# Unsupervised Learning: Word Embedding

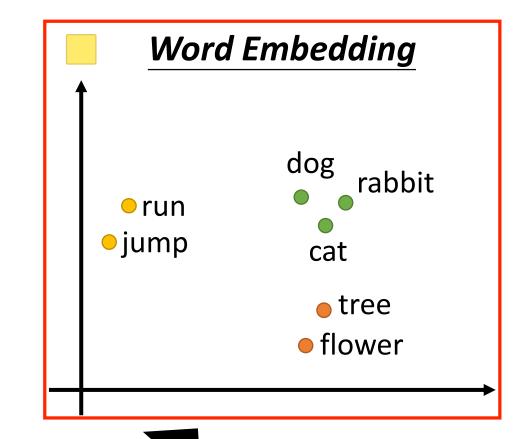
### 1-of-N Encoding

bag = 
$$[0 \ 1 \ 0 \ 0]$$

cat = 
$$[0 \ 0 \ 1 \ 0 \ 0]$$

$$dog = [0 \ 0 \ 0 \ 1 \ 0]$$

elephant =  $[0 \ 0 \ 0 \ 1]$ 



### **Word Class**

class 1

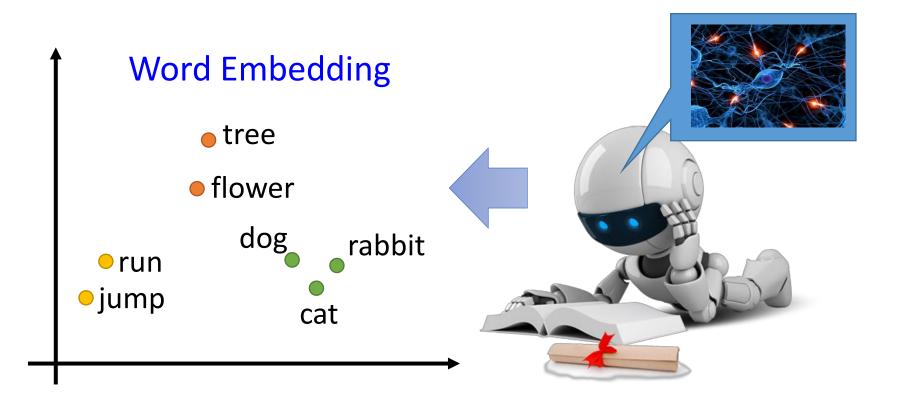
dog cat bird Class 2

ran jumped walk Class 3

flower

tree apple

 Machine learns the meaning of words from reading a lot of documents without supervision



- Machine learns the meaning of words from reading a lot of documents without supervision
- A word can be understood by its context

蔡英文、馬英九 are something very similar

You shall know a word by the company it keeps

馬英九 520宣誓就職

蔡英文 520宣誓就職



## How to exploit the context?

進行 Word Embedding (尋找一個 Word 的 Vector) 有兩種方式:

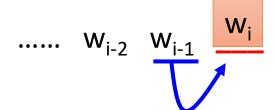
- Count based
  - If two words  $w_i$  and  $w_j$  frequently co-occur,  $V(w_i)$  and  $V(w_i)$  would be close to each other
  - E.g. Glove Vector: http://nlp.stanford.edu/projects/glove/

