



Information Retrieval & Text Mining

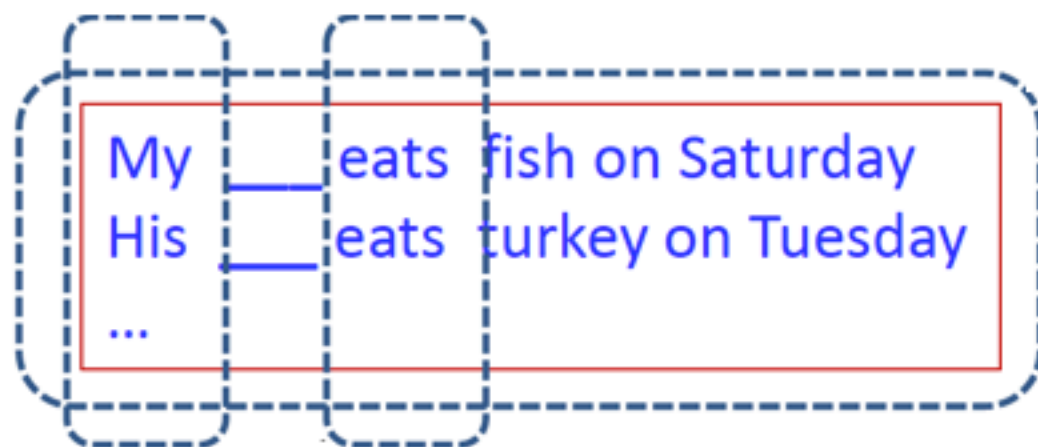
Paradigmatic Relation Discovery

Dr. Iqra Safder

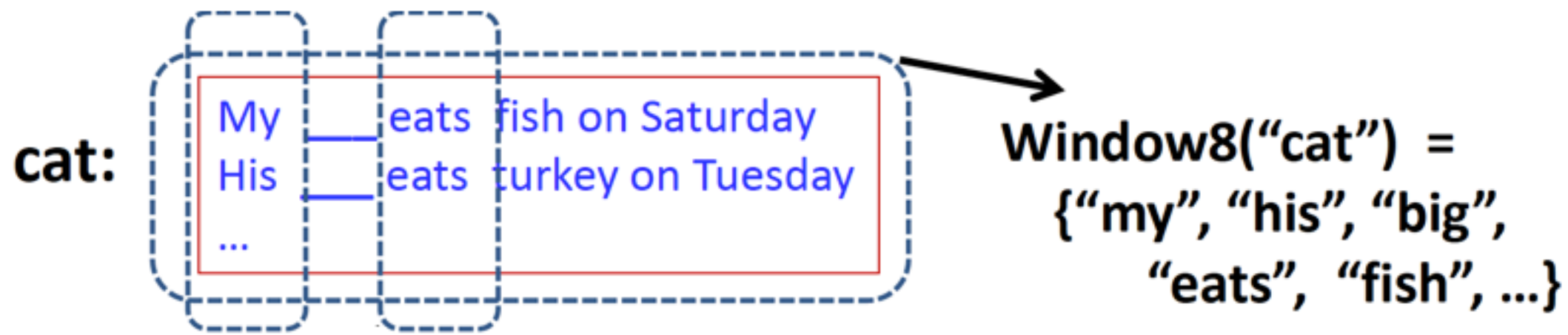
Information Technology University

Word Context as “Pseudo Document”

cat:

