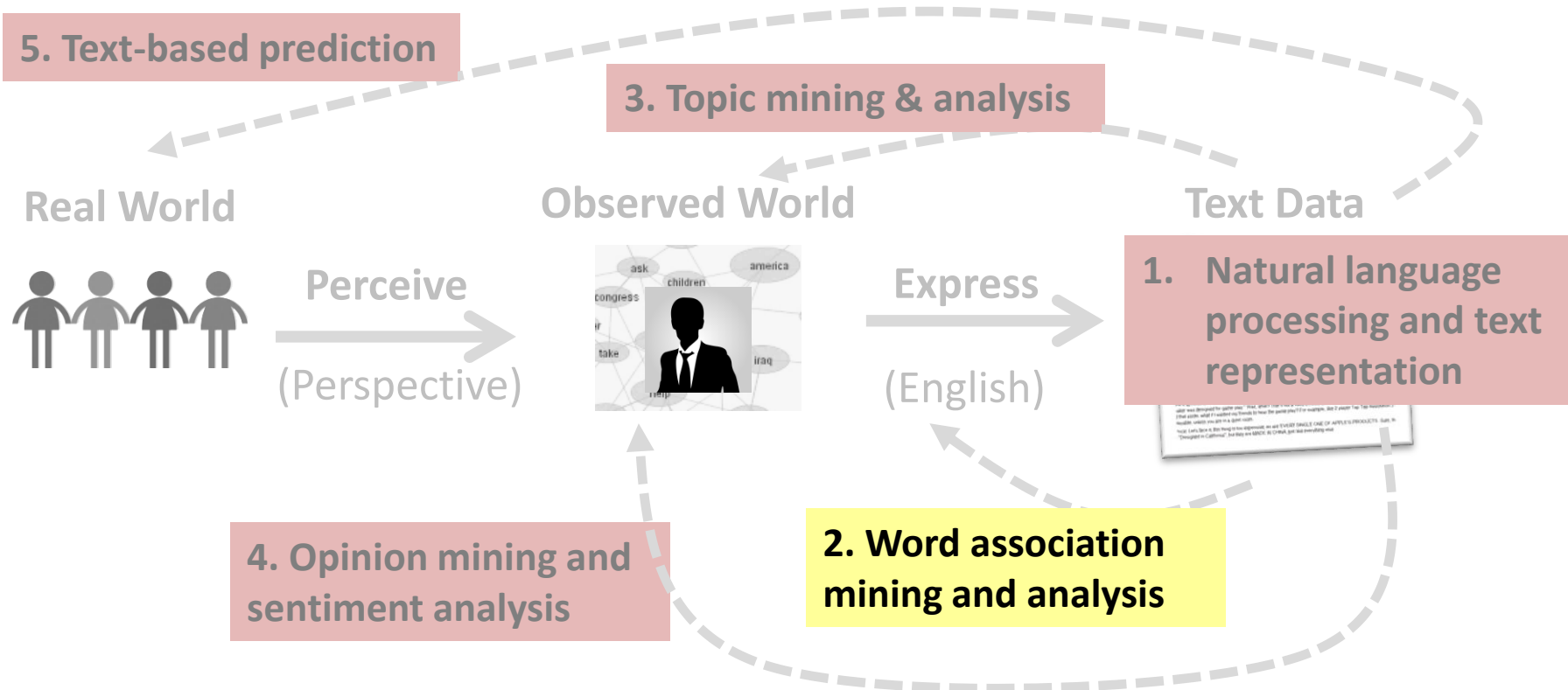# Paradigmatic Relation Discovery

Parts 1-3
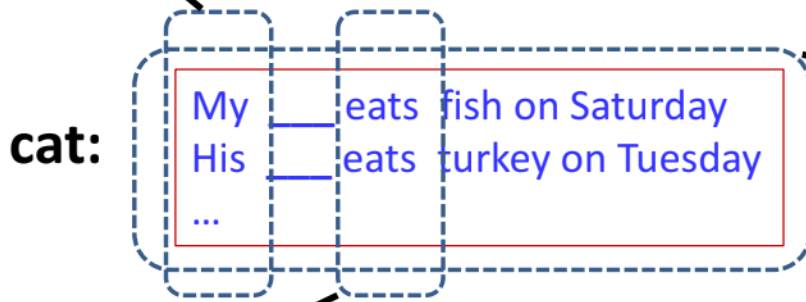
ChengXiang "Cheng" Zhai
Department of Computer Science
University of Illinois at Urbana-Champaign

