET



http://wordnetweb.princeton.edu/perl/webwn

Distance-based: PATH (Rada, Mili, Bicknell, & Blettner, 1989)

# WORD-SENSE DISAMBIGUATION(2/2)

Information Content-based: RES (Resnik, 1995)

# WORD-SENSE DISAMBIGUATION(2/2)



- Gloss-based: VECTOR (Patwardhan, 2003)

| Cute | Cunning |
|---|---|
| 1. attractive especially by means of smallness or prettiness or quaintness | 1. attractive especially by means of smallness or prettiness or quaintness |
| 2. obviously contrived to charm | 2. marked by skill in deception |
|  | 3. showing inventiveness and skill |

# STANFORD PARSER



**Stanford Parser**

Please enter a sentence to be parsed:

My dog also likes eating sausage.

Language: English ♦   Sample Sentence                    Parse

**Your query**

*My dog also likes eating sausage.*

**Tagging**

My/PRP$ dog/NN also/RB likes/VBZ eating/VBG sausage/NN ./.

**Parse**

```
(ROOT
  (S
    (NP (PRP$ My) (NN dog))
    (ADVP (RB also))
    (VP (VBZ likes)
      (S
        (VP (VBG eating)
          (NP (NN sausage)))))
    (. .)))
```

**Universal dependencies**

```
nmod:poss(dog-2, My-1)
nsubj(likes-4, dog-2)
advmod(likes-4, also-3)
root(ROOT-0, likes-4)
xcomp(likes-4, eating-5)
dobj(eating-5, sausage-6)
```



The Stanford Natural Language Processing Group

people   publications   research blog   software   teaching   local

Software > Stanford Parser

## The Stanford Parser: A statistical parser

About | Citing | Questions | Download | Included Tools | Extensions | Release history | Sample output | Online | FAQ

### About

A natural language parser is a program that works out the grammatical **structure of sentences**, for instance, which groups of words go together (as "phrases") and which words are the **subject** or **object** of a verb. Probabilistic parsers use knowledge of language gained from hand-parsed sentences to try to produce the *most likely* analysis of new sentences. These statistical parsers still make some mistakes, but commonly work rather well. Their development was one of the biggest breakthroughs in natural language processing in the 1990s. You can try out our parser online.

### Package contents

This package is a Java implementation of probabilistic natural language parsers, both highly optimized PCFG and lexicalized dependency parsers, and a lexicalized PCFG parser. The original version of this parser was mainly written by Dan Klein, with support code and linguistic grammar development by Christopher Manning. Extensive additional work (internationalization and language-specific modeling, flexible input/output, grammar compaction, lattice parsing, *k*-best parsing, typed dependencies output, user support, etc.) has been done by Roger Levy, Christopher Manning, Teg Grenager, Galen Andrew, Marie-Catherine de Marneffe, Bill MacCartney, Anna Rafferty, Spence Green, Huihsin Tseng, Pi-Chuan Chang, Wolfgang Maier, and Jenny Finkel.

The lexicalized probabilistic parser implements a factored product model, with separate PCFG phrase structure and lexical dependency experts, whose preferences are combined by efficient exact inference, using an A* algorithm. Or the software can be used simply as an accurate unlexicalized stochastic context-free grammar parser. Either of these yields a good performance statistical parsing system. A GUI is provided for viewing the phrase structure tree output of the parser.

18

https://nlp.stanford.edu/software/lex-parser.shtml#Sample

# PARSE TREE

# 廣義知網知識本體



http://ehownet.iis.sinica.edu.tw/ehownet.php

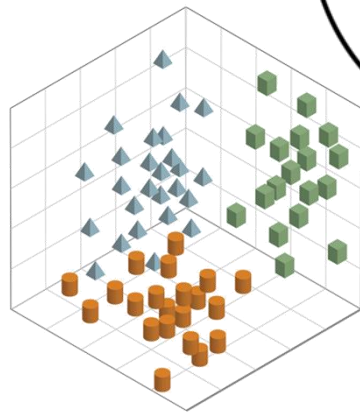# SEMANTIC SIMILARITY MEASURES

文字檔案

word2vec

將被拆解成多個字元

透過向量比對
找出相似的資料

Input:
one document

*Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et*

word vectors

Model:

kite

space

dog

netherlands
france
spain
belgium
italy

water

house

vector space

解析成多元維度的向量

most_similar('france'):

| spain | 0.678515 |
| belgium | 0.665923 |
| netherlands | 0.652428 |
| italy | 0.633130 |

highest cosine
distance values
in vector space
of the nearest
words

# VECTOR REPRESENTATION

| | $w_1$ | $w_2$ | $w_3$ | .. | .. | .. | $w_{n-1}$ | $w_n$ | label |
|---|---|---|---|---|---|---|---|---|---|
| $D_1$ | 0.11 | 0.23 | 0 | .. | .. | .. | 0.57 | 0 | 0 |
| $D_2$ | 0 | 0 | 0 | .. | .. | .. | 0.29 | 0.7 | 1 |
| $D_3$ | 0 | 0.81 | 0.44 | .. | .. | .. | 0 | 0 | 0 |
| $D_4$ | 0 | 0.37 | 0 | .. | .. | .. | 0 | 0.16 | 1 |
| .. | .. | .. | .. | .. | .. | .. | .. | .. | .. |
| $D_k$ | .. | .. | .. | .. | .. | .. | .. | .. | 1 |

Machine learning

# TF-IDF

# TF-IDF

- TF: term frequency: $\quad tf_{i,j} = \dfrac{n_{i,j}}{\sum_k n_{k,j}}$

- IDF: inverse document frequency: $\quad idf_i = \log \dfrac{|D|}{|\{j : t_i \in d_j\}|}$

where:

- |D|: total number of documents in the corpus
- $|\{j : t_i \in d_j\}|$ : number of documents where term $t_i$ appears

Then:

- $tfidf_{i,j} = tf_{i,j} \times idf_i$

| Document 1 | | Document 2 | |
|---|---|---|---|
| **Term** | **Term Count** | **Term** | **Term Count** |
| this | 1 | this | 1 |
| is | 1 | is | 1 |
| a | 2 | another | 2 |
| sample | 1 | example | 3 |

- The calculation of tf–idf for the term "this" is performed as follows:

  - 

$$tf("this", d_1) = \frac{1}{5} = 0.2$$
$$tf("this", d_2) = \frac{1}{7} \approx 0.14$$

$$idf("this", D) = \log\left(\frac{2}{2}\right) = 0$$

- So tf–idf is zero for the word "this", which implies that the word is not very informative as it appears in all documents.

$$tfidf("this", d_1) = 0.2 \times 0 = 0$$
$$tfidf("this", d_2) = 0.14 \times 0 = 0$$

26

| Document 1 | | Document 2 | |
| --- | --- | --- | --- |
| Term | Term Count | Term | Term Count |
| this | 1 | this | 1 |
| is | 1 | is | 1 |
| a | 2 | another | 2 |
| sample | 1 | example | 3 |

- A slightly more interesting example arises from the word "example", which occurs three times only in the second document:

- $$\text{tf}(''\text{example}'', d_1) = \frac{0}{5} = 0$$
  $$\text{idf}(''\text{example}'', D) = \log\left(\frac{2}{1}\right) = 0.301$$

- $$\text{tf}(''\text{example}'', d_2) = \frac{3}{7} \approx 0.429$$

$$\text{tfidf}(''\text{example}'', d_1) = \text{tf}(''\text{example}'', d_1) \times \text{idf}(''\text{example}'', D) = 0 \times 0.301 = 0$$
$$\text{tfidf}(''\text{example}'', d_2) = \text{tf}(''\text{example}'', d_2) \times \text{idf}(''\text{example}'', D) = 0.429 \times 0.301 \approx 0.13$$

