

`king` + `woman` - `man` \approx `queen` ?

Word Vector in Natural Language Processing

Zhenghui Wang



**SHANGHAI JIAO TONG
UNIVERSITY**

What is NLP?

Sentiment Analysis

📄 API TEST TOOL

English ▼

Sentiment ▼

Graphical ▼

I ¹ really enjoyed using the ¹ Canon Ixus in Madrid on March 4. The ² Panasonic Lumix ² is a bit disappointing, but the ³ Canon ³ camera is ³ not bad at all. All I want when taking photos is point it and then just press the button. For only 200 dollars, a ⁴ really fair ⁴ price, this ⁵ camera is ⁵ perfect for me. Besides, I have had a ⁶ good ⁶ customer ⁶ service ⁶ experience. ⁷ John Faraday was ⁷ very nice!

ANALYZE TEXT ▶

RESET ↺

Neural Machine Translation



翻译 关闭即时翻译

英语 中文 德语 检测语言 ▼

↔

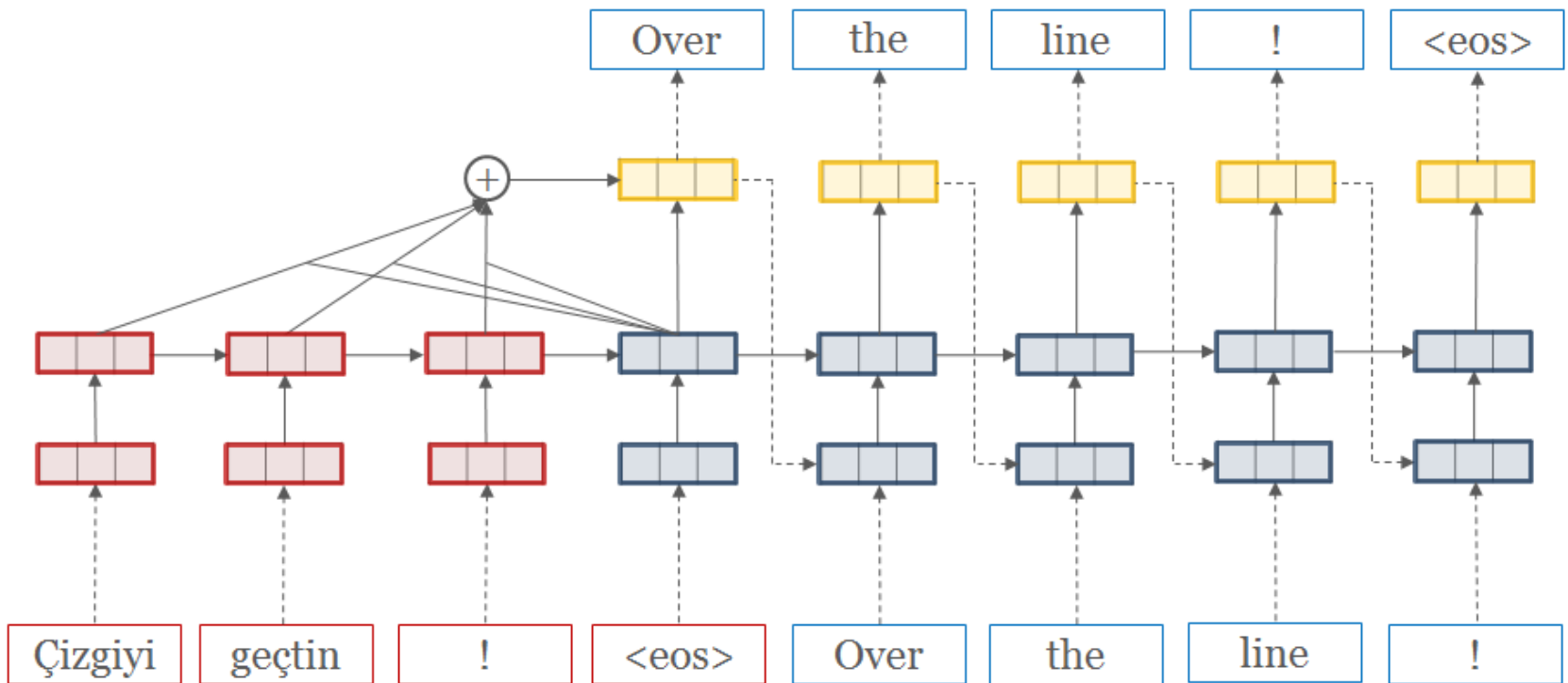
中文(简体) 英语 日语 ▼

翻译

谷歌是一家美国公司 ×
  9/5000

Google is a U.S. company
  

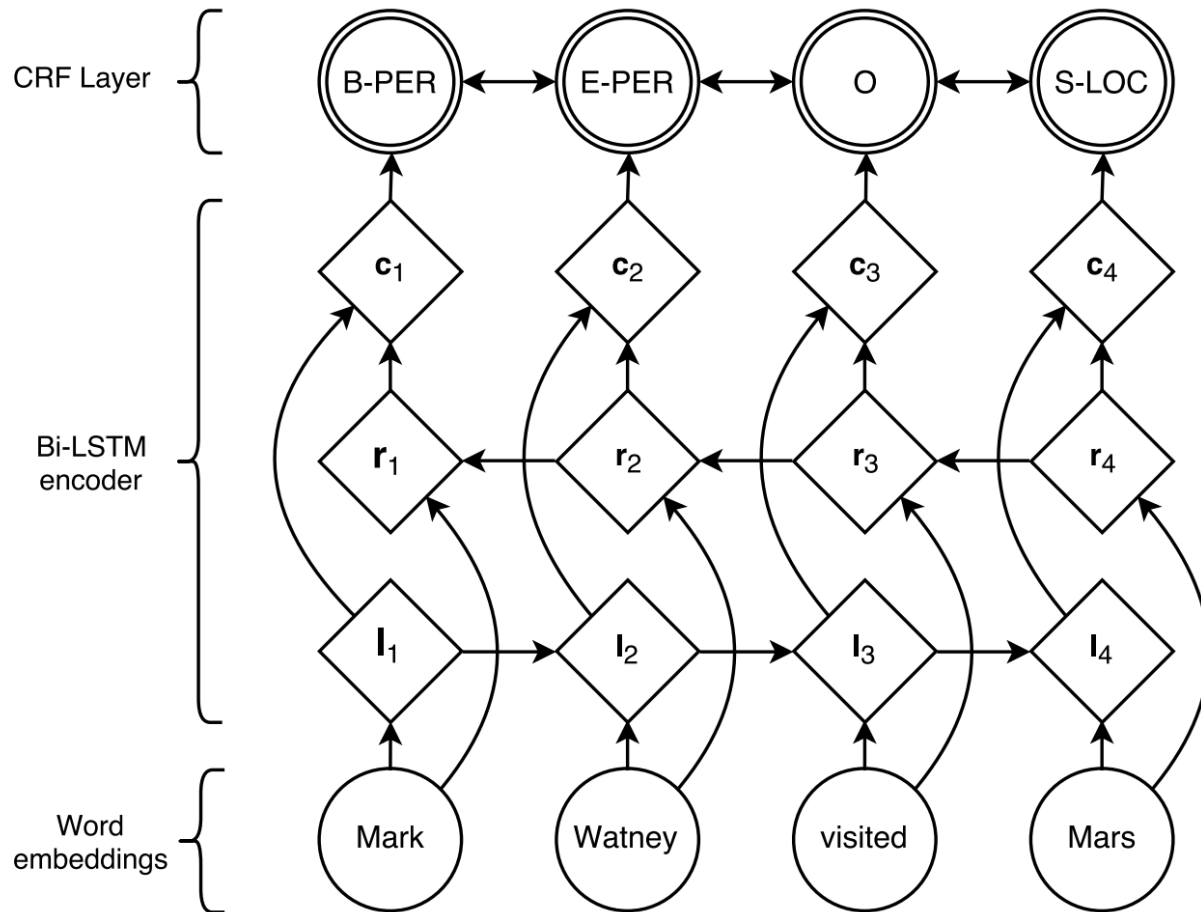
Neural Machine Translation



Named Entity Recognition



Named Entity Recognition



[Lample G, Ballesteros M, Subramanian S, et al. Neural architectures for named entity recognition[J]. arXiv preprint arXiv:1603.01360, 2016.]

What is word vector?

- A way to **represent** the **meaning** of words by **fixed-size vectors**