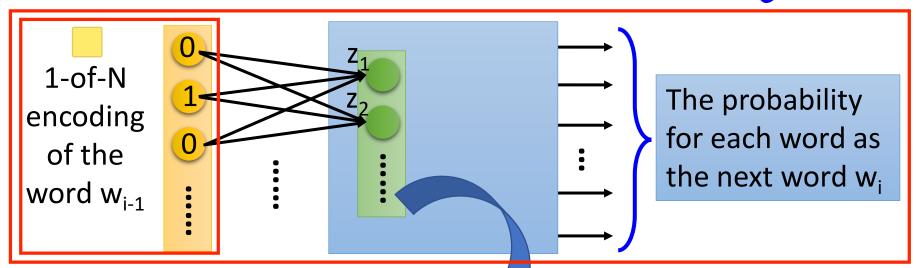
$$V(w_i) \cdot V(w_j)$$

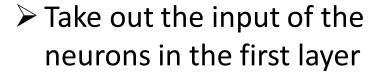
Inner product

Number of times  $w_i$  and  $w_j$  in the same document

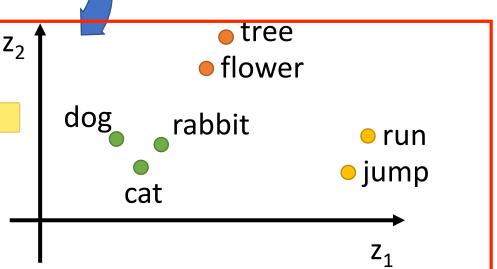
Prediction based



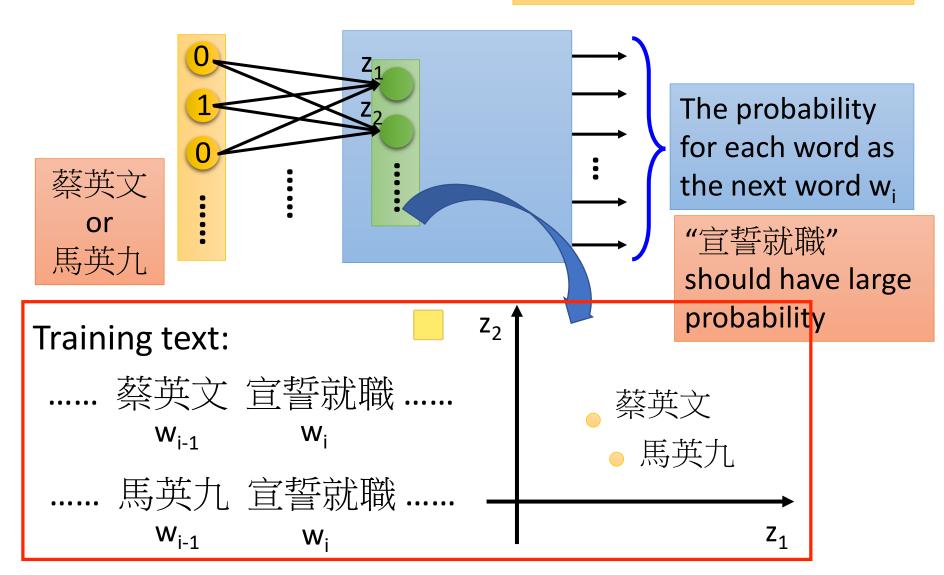




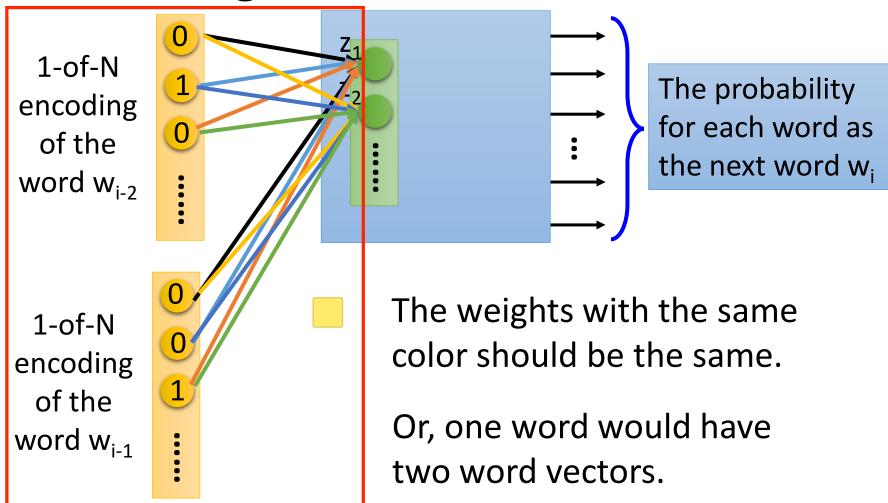
- Use it to represent a word w
- Word vector, word embedding feature: V(w)



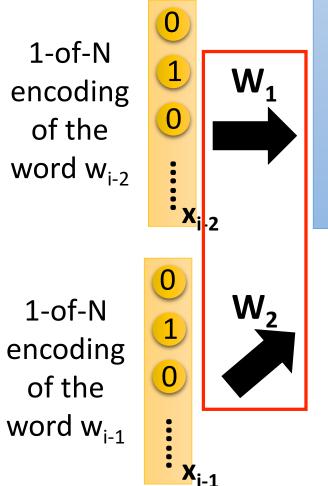
You shall know a word by the company it keeps

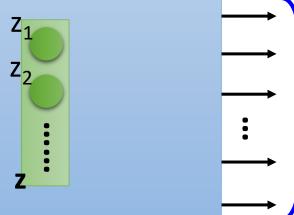


## Sharing Parameters



# Sharing Parameters





The probability for each word as the next word w<sub>i</sub>

The length of  $\mathbf{x_{i-1}}$  and  $\mathbf{x_{i-2}}$  are both |V|. The length of  $\mathbf{z}$  is |Z|.

$$z = W_1 x_{i-2} + W_2 x_{i-1}$$

The weight matrix  $W_1$  and  $W_2$  are both |Z|X|V| matrices.

