

# Lecture 5 : Sequence Tagging

楊立偉教授  
台灣大學工管系

wyang@ntu.edu.tw

# Class Topics We Have

---

- L1: Document Similarity
  - "a Bag of Words" Model, Term Weighting (TF-IDF,  $\chi^2$ , MI)
  - Vector Space Model (VSM)
- L2: Co-occurrence
  - Association
  - Link analysis : Co-citation & Coupling
- L3: Classification
  - Naïve Bayes, k Nearest Neighbors, Support Vector Machine
  - Decision Tree, Bagging and Boosting, Random Forest
  - Neural Networks

- L4: Clustering
  - k Means, DBSCAN, Hierarchical Agglomerative Clustering
  - Topic Modeling
- L5: Sequence Tagging
  - Language model, HMM, CRF, RNN
- L6: Language Processing
  - Chinese Processing and other language issues
  - Word Embedding : LSA, Word2Vec, GloVe
  - Word window classification

# Sequence Tagging

---

- To assign of a label to each member of a sequence of observed values. For example:
  - part of speech tagging and voice recognition in language processing 語意分析中的詞性標記及語音辨識
  - Applications of sequence analysis or prediction in finance / bioinformatics 財務或生物資訊上的序列分析應用
- This can be done...
  - as ***a set of independent classification tasks***. For example, step forward one per member of the sequence a time.
  - by making the optimal label for a given element dependent on the choices of ***nearby*** elements.

# Predict the next element ?

- Try to predict the next element in a sequence of data

current (and previous ones)	candidates of the next one
馬	(馬)上、(馬)虎、(馬)腳、(馬)達、(馬)屁、(馬)克...
馬 英	(馬英)九

# Example (1)

- "美國是個自由的國家"
  - moving a cursor from left to right, use the last  $n$  words to tag the current one.
  - as a classification task, for example, to Decision Tree or SVM.

$W_{i-4}$	$W_{i-3}$	$W_{i-2}$	$W_{i-1}$	$W_i$
美	國	是	個	自
國	是	個	自	由
是	個	自	由	的
個	自	由	的	國
自	由	的	國	家

Use the last 4 words to tag the current one

- "美國是個自由的國家"

$W_{i-4}$	$W_{i-3}$	$W_{i-2}$	$W_{i-1}$	$W_i$
美	國	是	個	自
國	是	個	自	由
是	個	自	由	的
個	自	由	的	國
自	由	的	國	家

forward predicting ↑

$W_{i+1}$	$W_{i+2}$	$W_{i+3}$	$W_{i+4}$	$W_i$
國	是	個	自	美
是	個	自	由	國
個	自	由	的	是
自	由	的	國	個
由	的	國	家	自

↑ backward predicting

$W_{i-2}$	$W_{i-1}$	$W_i$	$W_{i+1}$	$W_{i+2}$
美	國	是	個	自
國	是	個	自	由
是	個	自	由	的
個	自	由	的	國
自	由	的	國	家

← bidirectional predicting

# Example (1)

- "美國是個自由的國家"
  - to obtain probability  $P(W_i | W_{i-1}, W_{i-2}, \dots, W_{i-n})$
  - for the sequence, when  $n=1$ , the probability is

$P(\text{"美國是個自由的國家"}) =$

$P(\text{國} | \text{美}) \times P(\text{是} | \text{國}) \times P(\text{個} | \text{是}) \times P(\text{自} | \text{個}) \times P(\text{由} | \text{自}) \times P(\text{的} | \text{由}) \times P(\text{國} | \text{的}) \times P(\text{家} | \text{國})$

\* instead of product of the probabilities, sum of log is used usually in programming.



# Example (1)

---

- "美國是個自由的國家"
  - in voice recognition (語音識別) and some intelligent input methods (輸入法), the probabilities of various candidates are evaluated, and the highest one is chosen.