“One-hot” Representation

motel [0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0]

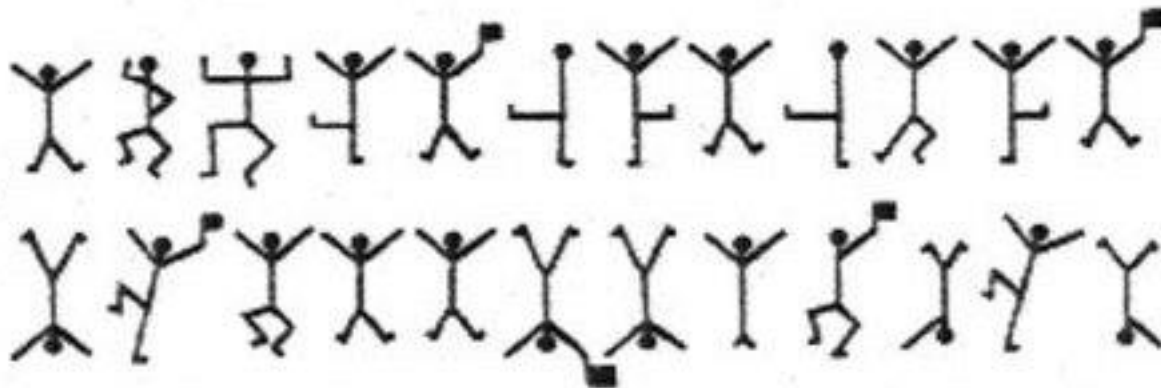
- No natural notion of similarity:

motel [0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0]^T
hotel [0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0] = 0

- Scalability: what if we want to add some new words?

**How does computer understand
natural language?**

Holmes: Dancing Man



Main Idea of word2vec

Predict between every word and its context words!

“You shall know a word by the company it keeps”

(J. R. Firth 1957: 11)

“I like playing X.”

**“The X had an attractive library
of games I wanted.”**

**“Many good PS1 games are
available on either PS3 or X.”**

**“The Sony X was blessed with some
of the biggest names in games.”**

X =



or



?

$$\frac{\begin{array}{c} \rightarrow \\ \text{PSP-3000} \end{array} \cdot \begin{array}{c} \rightarrow \\ \text{PSP-3000} \end{array}}{\begin{array}{c} \rightarrow \\ \text{PSP-3000} \end{array} \cdot \begin{array}{c} \rightarrow \\ \text{PSP-3000} \end{array}} \approx 0.9$$

Word2Vec

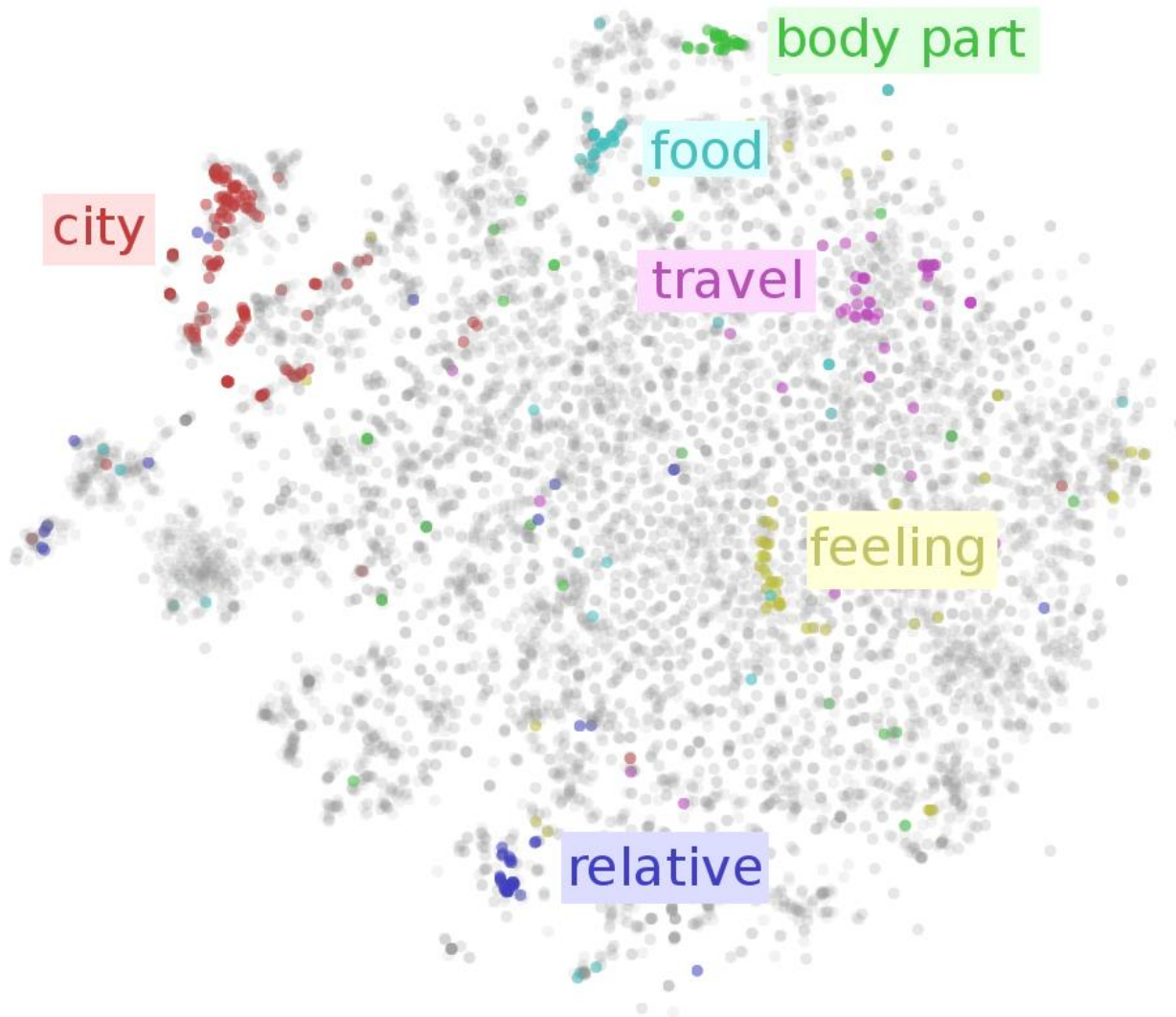
- Distributional similarity based representations