# 潛藏語意分析(LSA)

- 奇異值分解
  - Singular Value Decomposition (SVD)

# 潛藏語意分析(LSA)

- 文件分類/主題探勘
- 語意分析



| Index Words | Titles | | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | T1 | T2 | T3 | T4 | T5 | T6 | T7 | T8 | T9 |
| book | | | 1 | 1 | | | | | |
| dads | | | | | | 1 | | | 1 |
| dummies | | 1 | | | | | | 1 | |
| estate | | | | | | | 1 | | 1 |
| guide | 1 | | | | | 1 | | | |
| investing | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| market | 1 | | 1 | | | | | | |
| real | | | | | | | 1 | | 1 |
| rich | | | | | | 2 | | | 1 |
| stock | 1 | | 1 | | | | | 1 | |
| value | | | | 1 | 1 | | | | |

# 以語言學習輔助工具為例

# 以語言學習輔助工具為例

| | computer | data | pinch | result | sugar |
|---|---|---|---|---|---|
| aprocot | 0 | 0 | 1 | 0 | 1 |
| pineapple | 0 | 0 | 1 | 0 | 1 |
| digital | 2 | 1 | 0 | 1 | 0 |
| information | 1 | 6 | 0 | 4 | 0 |

$$P(x = information, y = data) = \frac{6}{19} = 0.32$$

$$P(x = information) = \frac{6 + 4 + 1}{19} = \frac{11}{19} = 0.58$$

$$P(y = data) = \frac{6 + 1}{19} = \frac{7}{19} = 0.37$$

$$pmi(x = information, y = data)$$
$$= log \frac{P(x = information, y = data)}{P(x = information) \times P(y = data)}$$
$$= log 1.49$$
$$= 0.57$$