$P(\text{"美國是個自由的國家"})=0.08$  ← chosen, for example

$P(\text{"美國是個製油的國家"})=0.07$

$P(\text{"美國似個自由的國家"})=0.05$

$P(\text{"美國事個自遊的國家"})=0.04$

# Variant (1)

- Larger size of the sliding window

for window size=2

$P(\text{"美國是個自由的國家"}) =$

$$\begin{aligned}
 &P(\text{是} | W_{i-1}=\text{國} \cap W_{i-2}=\text{美}) \times P(\text{個} | W_{i-1}=\text{是} \cap W_{i-2}=\text{國}) \times P(\text{自} | W_{i-1}=\text{個} \cap W_{i-2}=\text{是}) \times \\
 &P(\text{由} | W_{i-1}=\text{自} \cap W_{i-2}=\text{個}) \times P(\text{的} | W_{i-1}=\text{由} \cap W_{i-2}=\text{自}) \times P(\text{國} | W_{i-1}=\text{的} \cap W_{i-2}=\text{由}) \times \\
 &P(\text{家} | W_{i-1}=\text{國} \cap W_{i-2}=\text{的})
 \end{aligned}$$

可經由大量語料做次數統計而得

## Variant (2)

- Word/n-gram instead of Character

$P(\text{"美國是個自由的國家"}) =$

$$P(\text{國是} | \text{美國}) \times P(\text{是個} | \text{國是}) \times P(\text{個自} | \text{是個}) \times P(\text{自由} | \text{個自}) \times P(\text{由的} | \text{自由}) \times P(\text{的國} | \text{由的}) \times P(\text{國家} | \text{的國})$$

Sliding window size = 2 與 2-gram 哪個方法好？

→ 2-gram效果類似猜連續2字，組合數(類別)變多，需要的訓練資料要足夠

# Test the language model in your brain

- use the surrounding **unigrams** to predict

這□天□氣□冷□，所□都□到□較□才□門，  
不□夠□間□早□。

每個空格都有若干候選字

- use the surrounding **bigrams** to predict

這幾□天氣□冷了，所以□睡到□較晚□出門，  
不太□時間□早餐。

大幅減少空格中的候選字數

- one of the possible answers

這幾天天氣變冷了，所以都睡到比較晚才出門，  
不太夠時間吃早餐。



# Test the language model in your brain

- the window vs. the order

最近的研究表示，漢字序順並不一定影響閱讀。

- for auto-correction : can you read this ?

Aoccdrnig to a rscheearch at Cmabrigde Uinervtisy,  
it deosn't mttar in waht oredr the ltteers in a wrod  
are, the olny iprmoetnt tihng is taht the frist and lsat  
ltteer be at the rghit pclae. The rset can be a toatl  
mses and you can sitll raed it wouthit porbelm. Tihs  
is bcuseae the human mnid deos not raed ervey  
lteter by istlef, but the wrod as a wlohe.

From : <http://www.mrc-cbu.cam.ac.uk/people/matt.davis/Cmabrigde/>

# Test the language model in your brain

- ***"You shall know a word by the company it keeps"***  
(Firth, J. R. 1957:11)
- answer what "令和" is, after reading the following :

昨天日本公布新年號「令和」後，掀起一股新風潮。商家開始大搶「令和」商機，還有店家祭出，只要名字中有其中一個字，就有折價。  
日本民眾：「咦？是什麼？令和？！令和，令和元年。」  
日本學生：「新年號，令和。」

for example

2-gram	occurrences in bi-dir. window of 5
年號	2
商機	1
元年	1

## Example (2) for financial applications

- 若想利用股票今日價量尋找明日價格漲跌之關係，採用序列模型編碼如下

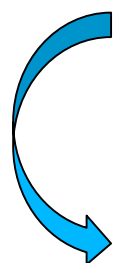
	Day <sub>1</sub>	Day <sub>2</sub>	Day <sub>3</sub>	Day <sub>4</sub>	Day <sub>5</sub>	Day <sub>6</sub>	Day <sub>7</sub>	Day <sub>8</sub>
價格	升	升	平	升	降	降	降	平
成交量	升	升	降	降	升	降	平	平

