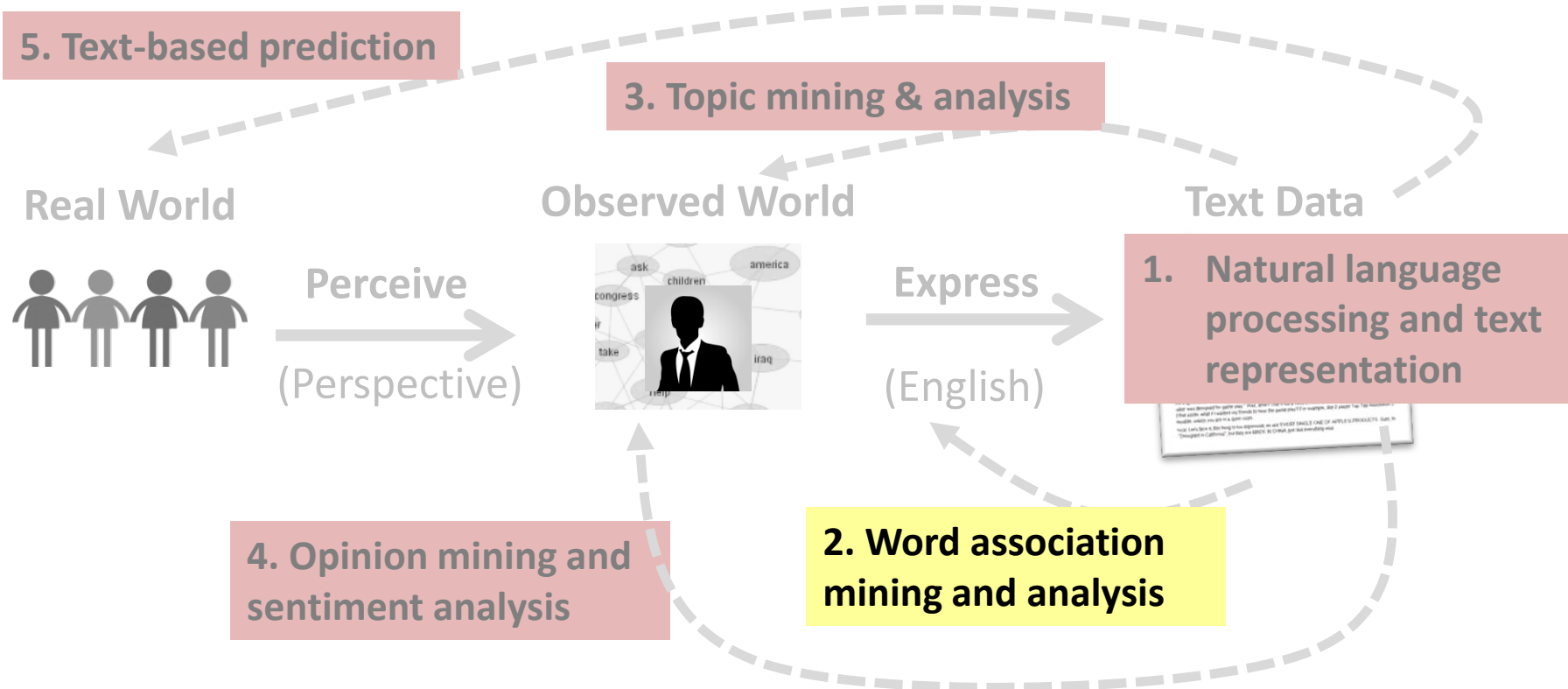# Paradigmatic Relation Discovery

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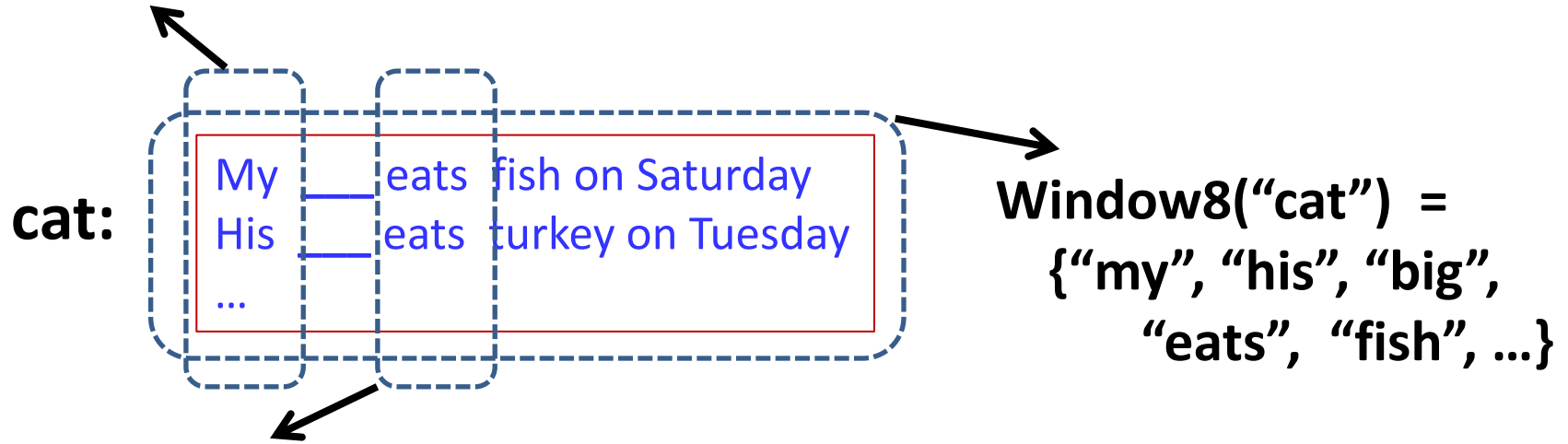
ChengXiang "Cheng" Zhai
Department of Computer Science
University of Illinois at Urbana-Champaign

# Paradigmatic Relation Discovery



**5. Text-based prediction**

**3. Topic mining & analysis**

Real World

Observed World

Text Data

**Perceive**

**Express**

1. **Natural language processing and text representation**

(Perspective)

(English)

**4. Opinion mining and sentiment analysis**

**2. Word association mining and analysis**

# Word Context as "Pseudo Document"

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

**cat:**

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

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

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

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

3

# Measuring Context Similarity

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

**Sim(Left1("cat"), Left1("dog"))**

**+ Sim(Right1("cat"), Right1("dog")) +**

**…**

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

**High** sim(word1, word2)
➔ word1 and word2 are **paradigmatically related**

# Bag of Words ➜ Vector Space Model (VSM)



pseudo doc ("cat")

$d1=(x_1, \ldots x_N)$

**N: vocabulary size**

$d2=(y_1, \ldots y_N)$

pseudo doc ("dog")