# WORD2VEC

# One Hot Encoding

```
The -> [0. 1. 0. 0. 0. 0. 0. 0. 0. 0.]
cat -> [0. 0. 1. 0. 0. 0. 0. 0. 0. 0.]
jump -> [0. 0. 0. 1. 0. 0. 0. 0. 0. 0.]
over -> [0. 0. 0. 0. 1. 0. 0. 0. 0. 0.]
the -> [0. 0. 0. 0. 0. 1. 0. 0. 0. 0.]
dog -> [0. 0. 0. 0. 0. 0. 1. 0. 0. 0.]
The -> [0. 1. 0. 0. 0. 0. 0. 0. 0. 0.]
dog -> [0. 0. 0. 0. 0. 0. 1. 0. 0. 0.]
ate -> [0. 0. 0. 0. 0. 0. 0. 1. 0. 0.]
my -> [0. 0. 0. 0. 0. 0. 0. 0. 1. 0.]
homework -> [0. 0. 0. 0. 0. 0. 0. 0. 0. 1.]
```
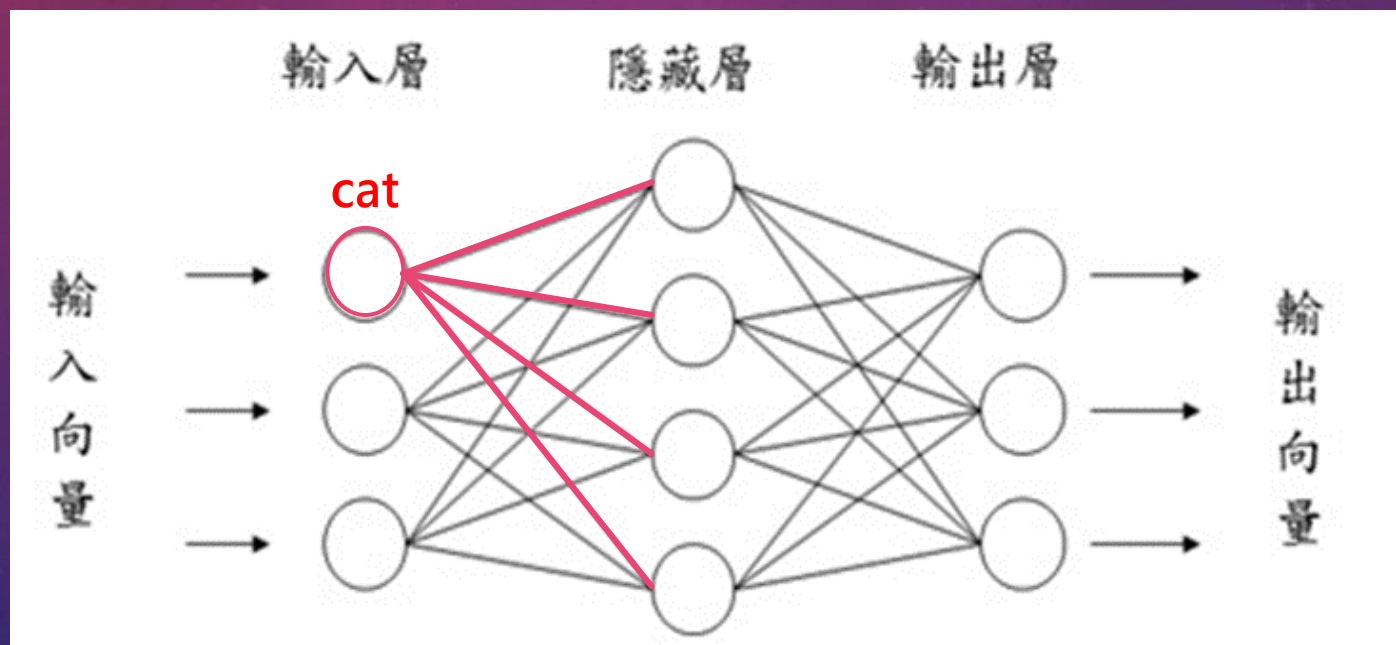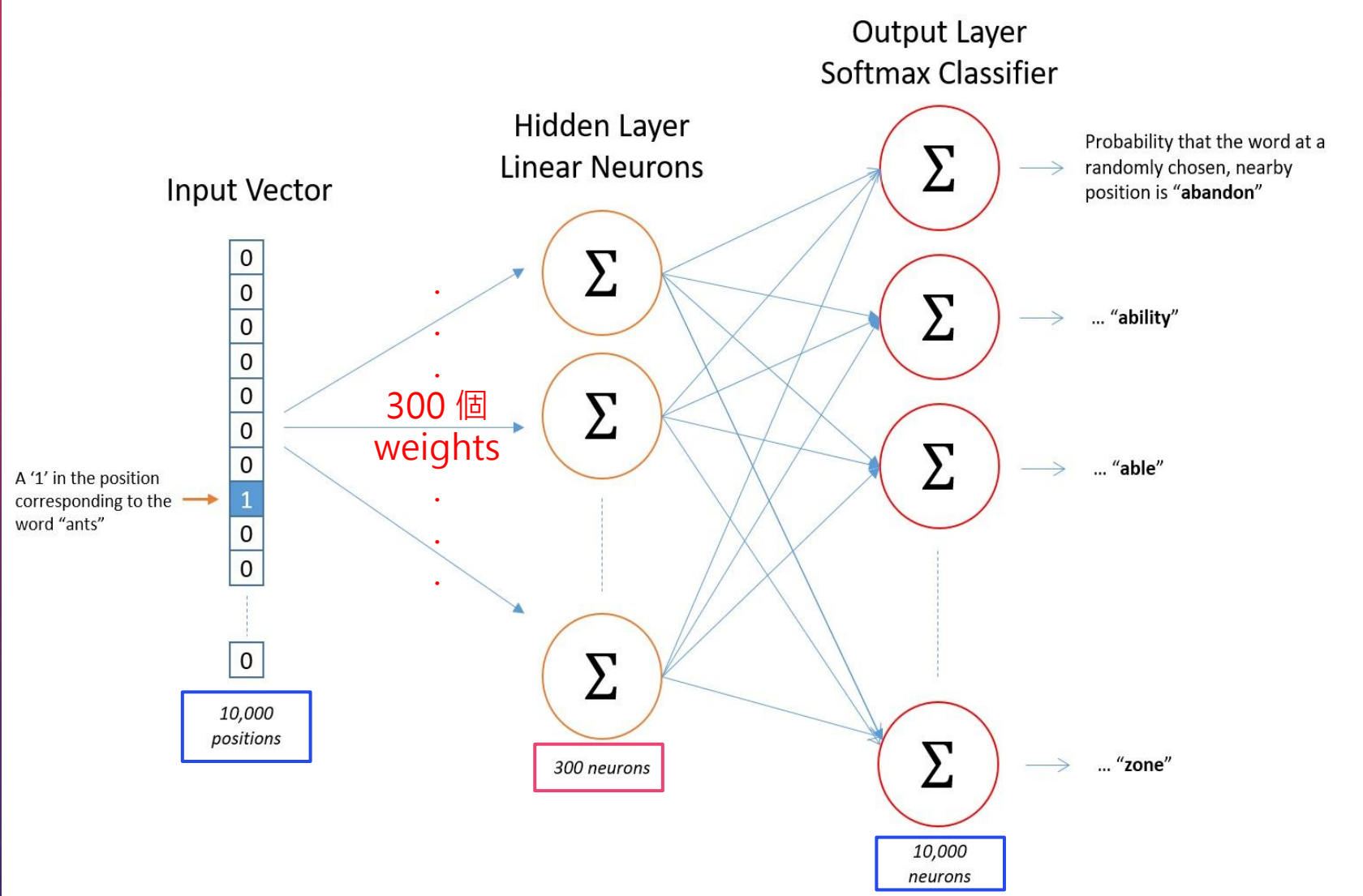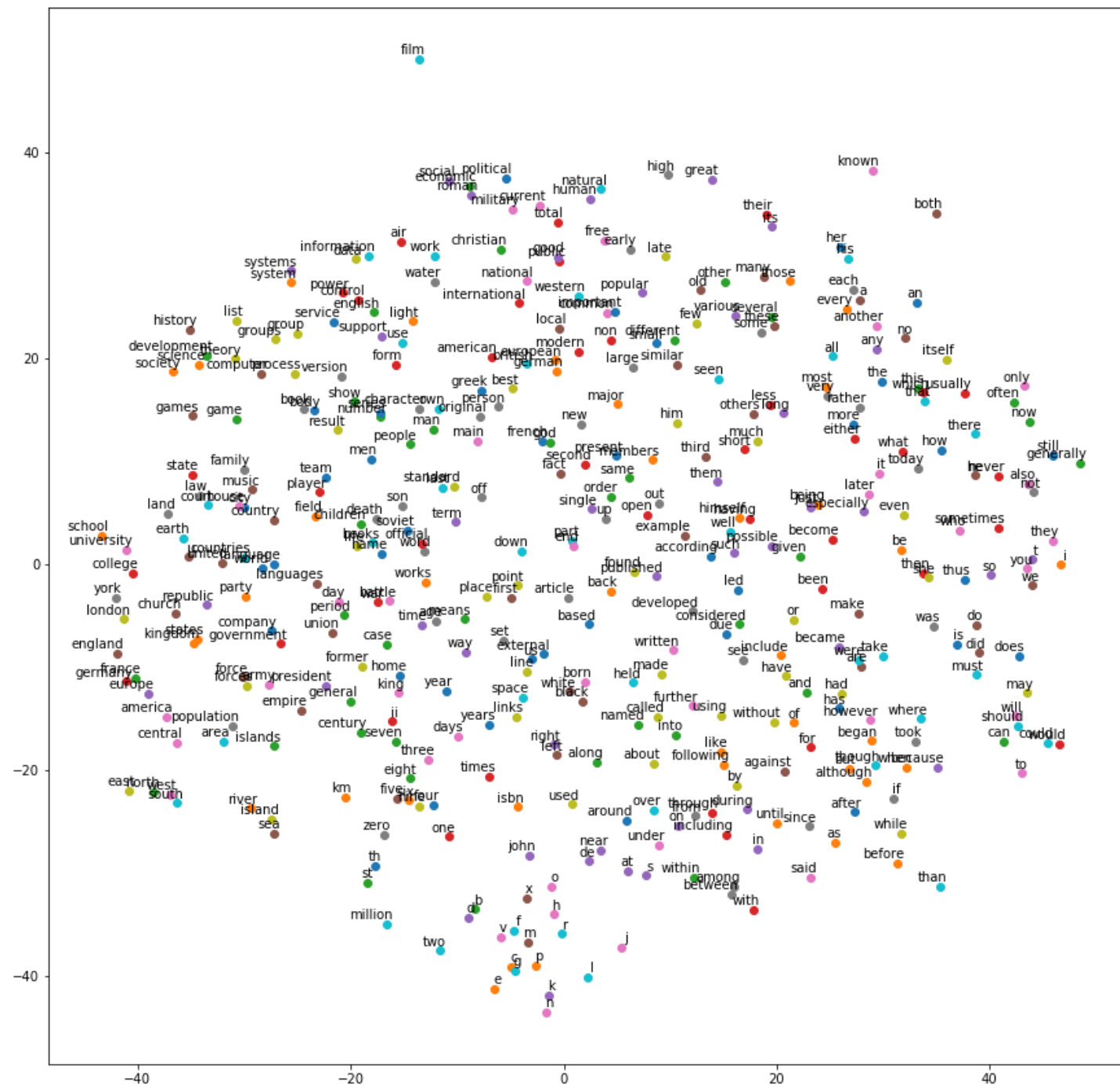
# Perceptron Linear Algorithm