- 同樣可以使用分類、或連續的條件機率進行分析

# Example (2)

- 使用分類演算法來預測
  - 假設 $n=2$ ，編碼方式為 $D_{i-2\text{價}}$ 、 $D_{i-2\text{量}}$ 、 $D_{i-1\text{價}}$ 、 $D_{i-1\text{量}}$ ，預測 $D_{i\text{價}}$

	Day <sub>1</sub>	Day <sub>2</sub>	Day <sub>3</sub>	Day <sub>4</sub>	Day <sub>5</sub>	Day <sub>6</sub>	Day <sub>7</sub>	Day <sub>8</sub>	Day <sub>9</sub>
價格	升	升	平	升	降	降	降	平	?
成交量	升	升	降	降	升	降	平	平	



i	$D_{i-2\text{價}}$	$D_{i-2\text{量}}$	$D_{i-1\text{價}}$	$D_{i-1\text{量}}$	$D_{i\text{價}}$
3	升	升	升	升	平
4	升	升	平	降	升
5	平	降	升	降	降
6	升	降	降	升	降
7	降	升	降	降	降
8	降	降	降	平	平
9	降	平	平	平	?

forward predicting,  
for example



## Example (2)

- 使用連續的條件機率來預測
  - 假設  $n=2$ ，編碼方式為  $D_{i-2\text{價}}$ 、 $D_{i-2\text{量}}$ 、 $D_{i-1\text{價}}$ 、 $D_{i-1\text{量}}$ ，預測  $D_{i\text{價}}$

	Day <sub>1</sub>	Day <sub>2</sub>	Day <sub>3</sub>	Day <sub>4</sub>	Day <sub>5</sub>	Day <sub>6</sub>	Day <sub>7</sub>	Day <sub>8</sub>	Day <sub>9</sub>
價格	升	升	平	升	降	降	降	平	?
成交量	升	升	降	降	升	降	平	平	

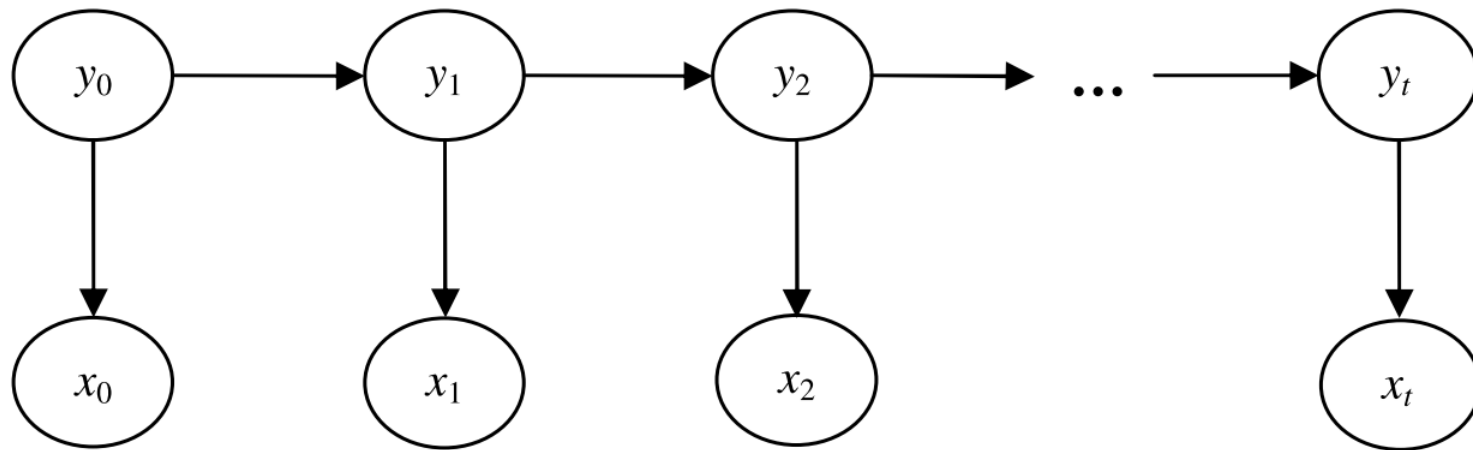
$$P(D_{i\text{價}}=\text{升} \mid D_{i-1\text{價}}=\text{平} \cap D_{i-1\text{量}}=\text{平} \cap D_{i-2\text{價}}=\text{降} \cap D_{i-2\text{量}}=\text{平})=?$$