word₁

word₂

word_N

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

# VSM for Paradigmatic Relation Mining

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

$\mathbf{word_1}$

$\mathbf{d1}=(x_1, \ldots x_N)$

$x_i = ?$

$\mathbf{d2}=(y_1, \ldots y_N)$

$y_j = ?$

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

$\mathbf{word_N}$

$\mathbf{word_2}$

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

6

# Expected Overlap of Words in Context (EOWC)

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

**Count of word wi in d1**

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

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

**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.

# Would EOWC Work Well?

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

# Expected Overlap of Words in Context (EOWC)

**Probability that a randomly picked word from d1 is wi** → **Count of word wi in d1** →

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

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

**Total counts of words in d1**

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

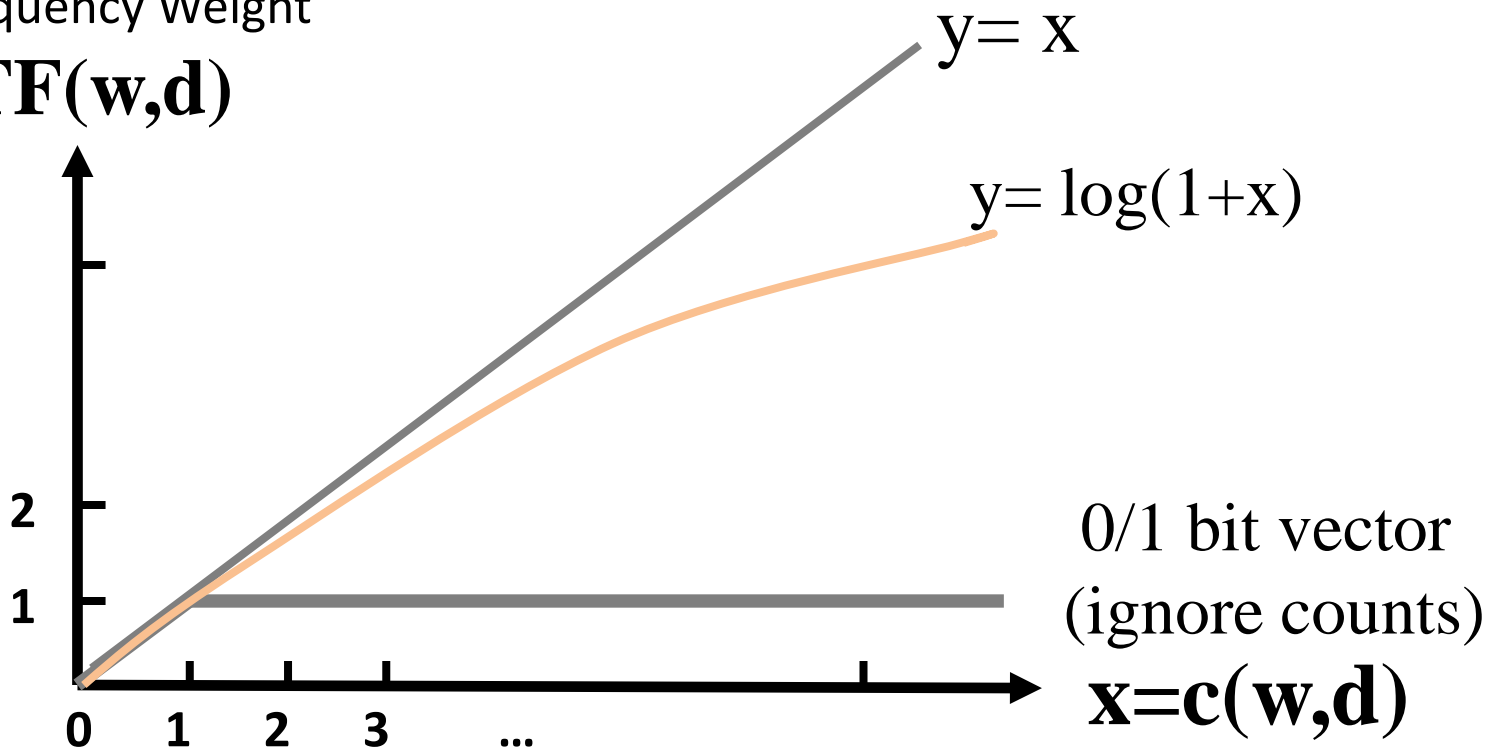Probability that two randomly picked words from d1 and d2, respectively, are identical.