Shanghai = [0.286, 0.792, \dots , -0.107, 0.109]

Beijing = [0.178, 0.490, \dots , -0.287, 0.201]

\dots

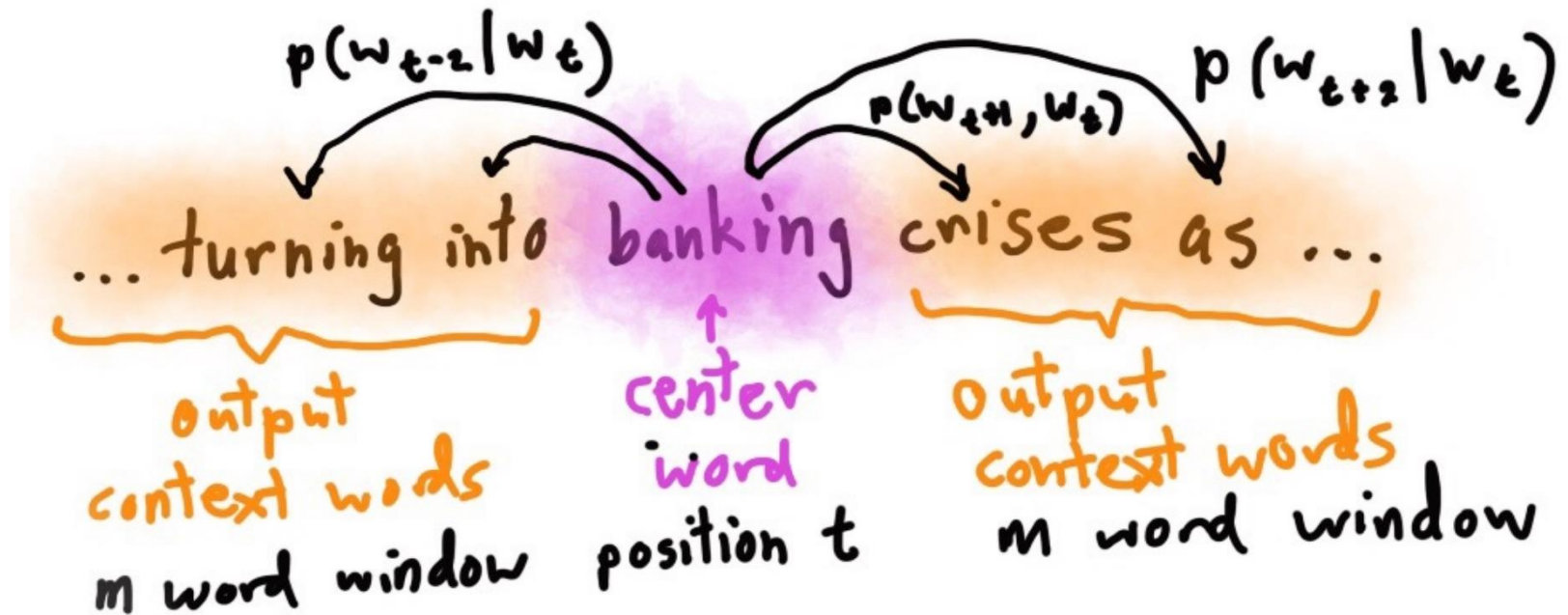
Washington = [0.334, -0.321, \dots , -0.233, 0.391]



How Do We Calculate?

- Two algorithms for word2vec
 - 1. **Skip-grams** (SG):
Predict context words given target (position independent)
 - 2. Continuous Bag of Words (CBOW) :
Predict target word from bag-of-words context

Skip-gram Prediction



$$p(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)}$$

Source Text

Training Samples

The quick brown fox jumps over the lazy dog. →

(the, quick)
(the, brown)

The quick brown fox jumps over the lazy dog. →

(quick, the)
(quick, brown)
(quick, fox)

The quick brown fox jumps over the lazy dog. →

(brown, the)
(brown, quick)
(brown, fox)
(brown, jumps)

The quick brown fox jumps over the lazy dog. →

(fox, quick)
(fox, brown)
(fox, jumps)
(fox, over)

$$p(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)}$$

- For each word $t = 1 \dots T$, predict surrounding words in a window of “radius” m of every word.