$$P(D_{i\text{價}}=\text{平} \mid D_{i-1\text{價}}=\text{平} \cap D_{i-1\text{量}}=\text{平} \cap D_{i-2\text{價}}=\text{降} \cap D_{i-2\text{量}}=\text{平})=? \quad \text{三者取機率大者為猜測}$$

$$P(D_{i\text{價}}=\text{降} \mid D_{i-1\text{價}}=\text{平} \cap D_{i-1\text{量}}=\text{平} \cap D_{i-2\text{價}}=\text{降} \cap D_{i-2\text{量}}=\text{平})=?$$

將條件機率各項展開成可算之項後求(近似)解，比較大小

# Hidden Markov Model



**Fig. 3** Hidden Markov model

We have

$Y = \langle y_0, y_1, \dots, y_t \rangle =$  hidden state sequence

$X = \langle x_0, x_1, \dots, x_t \rangle =$  observation sequence

Ref. Lei Zhang and Bing Liu, "Aspect and Entity Extraction for Opinion Mining", 2014

HMM models a sequence of observations  $X$  by assuming that there is a *hidden* sequence of states  $Y$ . Observations are dependent on states. Each state has a probability distribution over the possible observations. To model the joint distribution  $p(y, x)$  tractably, two independence assumptions are made. First, it assumes that state  $y_t$  only depends on its immediate predecessor state  $y_{t-1}$ .  $y_t$  is independent of all its ancestor  $y_1, y_2, y_3, \dots, y_{t-2}$ . This is also called the *Markov* property. Second, the observation  $x_t$  only depends on the current state  $y_t$ . With these assumptions, we can specify HMM using three probability distributions:  $p(y_0)$  over initial state, state transition distribution  $p(y_t | y_{t-1})$  and observation distribution  $p(x_t | y_t)$ . That is, the joint probability of a state sequence  $Y$  and an observation sequence  $X$  factorizes as follows.

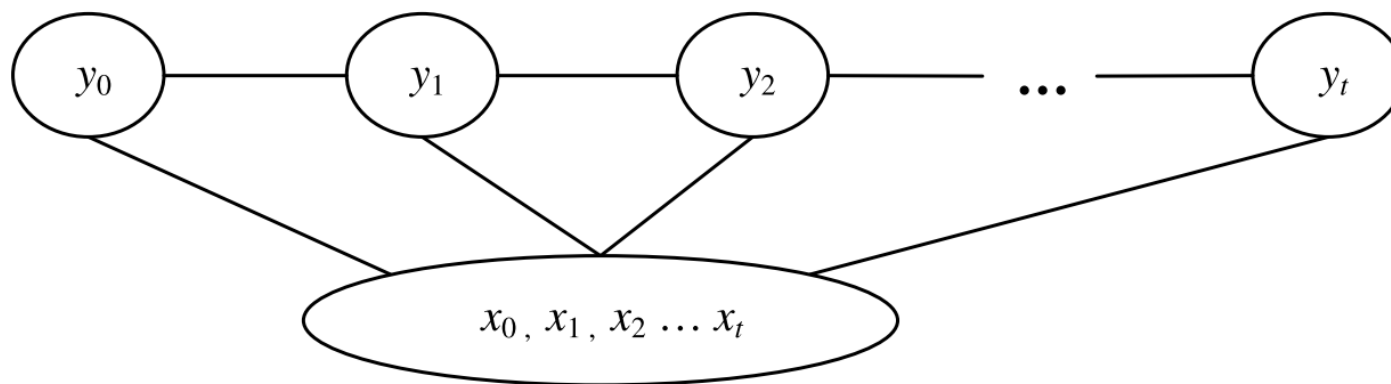
$$p(Y, X) = \prod_{t=1}^t p(y_t | y_{t-1}) p(x_t | y_t) \quad (2)$$

where we write the initial state distribution  $p(y_1)$  as  $p(y_1 | y_0)$ .