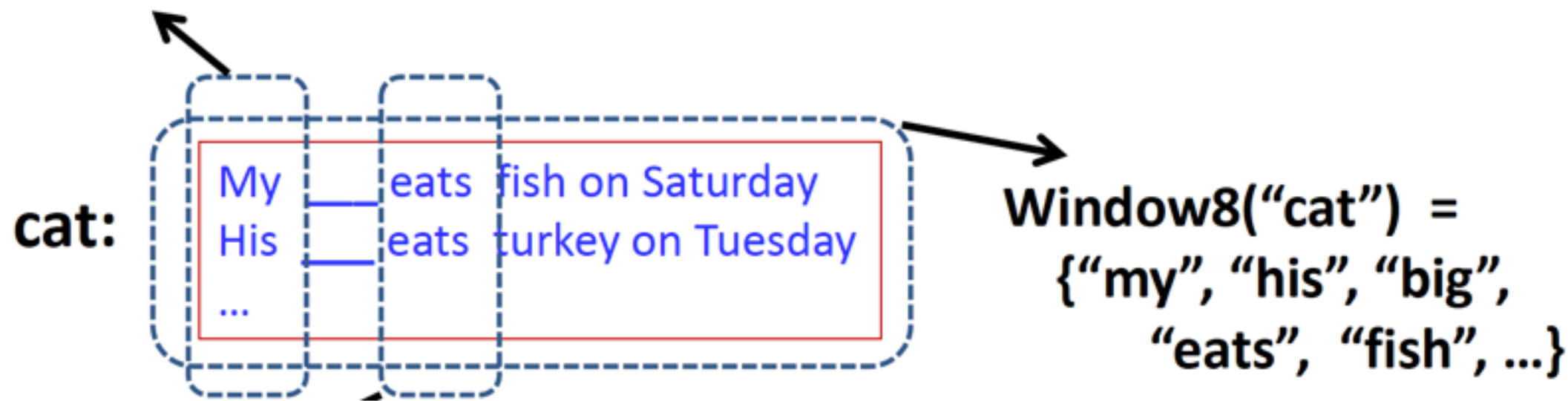


Word Context as “Pseudo Document”



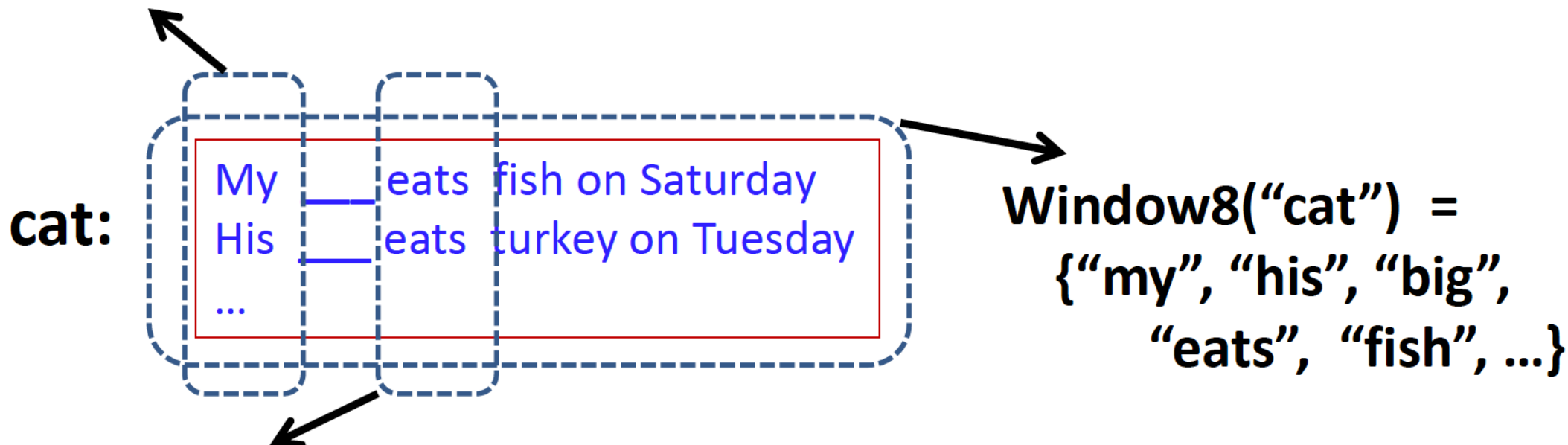
Word Context as “Pseudo Document”

$\text{Left1}(\text{"cat"}) = \{\text{"my"}, \text{"his"}, \text{"big"}, \text{"a"}, \text{"the"}, \dots\}$



Word Context as “Pseudo Document”