- $p(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)}$

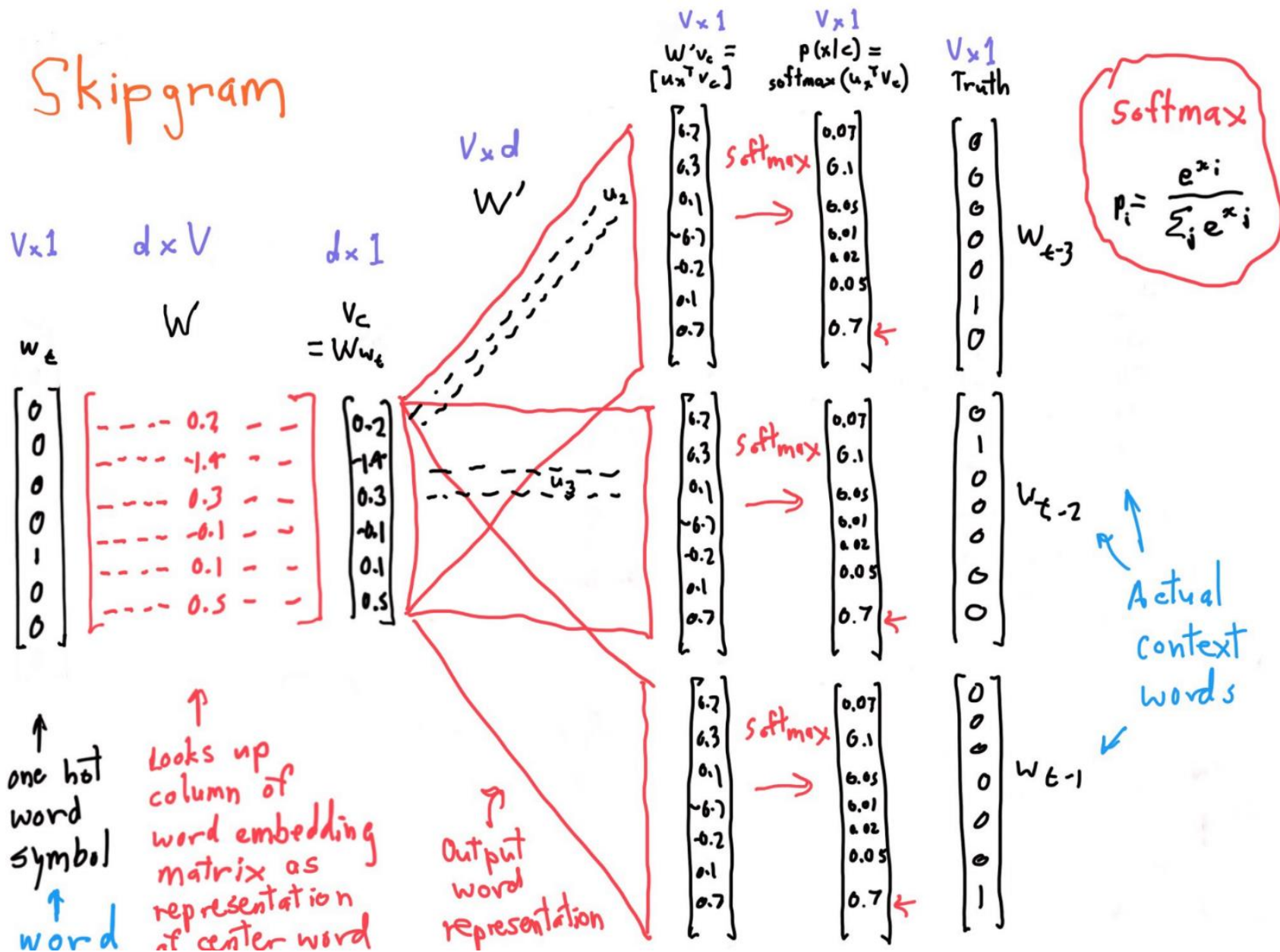
- Objective function: Maximize the probability of any context word given the current center word:

- $J'(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} p(w_{t+j} | w_t)$ Maximize



- $J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log p(w_{t+j} | w_t)$ Minimize

Skipgram



Objective Function

$$\text{Maximize } J'(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} p(w'_{t+j} | w_t; \theta)$$

Or minimize
neg. log
likelihood

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log p(w'_{t+j} | w_t)$$

[negate to minimize;
log is monotone]

↑
text
length

↑
window
size

where

$$p(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)}$$

word IDs ↗

We now take derivatives to work out minimum

Each word type
(vocab entry)
has two word
representations:
as center word
and context word

$$\frac{\partial}{\partial v_c} \log \frac{\exp(u_0^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)}$$

$$= \underbrace{\frac{\partial}{\partial v_c} \log \exp(u_0^T v_c)}_{(1)} - \underbrace{\frac{\partial}{\partial v_c} \log \sum_{w=1}^V \exp(u_w^T v_c)}_{(2)}$$

$$(1) \quad \frac{\partial}{\partial v_c} \underbrace{\log \exp(u_0^T v_c)}_{\text{inverses}} = \frac{\partial}{\partial v_c} u_0^T v_c = u_0$$

Vector!
Not high
school
single
variable
calculus

You can do things one variable at a time,
and this may be helpful when things
get gnarly.

$$\forall j \quad \frac{\partial}{\partial (v_c)_j} u_0^T v_c = \frac{\partial}{\partial (v_c)_j} \sum_{i=1}^d (u_0)_i (v_c)_i \\ = (u_0)_j$$

Each term is zero except when $i=j$

$$\textcircled{2} \frac{\partial}{\partial v_c} \log \underbrace{\sum_{w=1}^v \exp(u_w^T v_c)}_{\substack{f \\ z = g(v_c)}} \\ = \frac{1}{\sum_{w=1}^v \exp(u_w^T v_c)} \cdot \frac{\partial}{\partial v_c} \sum_{x=1}^v \exp(u_x^T v_c)$$

Important to change index

$$\frac{\partial}{\partial v_c} f(\underbrace{g(v_c)}_z) = \frac{\partial f}{\partial z} \cdot \frac{\partial z}{\partial v_c}$$

$$\frac{\partial z}{\partial v_c}$$

Use chain rule

$$= \frac{1}{\sum_{w=1}^v \exp(u_w^T v_c)} \cdot \left(\sum_{x=1}^v \frac{\partial}{\partial v_c} \underbrace{\exp(u_x^T v_c)}_{\substack{f \\ z = g(v_c)}} \right)$$

Move deriv inside sum

$$\left(\sum_{x=1}^v \exp(u_x^T v_c) \frac{\partial}{\partial v_c} u_x^T v_c \right)$$

Chain rule

$$\left(\sum_{x=1}^v \exp(u_x^T v_c) u_x \right)$$

$$\frac{\partial}{\partial v_c} \log(p(o|c)) = u_o - \frac{1}{\sum_{w=1}^V \exp(u_w^T v_c)} \cdot \left(\sum_{x=1}^V \exp(u_x^T v_c) u_x \right)$$

$$= u_o - \sum_{x=1}^V \frac{\exp(u_x^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)} u_x$$

Distribute
term
across sum

$$= u_o - \underbrace{\sum_{x=1}^V p(x|c) u_x}_{\text{this an expectation: average over all context vectors weighted by their probability}}$$

= observed - expected

This is just the derivatives for the center vector parameters
Also need derivatives for output vector parameters
(they're similar)

Then we have derivative w.r.t. all parameters and can minimize

Finding the degree of similarity between two words

- Vectors trained on ~100M sentences
- word segmentation with **jieba**
 - github.com/fxsjy/jieba
- $sim(u, v) = \frac{u^T v}{|u| \cdot |v|}$
- '不洗头发',
 - ('不洗头', 0.761),
 - ('洗次头', 0.751),
 - ('洗会油', 0.729),
 - ('超油', 0.715),
 - ('就会油', 0.704)