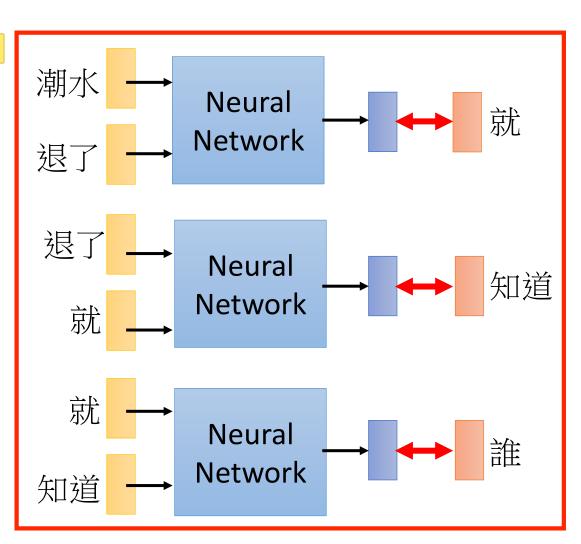
$$W_1 = W_2 = W$$
  $z = W (x_{i-2} + x_{i-1})$ 

# Prediction-based - Training

#### Collect data:

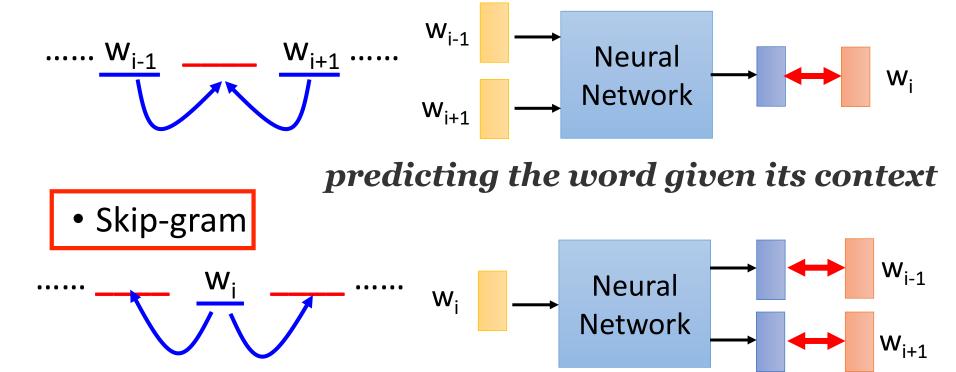
潮水 退了 就 知道 誰 ... 不爽 不要 買 ... 公道價 八萬 一 ...