$h(x) > 0$

$h(x) = 0$

$x_2$

$h(x) < 0$

$x_1$

- Features: $x = (x_1, x_2)$
- Target: $y = +1$ or $-1$
- $h(x) = w_0 + w_1 x_1 + w_2 x_2$



35

# Perceptron Linear Algorithm

h(x) > 0

h(x) = 0

$x_2$

h(x) < 0

$x_1$

$$h(x) = w_0 + w_1x_1 + w_2x_2$$

$$scores = \sum_{i}^{N} w_i x_i + b$$

$$scores = \sum_{i}^{N+1} w_i x_i$$

- 若 $scores \geq 0$ , 则 $\hat{y} = 1$
- 若 $scores < 0$ , 则 $\hat{y} = -1$

# Perceptron Linear Algorithm



- 若 $scores \geq 0$ ，則 $\hat{y} = 1$
- 若 $scores < 0$ ，則 $\hat{y} = -1$

$$w_{t+1} = w_t + y_t x_t$$

$$w_{t+1} = w_t + y_t x_t$$

[Case 1]
y = 1 錯分成 y = -1

[Case 2]
y = -1 錯分成 y = 1

# Perceptron Linear Algorithm

h(x) > 0

$x_2$

h(x) = 0

h(x) < 0

$x_1$

- 若 $scores \geq 0$ ，則 $\hat{y} = 1$
- 若 $scores < 0$ ，則 $\hat{y} = -1$

ya...正名為+1啦!

**+    -    +**

$w_{t+1} = w_t + y_t x_t$

修正

[Case 1]
y = 1 錯分成 y = -1

ya...正名為-1啦!