# 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**
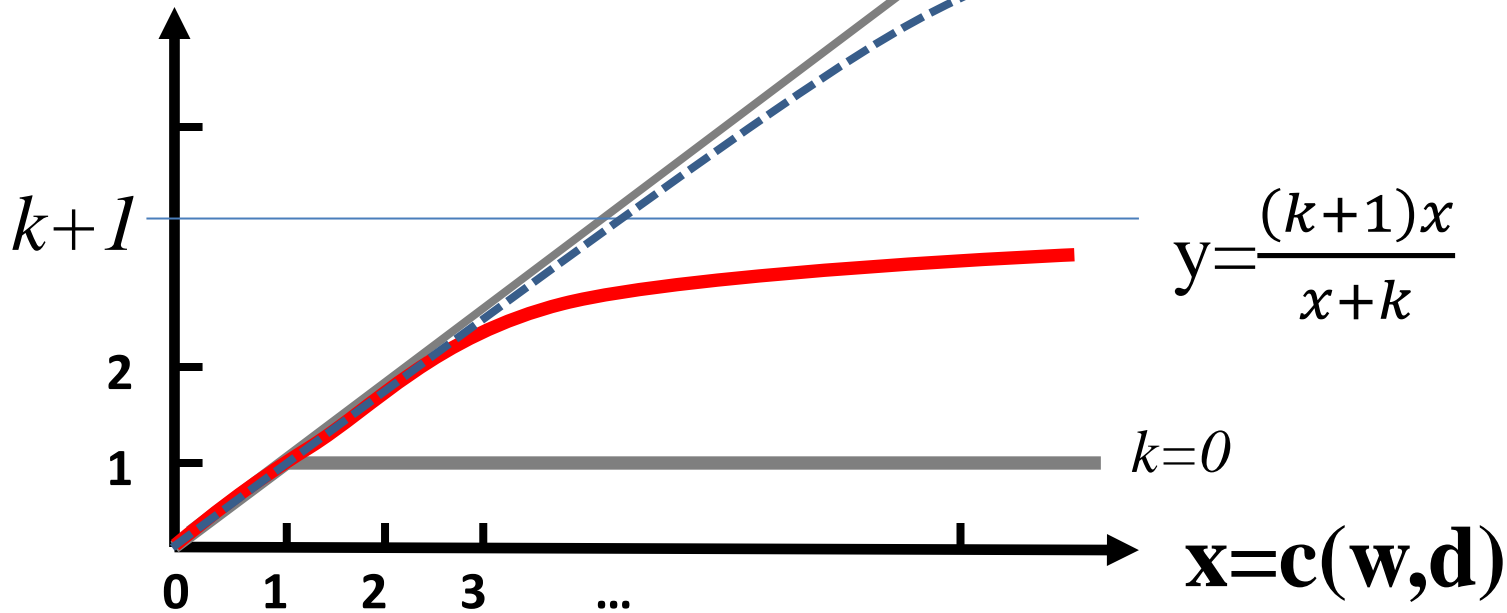
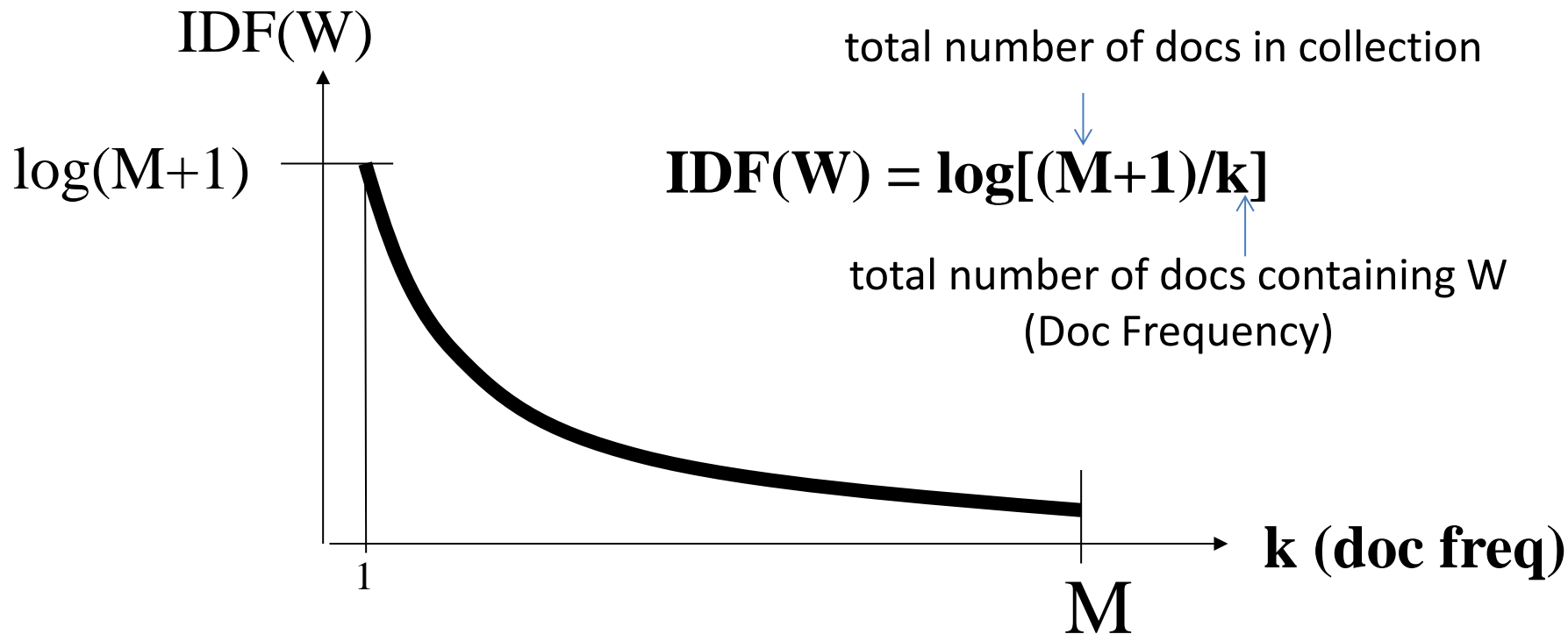# TF Transformation: c(w,d)➜TF(w,d)

Term Frequency Weight

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

$y= x$

$y= \log(1+x)$

2

1

0/1 bit vector
(ignore counts)

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

0   1   2   3      …

# TF Transformation: BM25 Transformation

Term Frequency Weight

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

$k+1$

$2$

$1$

$0 \quad 1 \quad 2 \quad 3 \quad \ldots$

*Very large k*

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

*k=0*

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

# IDF Weighting: Penalizing Popular Terms



$IDF(W)$

$\log(M+1)$

1

$M$

k (doc freq)

total number of docs in collection

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

total number of docs containing W
(Doc Frequency)

# Adapting BM25 Retrieval Model for Paradigmatic Relation Mining

$d1 = (x_1, \ldots x_N)$

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

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

$$b \in [0, 1]$$
$$k \in [0, +\infty)$$

$d2 = (y_1, \ldots y_N)$    $y_i$ is defined similarly

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

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

**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"