Finding the degree of similarity between two words

- '送病人',
 - ('急救车', 0.673),
 - ('抬走', 0.656),
 - ('妇产科住院', 0.648),
 - ('转送', 0.638),
 - ('救护车', 0.637)
- '胃酸过多',
 - ('胃酸', 0.811),
 - ('胃炎', 0.786),
 - ('反流性胃炎', 0.767),
 - ('烧心', 0.757),
 - ('胃溃疡', 0.756)
- '经济问题',
 - ('断绝关系', 0.645),
 - ('愚孝', 0.593),
 - ('不同意', 0.593),
 - ('撕破脸', 0.591),
 - ('两头跑', 0.586)
- '心不在焉',
 - ('走神', 0.727),
 - ('认真听讲', 0.725),
 - ('一片空白', 0.716),
 - ('专心', 0.713),
 - ('东张西望', 0.712)

Finding the degree of similarity between two words

Nearest words to
frog:

1. frogs

2. toad

3. litoria

4. Leptodactylidae

5. rana

6. lizard

7. eleutherodactylus



litoria



leptodactylidae

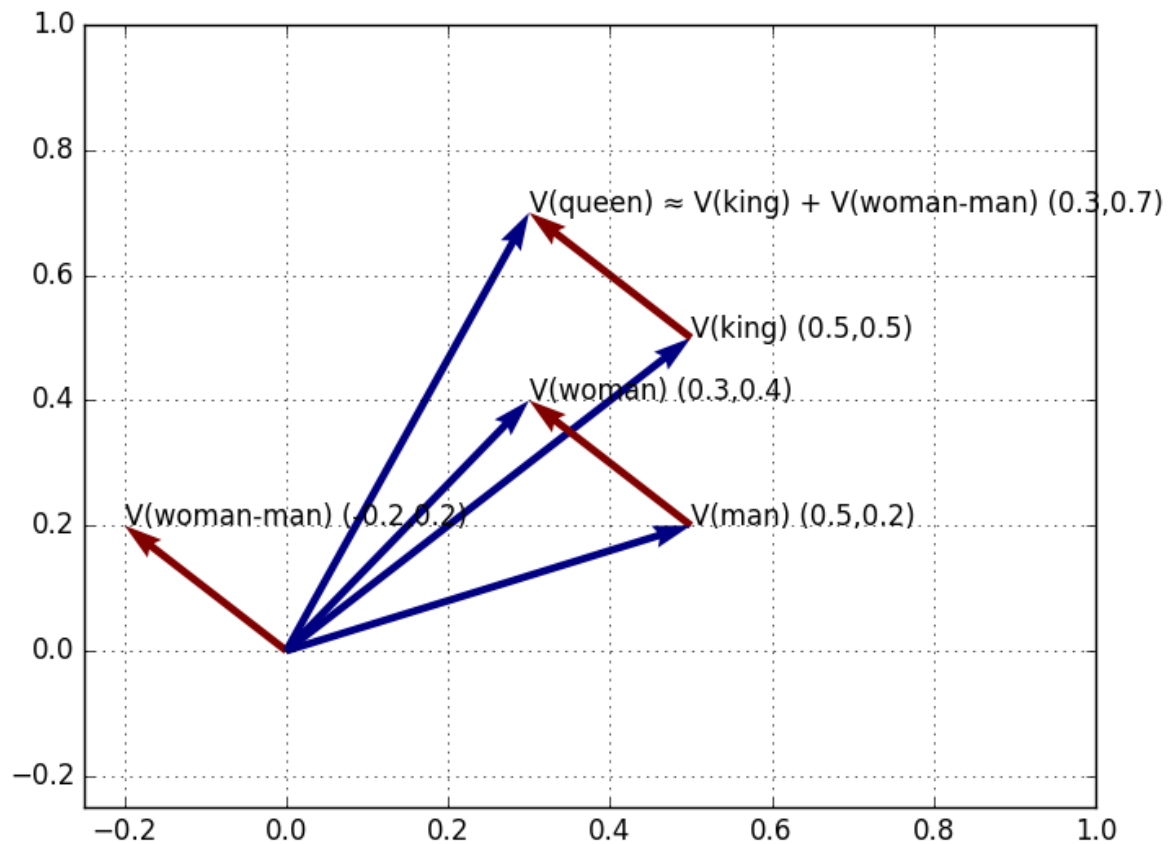


rana



eleutherodactylus

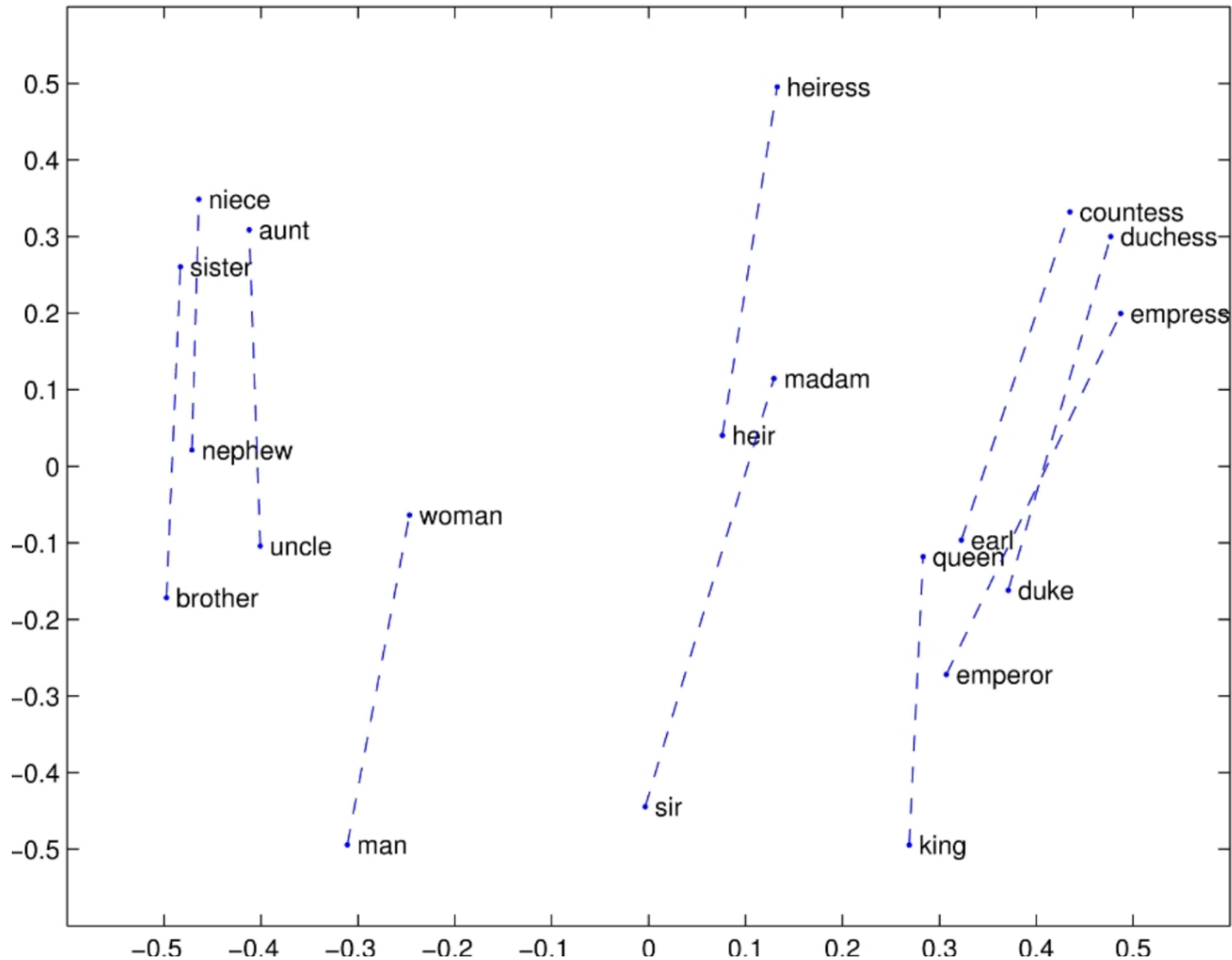
`king` - `man` + `woman` \approx `queen`