Minimizing cross entropy

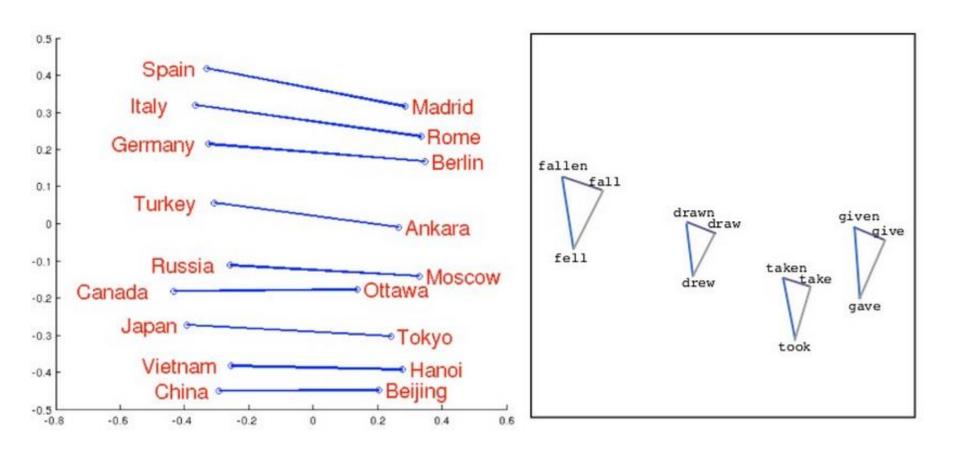


# Prediction-based — Various Architectures

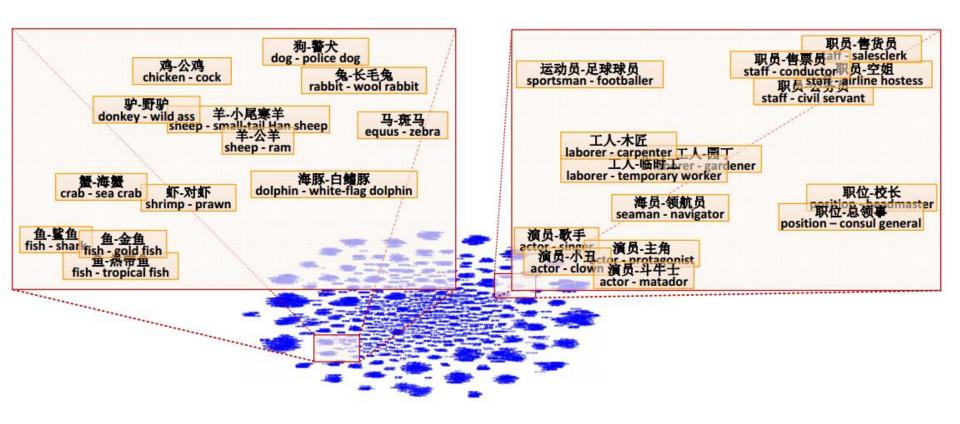
Continuous bag of word (CBOW) model



predicting the context given a word



Source: http://www.slideshare.net/hustwj/cikm-keynotenov2014



Fu, Ruiji, et al. "Learning semantic hierarchies via word embeddings." *Proceedings of the 52th Annual Meeting of the Association for Computational Linguistics: Long Papers*. Vol. 1. 2014.

• Characteristics V(Germany)• V(Berlin) - V(Rome) + V(Italy)  $V(hotter) - V(hot) \approx V(bigger) - V(big)$   $V(Rome) - V(Italy) \approx V(Berlin) - V(Germany)$   $V(king) - V(queen) \approx V(uncle) - V(aunt)$ 

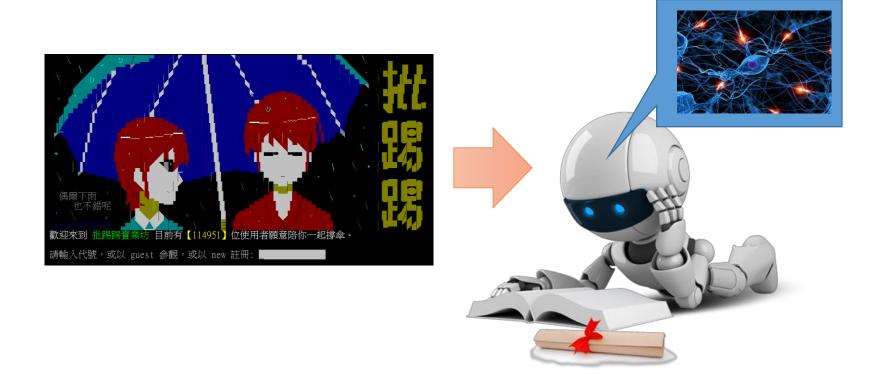
Solving analogies

Rome : Italy = Berlin : ?

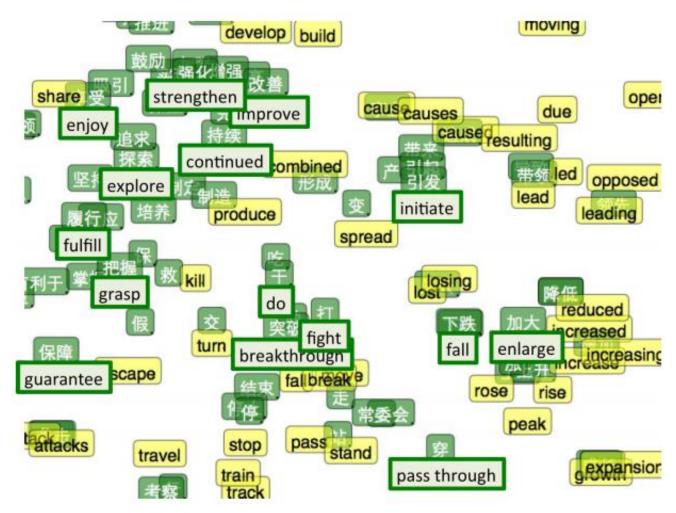
Compute V(Berlin) - V(Rome) + V(Italy)Find the word w with the closest V(w)