Given some observation sequences, we can learn the model parameter of HMM that maximizes the observation probability. That is, the learning of HMM can be done by building a model to best fit the training data. With the learned model, we can find an optimal state sequence for new observation sequences.

# Conditional Random Fields

One limitation of HMM is that its assumptions may not be adequate for real-life problems, which leads to reduced performance. To address the limitation, linear-chain Conditional Random fields (CRF) (Lafferty et al., 2001; Sutton and McCallum, 2006) is proposed as an undirected sequence model, which models a conditional probability  $p(Y|X)$  over hidden sequence  $Y$  given observation sequence  $X$ . That is, the conditional model is trained to label an unknown observation sequence  $X$  by selecting the hidden sequence  $Y$  which maximizes  $p(Y|X)$ . Thereby, the model allows relaxation of the strong independence assumptions made by HMM. The linear-chain CRF model is illustrated in Figure 4.



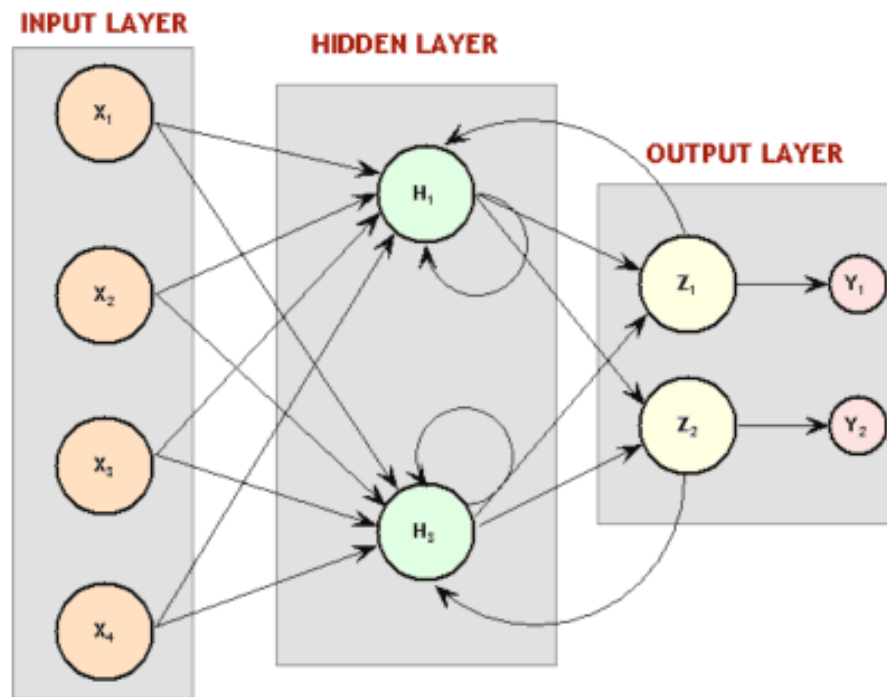
Ref. Lei Zhang and Bing Liu, "Aspect and Entity Extraction for Opinion Mining", 2014

---

- CRF implementation

- CRFsuite <http://www.chokkan.org/software/crfsuite/>
- CRF++ <https://taku910.github.io/crfpp/>
- MALLET <http://mallet.cs.umass.edu/>

- Recurrent Neural Networks (RNN)