`king` - `man` + `woman` \approx `queen`

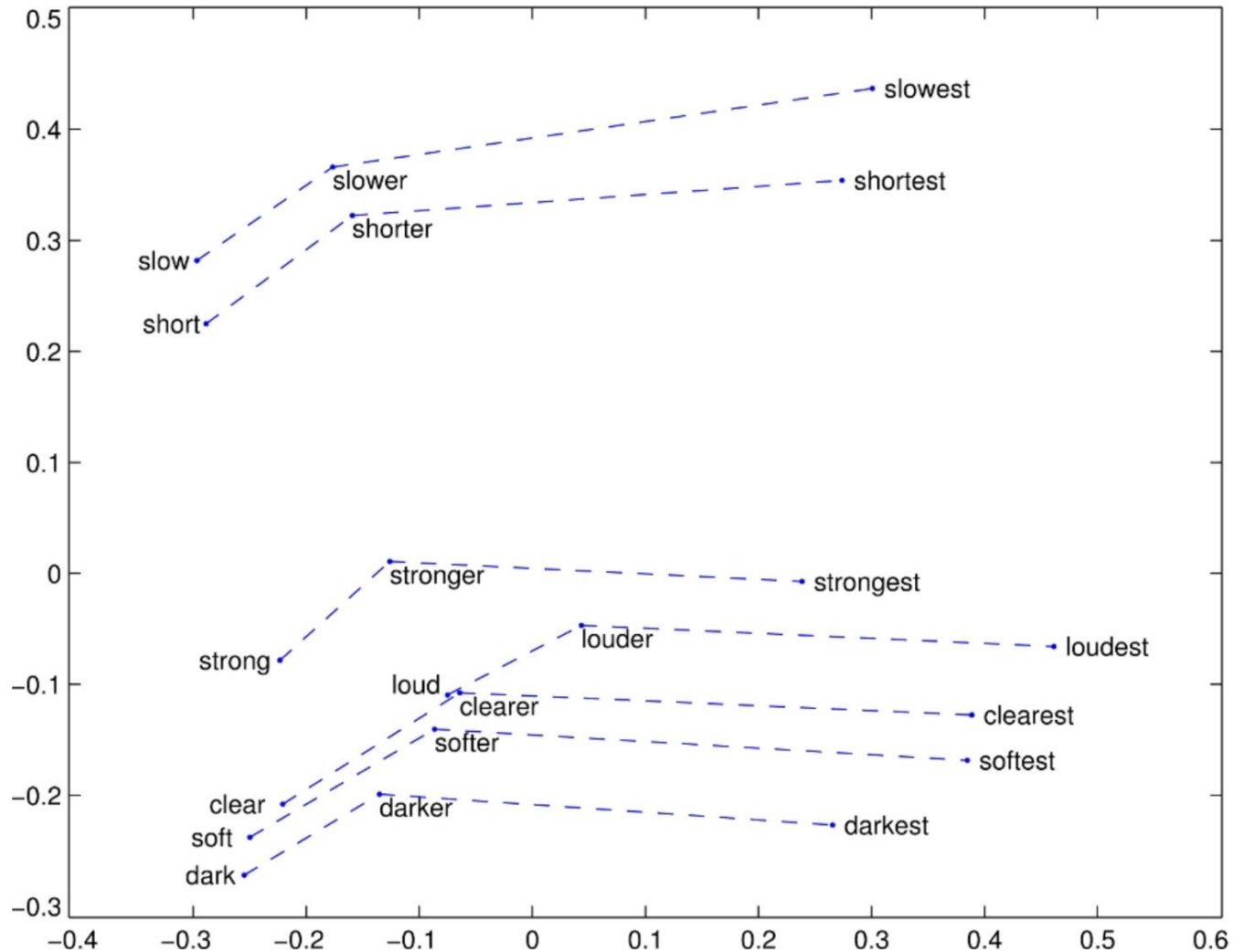
<i>Expression</i>	<i>Nearest token</i>
Paris - France + Italy	Rome
bigger - big + cold	colder
sushi - Japan + Germany	bratwurst
Cu - copper + gold	Au
Windows - Microsoft + Google	Android
Montreal Canadiens - Montreal + Toronto	Toronto Maple Leafs