$\text{Left1}(\text{"cat"}) = \{\text{"my"}, \text{"his"}, \text{"big"}, \text{"a"}, \text{"the"}, \dots\}$



$\text{Window8}(\text{"cat"}) = \{\text{"my"}, \text{"his"}, \text{"big"}, \text{"eats"}, \text{"fish"}, \dots\}$

$\text{Right1}(\text{"cat"}) = \{\text{"eats"}, \text{"ate"}, \text{"is"}, \text{"has"}, \dots\}$

Context = pseudo document = “bag of words”
Context may contain adjacent or non-adjacent words

Measuring Context Similarity

Sim("Cat", "Dog") =

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

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

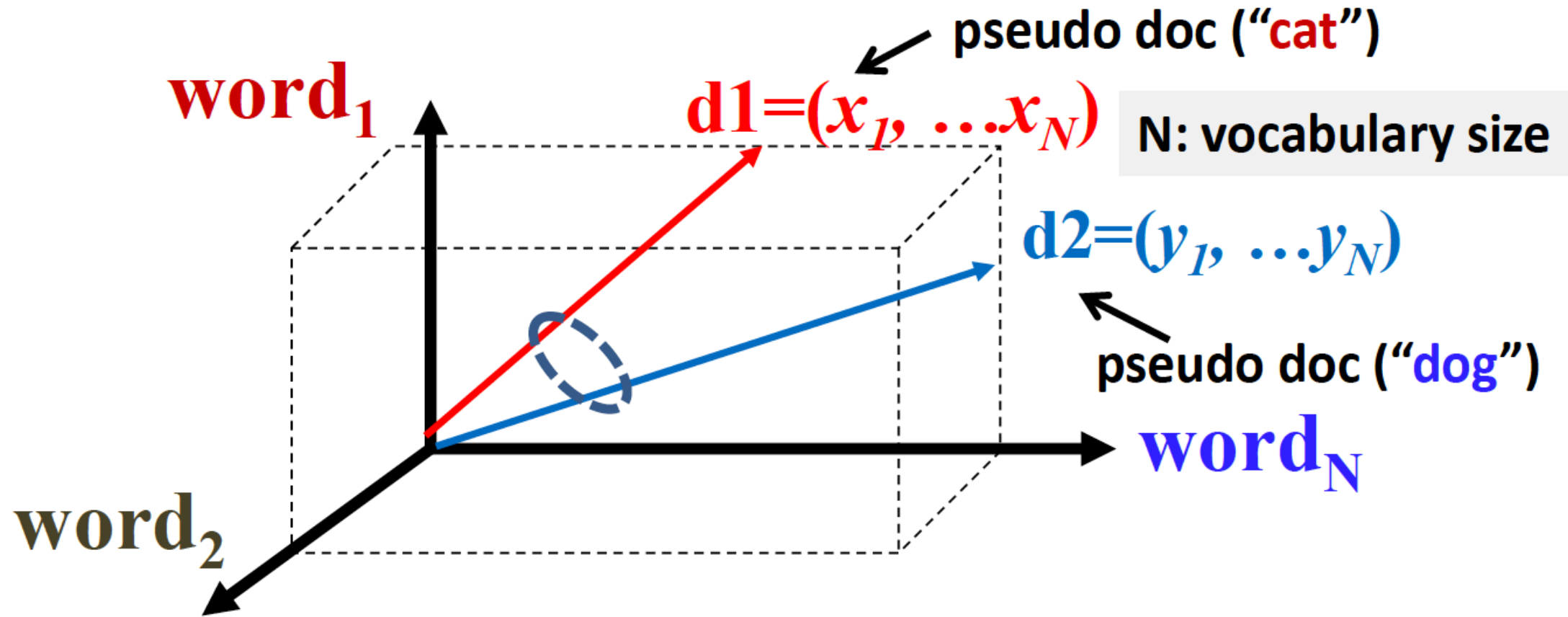
...