**-    +    -**

$w_{t+1} = w_t + y_t x_t$

修正

[Case 2]
y = -1 錯分成 y = 1

$$\tan \alpha = \lim_{\Delta x \to 0} \tan \varphi = \lim_{\Delta x \to 0} \frac{f(x_0 + \Delta x) - f(x_0)}{\Delta x}$$

線性關係

代價函數為凸函數
初始值隨機選也能降到全域最小值

$J(w)$

初始權重

代價函數
最小值

$J_{min}(w)$

$w$

# Multi-Layer Perceptron (MLP)



*線性組合 $w = a_1 v_1 + a_2 v_2 + a_3 v_3 + \cdots + a_n v_n$

# 梯度消失

\*現實世界的資料多為非線性，因此激活函數通常也是使用非線性函數(非凸函數)傳遞

PLA

MLP

```
The -> [0. 1. 0. 0. 0. 0. 0. 0. 0. 0.]
cat -> [0. 0. 1. 0. 0. 0. 0. 0. 0. 0.]
jump -> [0. 0. 0. 1. 0. 0. 0. 0. 0. 0.]
over -> [0. 0. 0. 0. 1. 0. 0. 0. 0. 0.]
the -> [0. 0. 0. 0. 0. 1. 0. 0. 0. 0.]
dog -> [0. 0. 0. 0. 0. 0. 1. 0. 0. 0.]
The -> [0. 1. 0. 0. 0. 0. 0. 0. 0. 0.]
dog -> [0. 0. 0. 0. 0. 1. 0. 0. 0. 0.]
ate -> [0. 0. 0. 0. 0. 0. 0. 1. 0. 0.]
my -> [0. 0. 0. 0. 0. 0. 0. 0. 1. 0.]
homework -> [0. 0. 0. 0. 0. 0. 0. 0. 0. 1.]
```