Linear Substructures: man-woman



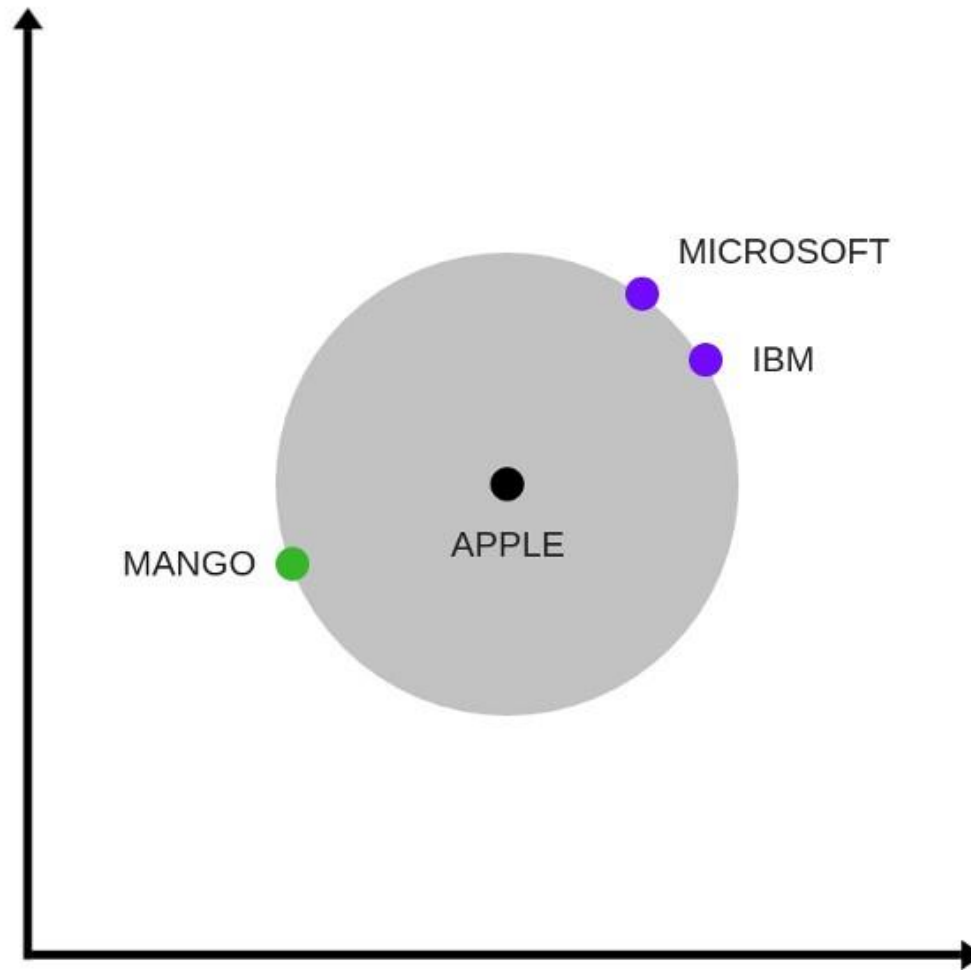
[Pennington J, Socher R, Manning C. Glove: Global vectors for word representation[C]//Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP). 2014: 1532-1543.]

Linear Substructures: comparative-superlative



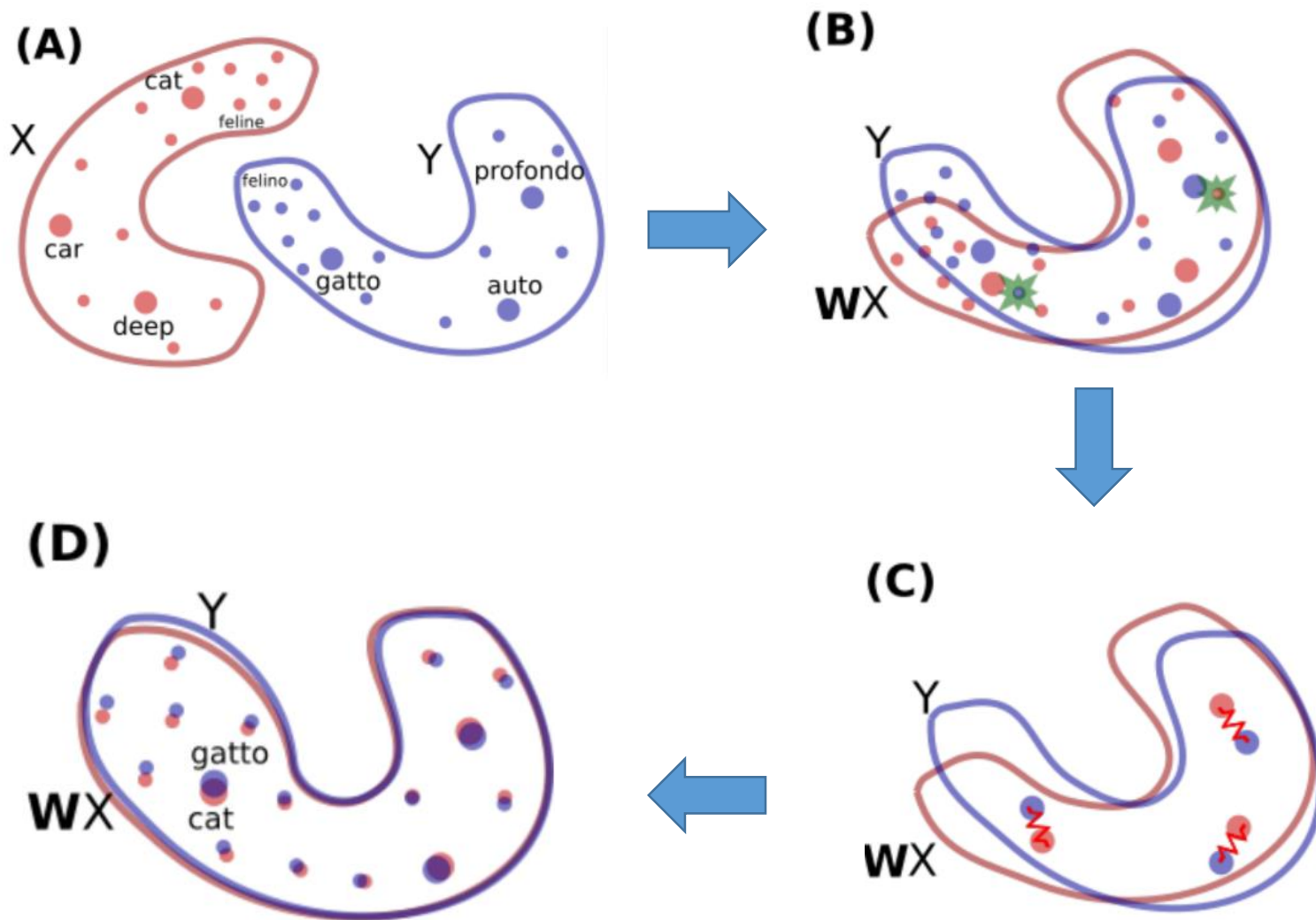
[Pennington J, Socher R, Manning C. Glove: Global vectors for word representation[C]//Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP). 2014: 1532-1543.]