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

High sim(word1, word2)

→ word1 and word2 are paradigmatically related

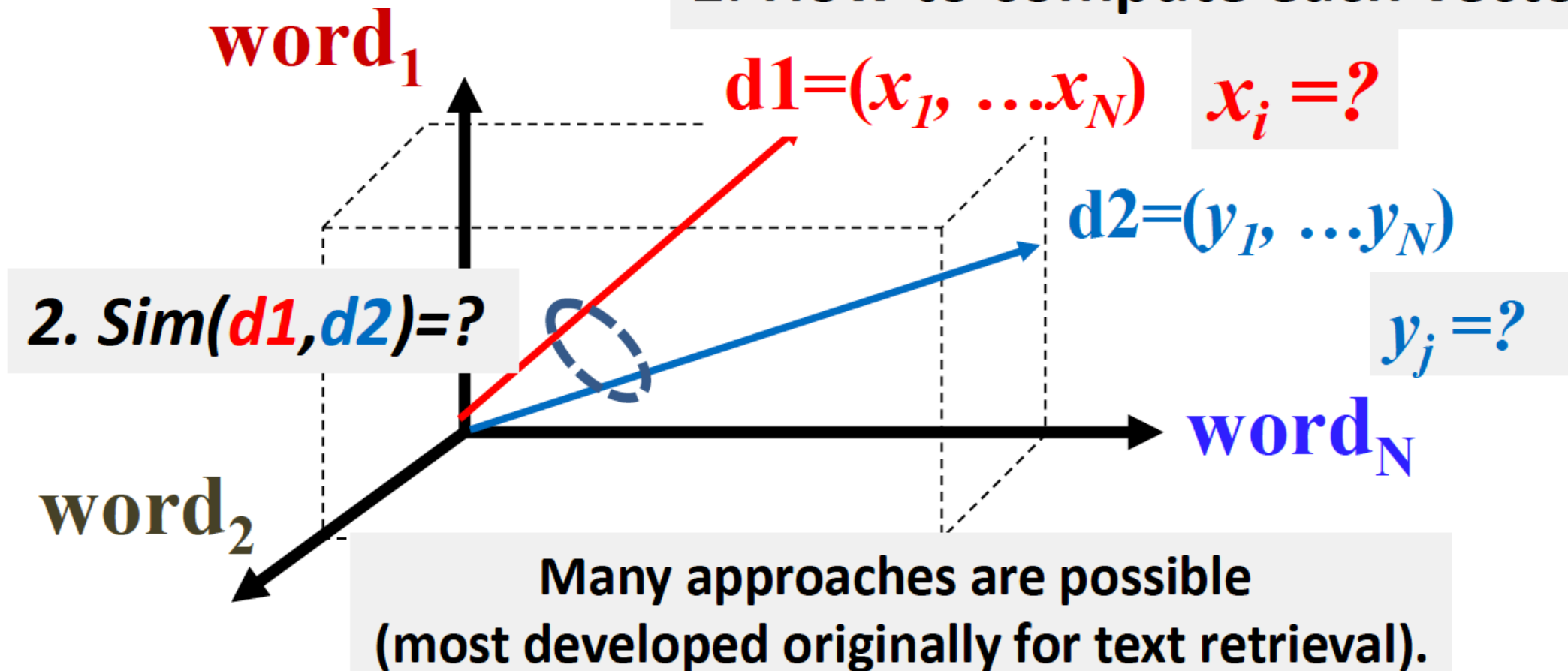
Bag of Words \rightarrow Vector Space Model (VSM)



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

VSM for Paradigmatic Relation Mining

1. How to compute each vector?



Expected Overlap of Words in Context (EOWC)

Probability that a randomly
picked word from d1 is w_i

Count of word w_i in d1

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

$$x_i = c(w_i, d1) / |d1|$$

$$d2 = (y_1, \dots, y_N)$$

$$y_i = c(w_i, d2) / |d2|$$

Total counts of
words in d1

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

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”).