輸入層　　　隱藏層　　　輸出層

cat

輸入向量　　　　　　　　　　　　輸出向量

# 機器翻譯
# MACHINE TRANSLATION

# GOOGLE 翻譯

# BING 翻譯



英文 (已偵測) ▼ ⇄ 繁體中文 ▼ 英文 義大利文

My dog also likes eating sausage.
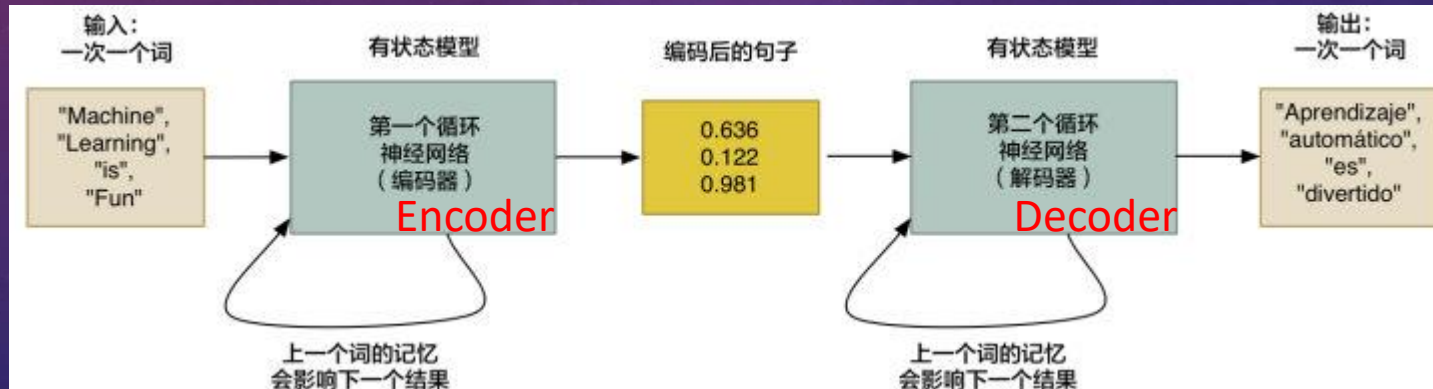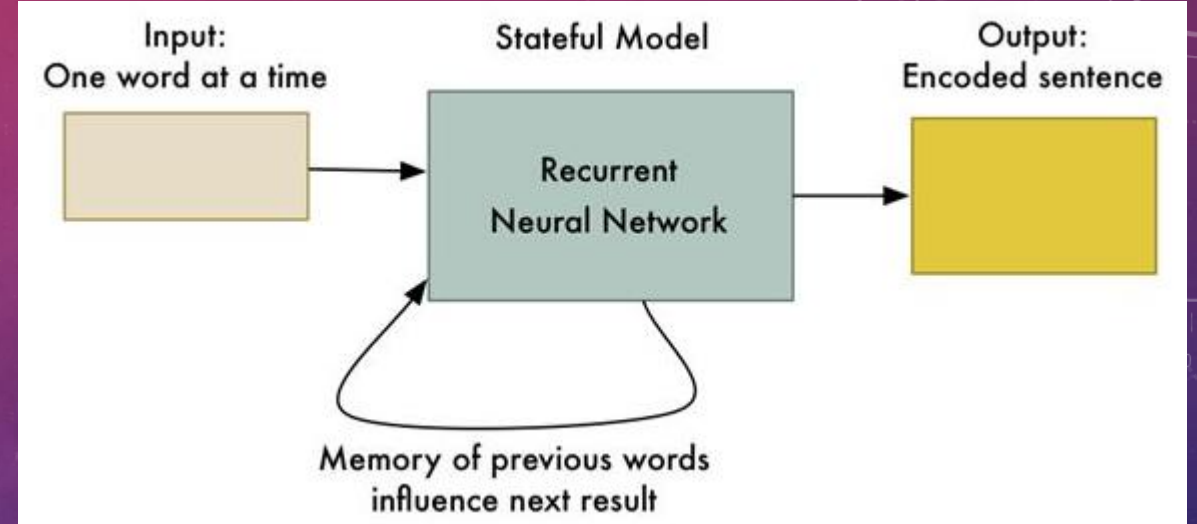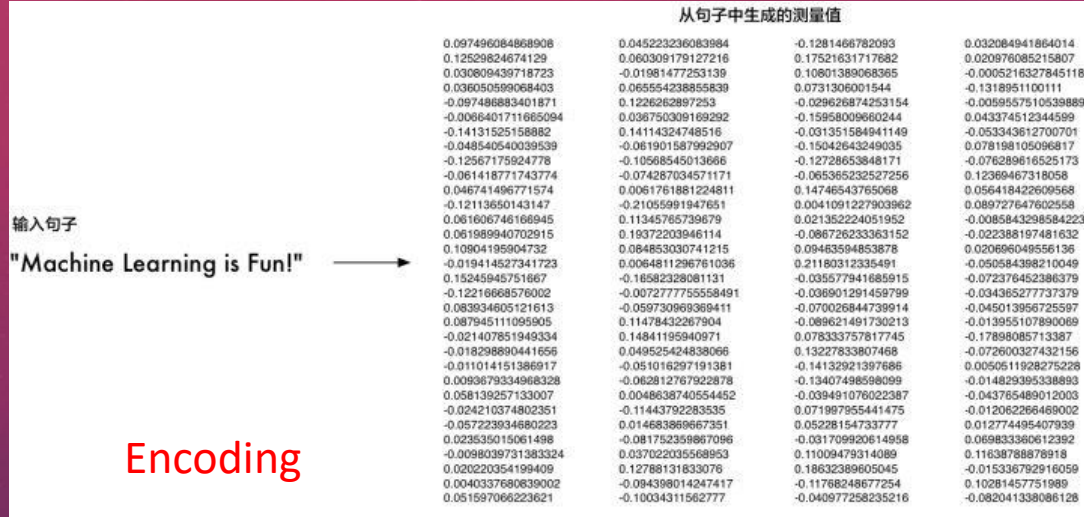
33/5000

我的狗也喜歡吃香腸。

wǒ de gǒu yě xǐ huān chī xiāng cháng.

# 有道翻譯

# 平行語料

# 統計式機器翻譯之原理

# 深度學習於機器翻譯之原理

# 深度學習於機器翻譯之原理



Encoding

# THANK YOU