# Paradigmatic Relation Discovery



5. Text-based prediction

3. Topic mining & analysis

Real World    Observed World    Text Data

Perceive    Express

(Perspective)    (English)

1. Natural language processing and text representation

4. Opinion mining and sentiment analysis

2. Word association mining and analysis

2

# Word Context as "Pseudo Document"

① **Left1("cat") = {"my", "his", "big", "a", "the",...}**

影响 的 words 在 cat 周围.

**cat:**

| My | ___ eats | fish on Saturday |
| His | ___ eats | turkey on Tuesday |
| ... | | |

**Window8("cat") =**
**{"my", "his", "big",**
**"eats", "fish", ...}**

② **Right1("cat") = {"eats", "ate", "is", "has", ....}**

> **Context = pseudo document = "bag of words"**
> **Context may contain adjacent or non-adjacent words**

# Measuring Context Similarity

计算 2个词的相似度.

**Sim("Cat", "Dog") =**

Combine
all the
perspective

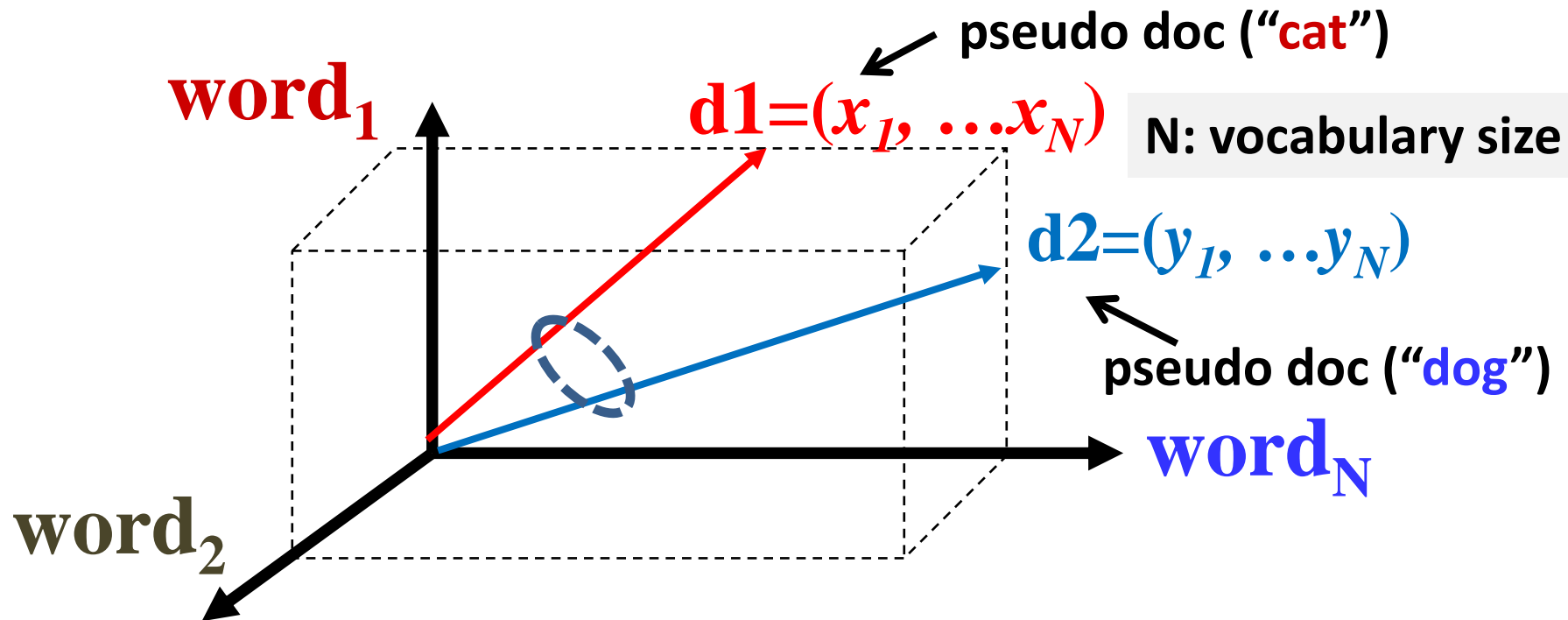$\qquad$ **Sim(Left1("cat"), Left1("dog"))**