- Long Short-Term Memory (LSTM) RNN

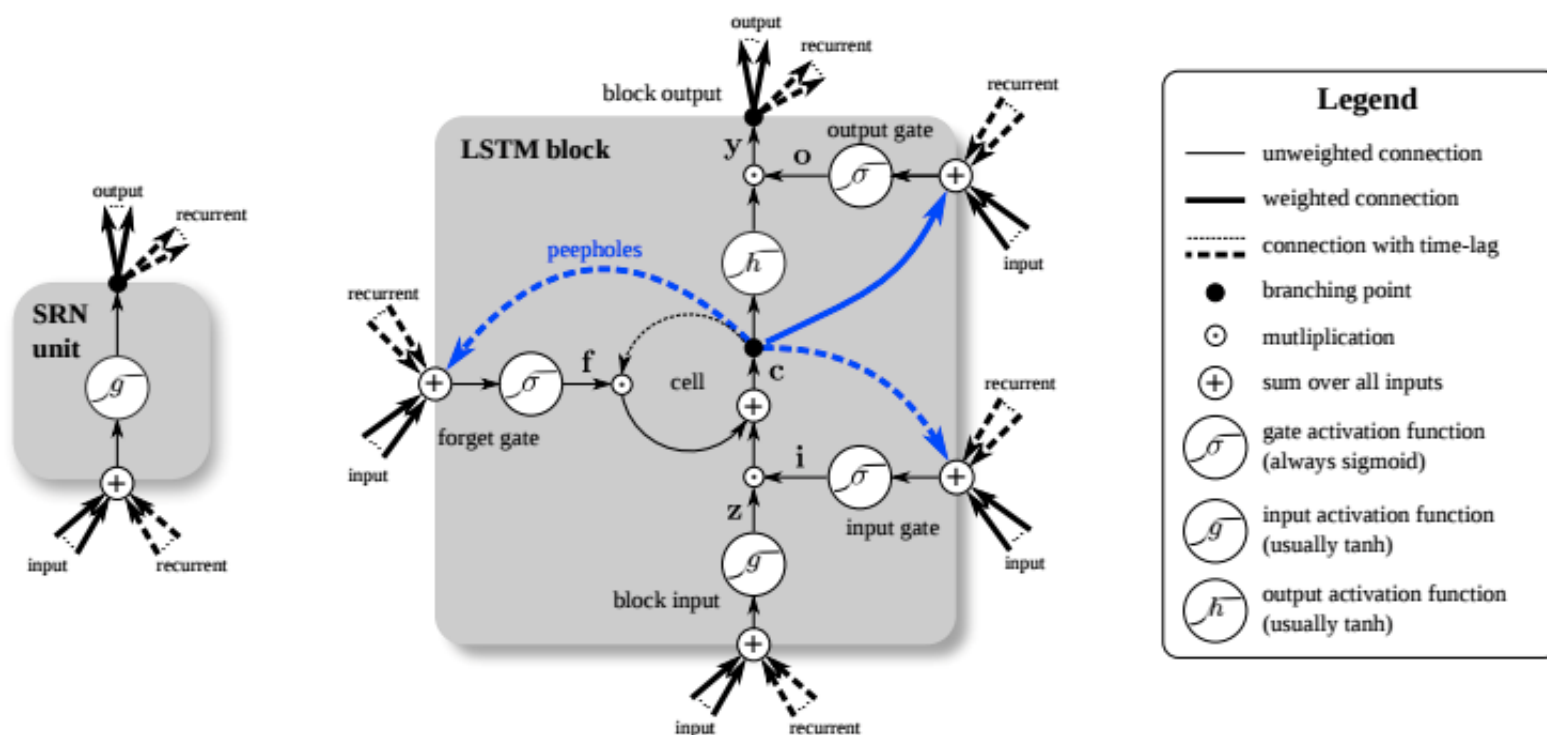


Figure 1. Detailed schematic of the Simple Recurrent Network (SRN) unit (left) and a Long Short-Term Memory block (right) as used in the hidden layers of a recurrent neural network.

# Comparison : ARIMA

---

- ARIMA model (Autoregressive Integrated Moving Average model) is well-known time series analysis approach invented in 1970's.
  - 透過數列的自迴歸進行預測，利用移動平均項數及差分來穩定數列，減少波動；大量應用在金融及經濟學領域
- 思考 1 : ARIMA與本次介紹兩種方法的優劣比較？
- 思考 2 : 若是非同期對應的序列資料，該如何處理？



# Discussions