Improving EOWC with Retrieval Heuristics

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

➔ Sublinear transformation of Term Frequency (TF)

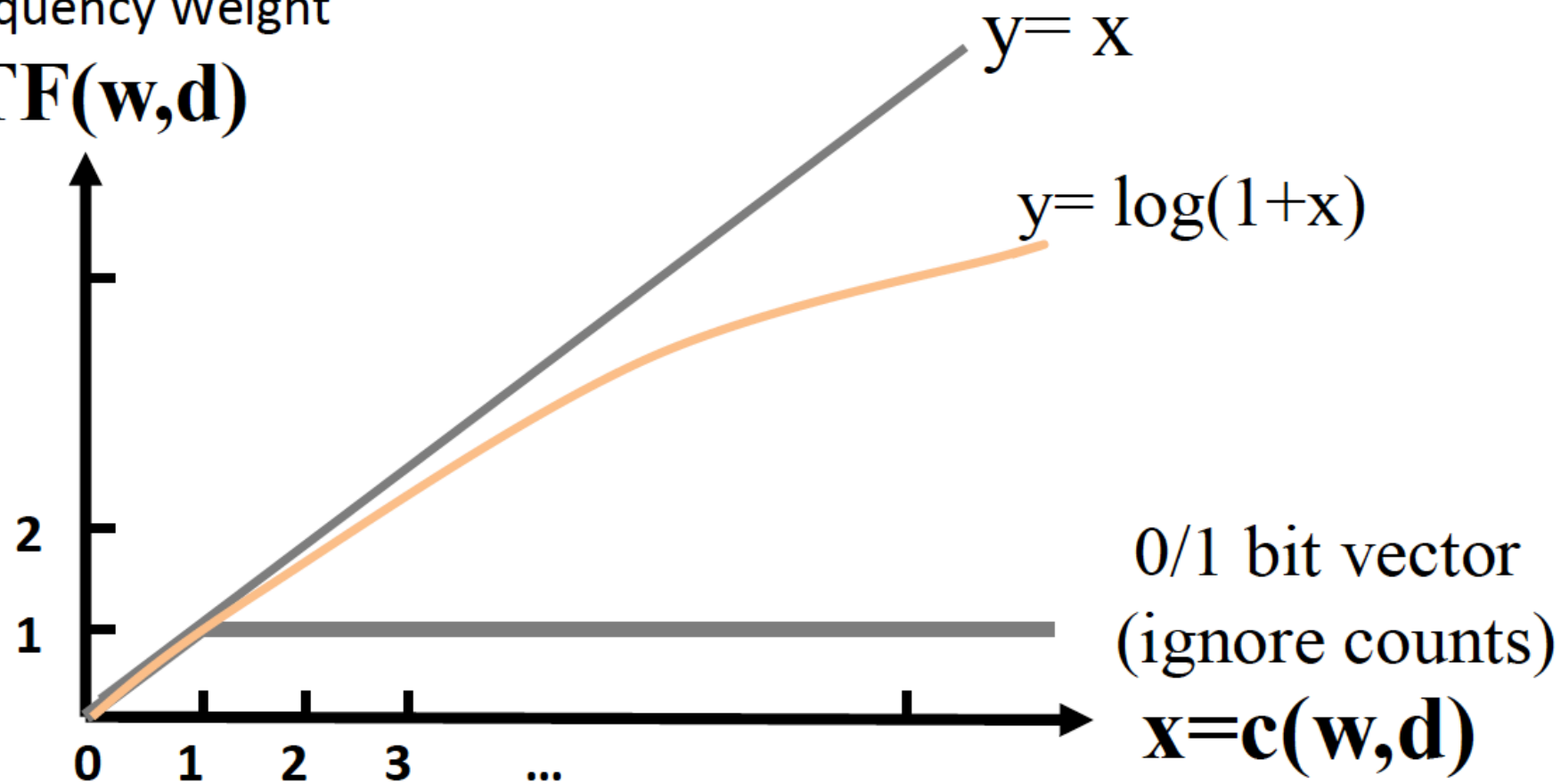
- It treats every word equally (overlap on “the” isn’t as so meaningful as overlap on “eats”).