$\qquad$ **+ Sim(Right1("cat"), Right1("dog")) +**

$\qquad$ **...**

$\qquad\qquad$ **+ Sim(Window8("cat"), Window8("dog"))=?**

**High** sim(word1, word2)     2个词相似度很高 班法 → 两者有深镁系

$\qquad$ ➜ word1 and word2 are **paradigmatically related**

# Bag of Words ➡ Vector Space Model (VSM)

pseudo doc ("cat")

**word₁** → $word_1$

$d1 = (x_1, \ldots x_N)$

**N: vocabulary size**

$d2 = (y_1, \ldots y_N)$

pseudo doc ("dog")

**word₂** → $word_2$

**word_N** → $word_N$

| Terms: | | "eats" | "ate" | "is" | "has" | .... | |
|--------|--|--------|-------|------|-------|------|--|
| Vector: | | ( 5, | 3, | 10, | 3 | .... | ) |

# VSM for Paradigmatic Relation Mining

**1. How to compute each vector?**

$$\mathbf{word_1}$$

$d1 = (x_1, \ldots x_N)$

$x_i = ?$

$d2 = (y_1, \ldots y_N)$

$y_j = ?$

**2. Sim($d1$, $d2$)=?**

$$\mathbf{word_N}$$

$$\mathbf{word_2}$$

**Many approaches are possible
(most developed originally for text retrieval).**

# Expected Overlap of Words in Context (EOWC)

**Probability that a randomly picked word from d1 is wi**

每个词的权重是？随机从d1从选取一代提书前词的标院率。也就是词频。

Count of word wi in d1

$$d1=(x_1, \ldots x_N) \quad x_i = c(w_i, d1)/|d1|$$

$$d2=(y_1, \ldots y_N) \quad y_i = c(w_i, d2)/|d2|$$

**Total counts of words in d1**

$$Sim(d1,d2)=d1.d2 = x_1 y_1 + \ldots + x_N y_N = \sum_{i=1}^{N} x_i y_i$$

Probability that two randomly picked words from d1 and d2, respectively, are identical. if they are 聚合关系

are not identical. if they are not 聚合关系.

# Would EOWC Work Well?

EOWC 核心思想

- Intuitively, it makes sense: The more overlap the two context documents have, the higher the similarity would be.

- However:
  - It favors matching one frequent term very well over matching more distinct terms. 对于出现次数少的词不友好.
  - It treats every word equally (overlap on "the" isn't as so meaningful as overlap on "eats").

这种词就算有 overlap 也对我们分析文本也没用

重点

# Expected Overlap of Words in Context (EOWC)

**Probability that a randomly picked word from d1 is wi**

**Count of word wi in d1**

$$d1 = (x_1, \ldots x_N) \qquad x_i = \mathbf{c}(w_i, \mathbf{d1})/|\mathbf{d1}|$$

$$\mathbf{d2} = (y_1, \ldots y_N) \qquad y_i = \mathbf{c}(w_i, \mathbf{d2})/|\mathbf{d2}|$$

**Total counts of words in d1**

$$Sim(d1, d2) = d1 \cdot d2 = x_1 y_1 + \ldots + x_N y_N = \sum_{i=1}^{N} x_i y_i$$

Probability that two randomly picked words from d1 and d2, respectively, are identical.

# Improving EOWC with Retrieval Heuristics

- It favors matching one frequent term very well over matching more distinct terms.