## Demo

 Machine learns the meaning of words from reading a lot of documents without supervision



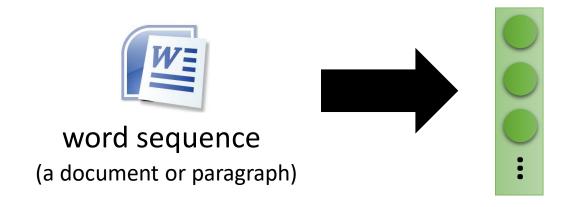
# Multi-lingual Embedding



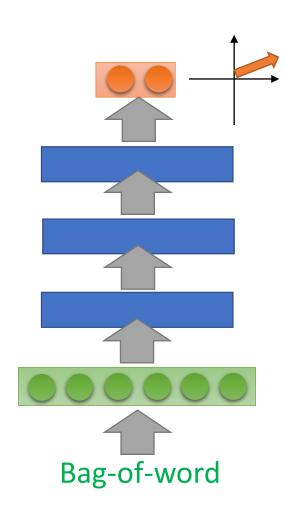
Bilingual Word Embeddings for Phrase-Based Machine Translation, Will Zou, Richard Socher, Daniel Cer and Christopher Manning, EMNLP, 2013

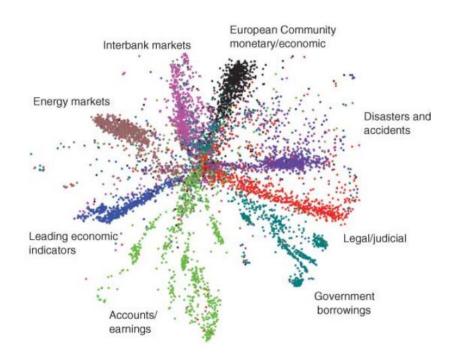
# Document Embedding

- word sequences with different lengths → the vector with the same length
  - The vector representing the meaning of the word sequence
  - A word sequence can be a document or a paragraph



## Semantic Embedding

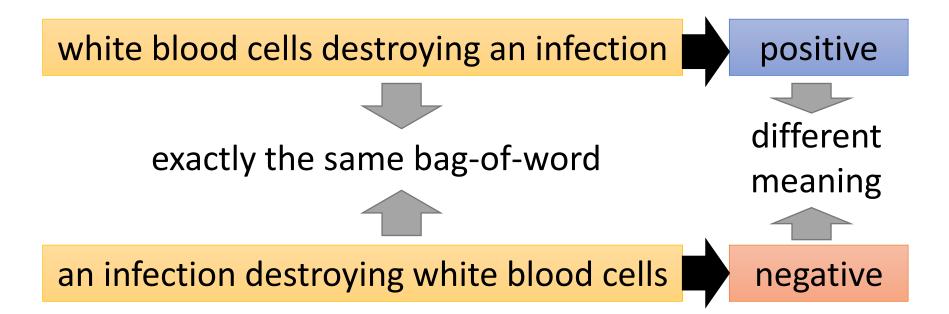




Reference: Hinton, Geoffrey E., and Ruslan R. Salakhutdinov. "Reducing the dimensionality of data with neural networks." *Science* 313.5786 (2006): 504-507

# Beyond Bag of Word

 To understand the meaning of a word sequence, the order of the words can not be ignored.



# Beyond Bag of Word

- Paragraph Vector: Le, Quoc, and Tomas Mikolov.
   "Distributed Representations of Sentences and Documents." ICML, 2014
- **Seq2seq Auto-encoder**: Li, Jiwei, Minh-Thang Luong, and Dan Jurafsky. "A hierarchical neural autoencoder for paragraphs and documents." arXiv preprint, 2015
- **Skip Thought**: Ryan Kiros, Yukun Zhu, Ruslan Salakhutdinov, Richard S. Zemel, Antonio Torralba, Raquel Urtasun, Sanja Fidler, "Skip-Thought Vectors" arXiv preprint, 2015.