➔ Reward matching a rare word: IDF term weighting

TF Transformation: $c(w,d) \rightarrow TF(w,d)$

Term Frequency Weight

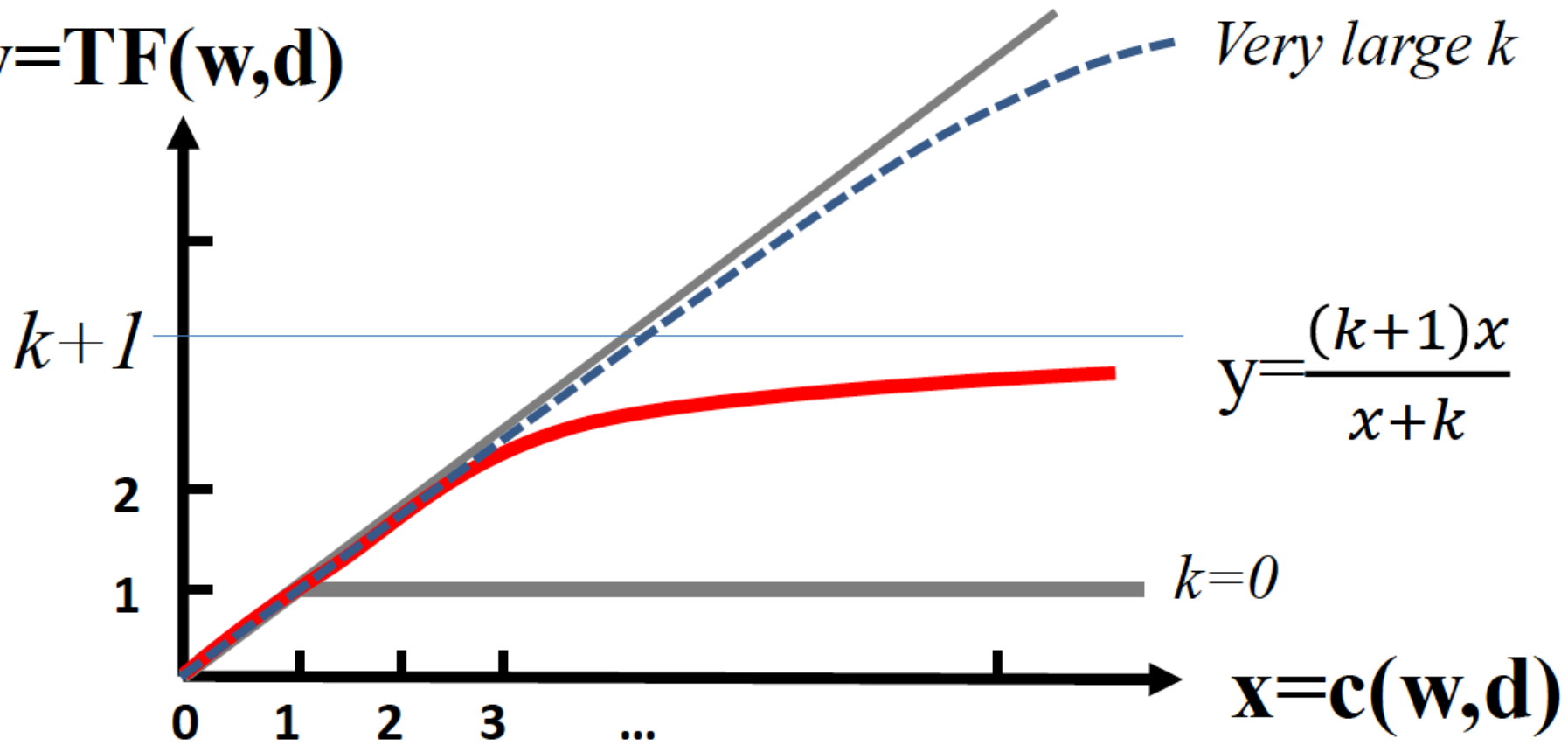
$$y = TF(w,d)$$



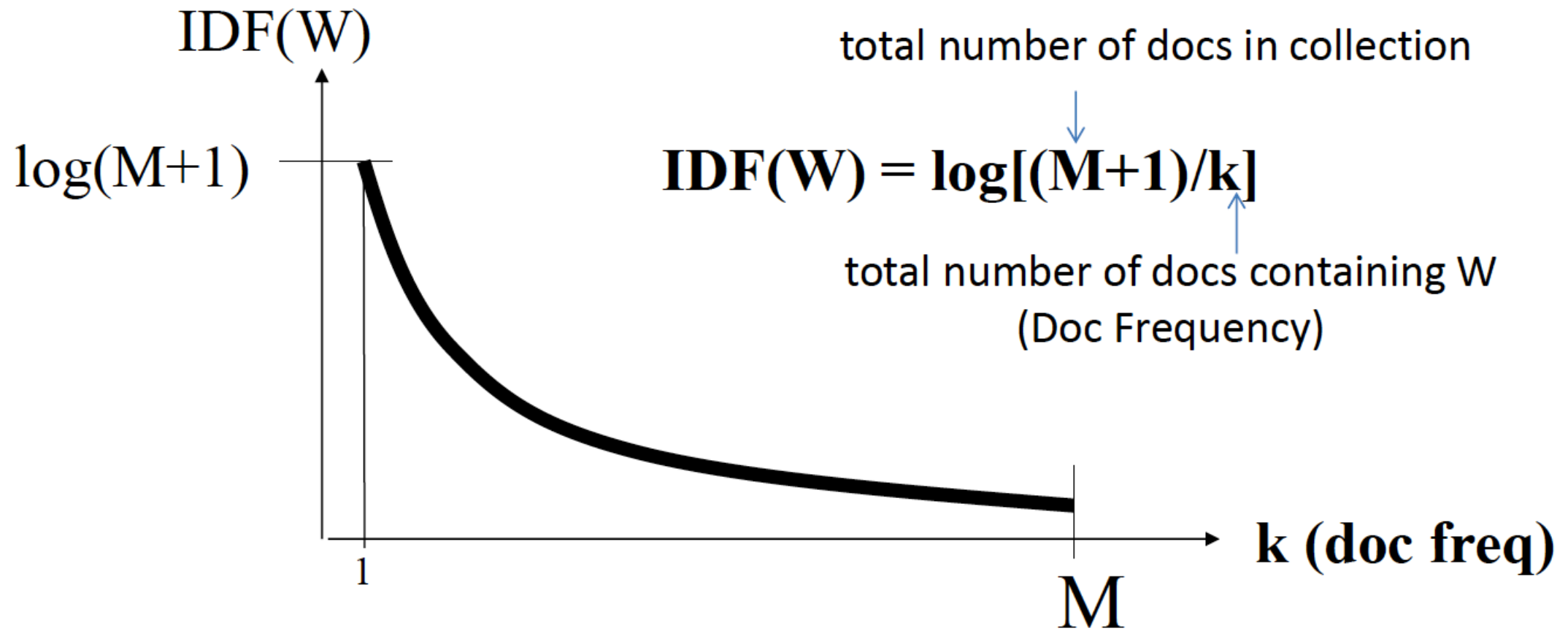
TF Transformation: BM25 Transformation

Term Frequency Weight

$$y = \text{TF}(w, d)$$



IDF Weighting: Penalizing Popular Terms



Adapting BM25 Retrieval Model for Paradigmatic Relation Mining

$$d1=(x_1, \dots x_N) \quad BM25(w_i, d1) = \frac{(k+1)c(w_i, d1)}{c(w_i, d1) + k(1-b+b*|d1|/avdl)}$$

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

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

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

$$d2=(y_1, \dots y_N) \quad y_i \text{ is defined similarly}$$

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

BM25 can also Discover Syntagmatic Relations

$$d1=(x_1, \dots x_N) \quad \text{BM25}(w_i, d1) = \frac{(k+1)c(w_i, d1)}{c(w_i, d1) + k(1-b+b*|d1|/avdl)}$$

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

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

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

$$\text{IDF-weighted } d1=(x_1 * \text{IDF}(w_1), \dots, x_N * \text{IDF}(w_N))$$

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”