  ➜ **Sublinear transformation of Term Frequency (TF)**

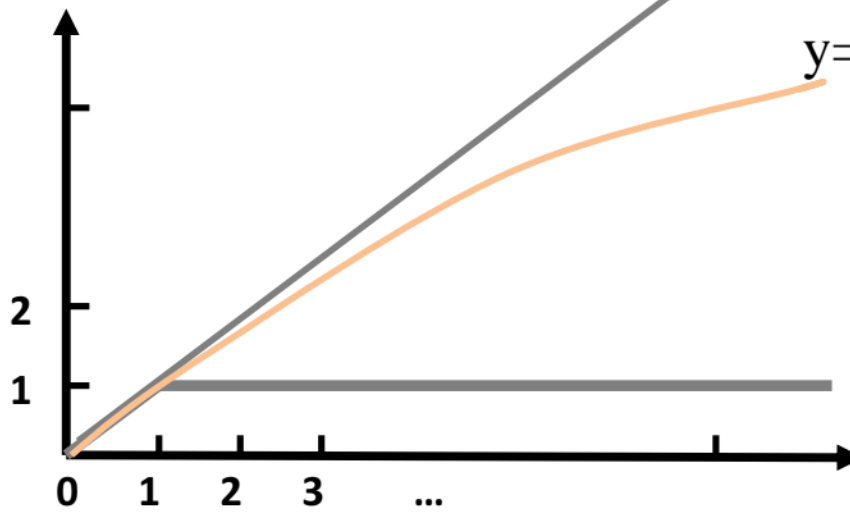- It treats every word equally (overlap on "the" isn't as so meaningful as overlap on "eats").

  ➜ **Reward matching a rare word:  IDF term weighting**

  *put more weight on rare word*

# TF Transformation: c(w,d)➜TF(w,d)

TF：出现次数，会比现
三卷的问题.

**Term Frequency Weight**

$$y=TF(w,d)$$

$y= x$

$y= \log(1+x)$

出现次数多的
词会在TF上有
惩罚.

出现时，则为1.没出现时，则为
0.

0/1 bit vector
(ignore counts)

$$x=c(w,d)$$

0    1    2    3    ...

2

1

4

# TF Transformation: BM25 Transformation

Term Frequency Weight

$y=TF(w,d)$

*Very large k*

$k+1$

$y=\dfrac{(k+1)x}{x+k}$

临界值为k+1.

2

1

$k=0$

$x=c(w,d)$

0  1  2  3  ...

# IDF Weighting: Penalizing Popular Terms



IDF(W)

log(M+1)

total number of docs in collection

$$IDF(W) = \log[(M+1)/k]$$

reward rare word

penal frequent word

total number of docs containing W
(Doc Frequency)

1

M

k (doc freq)

# Adapting BM25 Retrieval Model for Paradigmatic Relation Mining

$d1=(x_1, \ldots x_N)$

$$BM25(w_i, d1) = \frac{(k+1)c(w_i, d1)}{c(w_i, d1) + k(1 - b + b * |d1| / avdl)}$$

*该 doc 长度*

*平均长度 for all docs*

*增加时变为 no normalization*

$$x_i = \frac{BM25(w_i, d1)}{\sum_{j=1}^{N} BM25(w_j, d1)}$$

*normalized,
为了 $\sum_{i=1}^{N} x_i = 1$*

$b \in [0,1]$

$k \in [0, +\infty)$

$d2=(y_1, \ldots y_N)$

$y_i$ is defined similarly

$$Sim(d1, d2) = \sum_{i=1}^{N} IDF(w_i) x_i \, y_i$$

7

# BM25 can also Discover Syntagmatic Relations

$$d1=(x_1, \ldots x_N)$$

$$BM25(w_i, d1) = \frac{(k+1)c(w_i, d1)}{c(w_i, d1) + k(1 - b + b*|d1|/avdl)}$$

$$x_i = \frac{BM25(w_i, d1)}{\sum_{j=1}^{N} BM25(w_j, d1)}$$

$$b \in [0,1]$$

$$k \in [0, +\infty)$$

IDF-weighted $d1 = (x_1*IDF(w_1), \ldots, x_N*IDF(w_N))$

使用了IDF 从而 综合了 词义以及词的稀有度.

The highly weighted terms in the context vector of word w are likely syntagmatically related to w.

# Summary

- Main idea for discovering paradigmatic relations:
  - Collecting the context of a candidate word to form a pseudo document (bag of words)
  - Computing similarity of the corresponding context documents of two candidate words
  - Highly similar word pairs can be assumed to have paradigmatic relations
- Many different ways to implement this general idea
- Text retrieval models can be easily adapted for computing similarity of two context documents
  - BM25 + IDF weighting represents the state of the art
  - Syntagmatic relations can also be discovered as a "by product"

9