Multiple Contexts of Apple



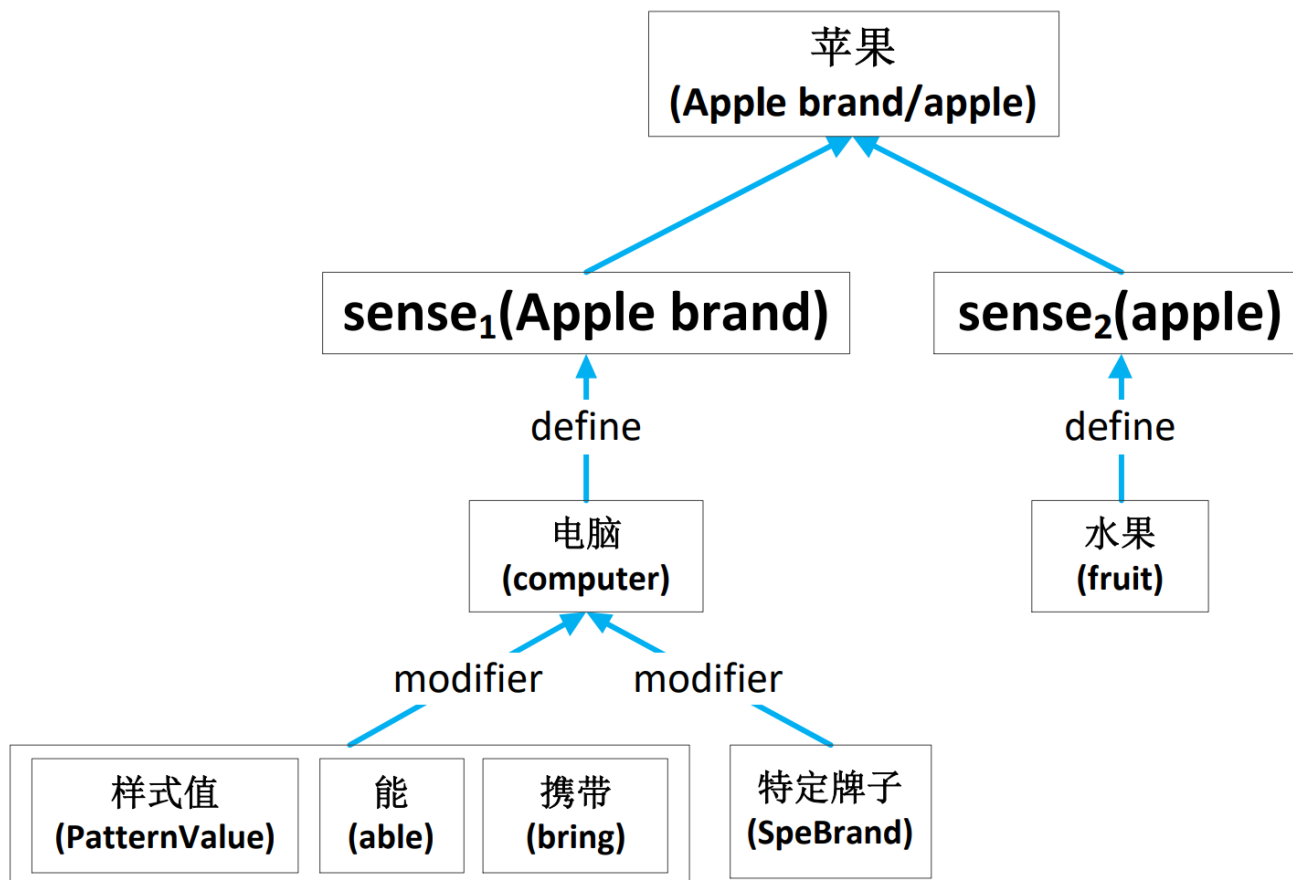
Research Frontier of Word Vector

Word translation without parallel data

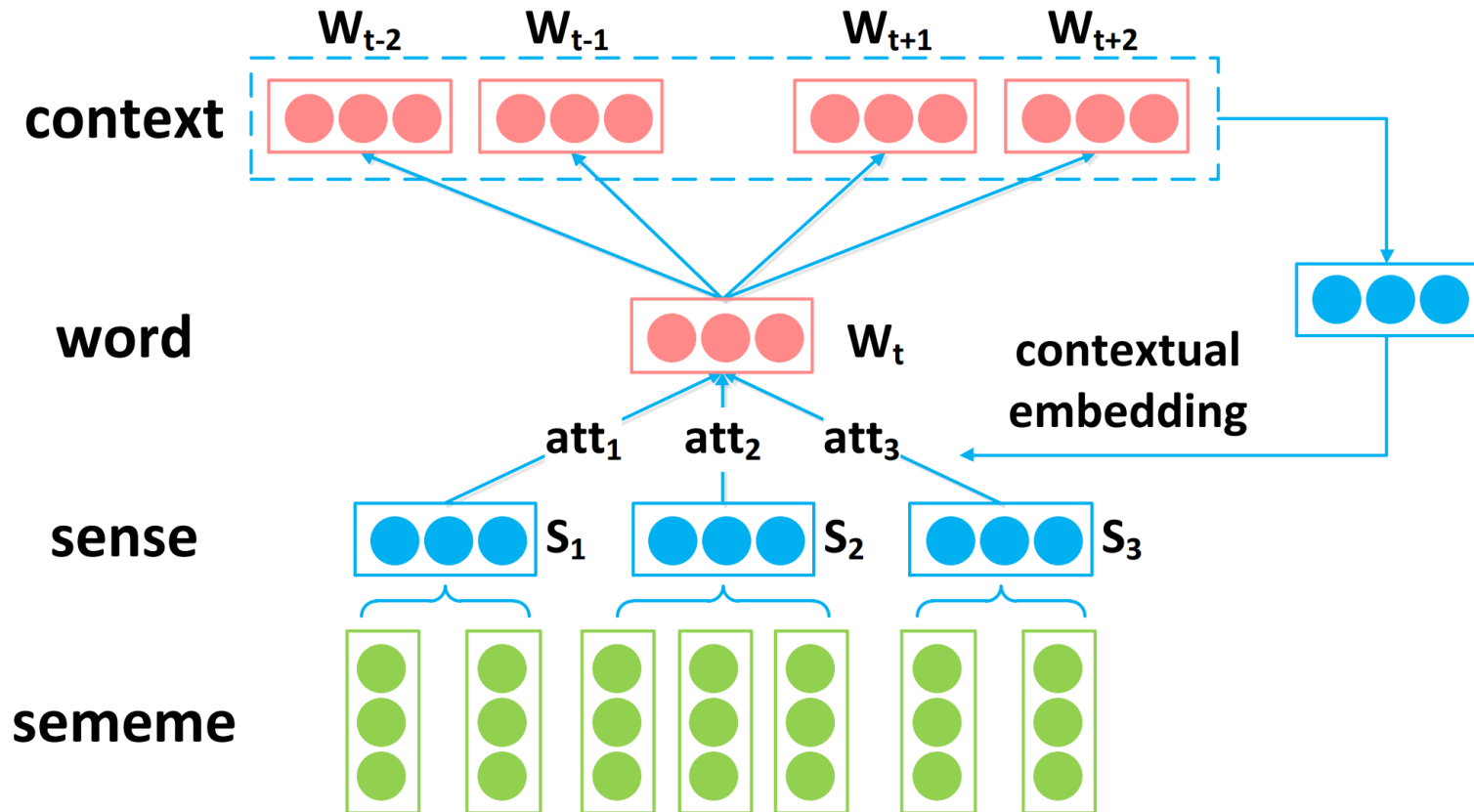


[Anonymous (2018). Word translation without parallel data. *International Conference on Learning Representations*, , .]

Improved Word Representation Learning with Sememes



Improved Word Representation Learning with Sememes



Q&A

Thanks!