



# 10

# Text Mining

Munawar, PhD

# Definition

- ❖ Text mining also is known as Text Data Mining (TDM) and Knowledge Discovery in Textual Database (KDT).[1]
- ❖ A process of identifying **novel information** from a collection of texts (also known as a corpus). [2]

# Data Mining vs. Text Mining

## ❖ Data Mining

- process directly
- Identify causal relationship
- Structured numeric transaction data residing in relational data warehouse

## ❖ Text Mining

- Linguistic processing or natural language processing (NLP)
- Discover heretofore unknown information[2]
- Applications deal with much more diverse and eclectic collections of systems and formats[4]

# Disputation

## ❖ Hearst: non-novel, novel.

- Text mining is not a simple extension of data mining applied to unstructured database.
- Text mining is the process of mining precious nuggets of ore from a mountain otherwise worthless rock.

## ❖ Kroeze: non-novel, semi-novel, and novel

- Non-novel: data/information retrieval
- Semi-novel: knowledge discovery (standard data-mining, metadata mining, and standard text mining)
- Novel: intelligent text mining

# Confusion

- ❖ Is text mining the same as information extraction? No!
- ❖ Information Extraction (IE)
  - Extract facts about pre-specified entities, events or relationships from unrestricted text sources.
  - **No novelty**: only information already present is extracted.

# Two Foundations

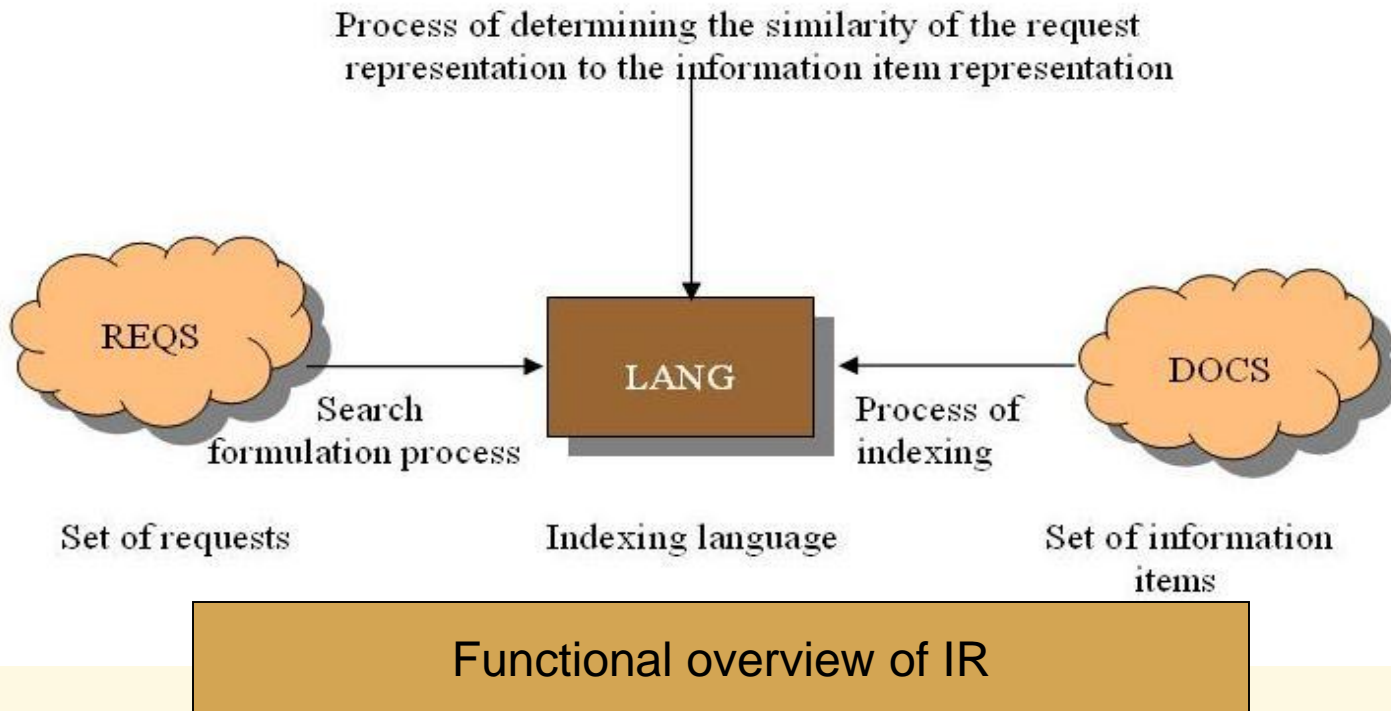
- ❖ Information Retrieval (IR)
- ❖ Artificial Intelligence (AI)

# Information Retrieval

- ❖ The science of searching for
  - Information in documents
  - Documents themselves
  - Metadata which describe documents
  - Text, sound, images or data, within database: relational stand-alone database or hypertext networked databases such as the Internet or intranets.

# Information Retrieval

## ❖ Gerard Salton



# Application I

## ❖ Semi-novel

## ❖ Text clustering: group similar documents for further examination -> create thematic overviews of text collections

### Issues:

- Information needs is vague
- Even if a topic were available, the words used to describe it may not be known to the user
- The words used to describe a topic may not be those used to discuss the topic and may thus fail to appear in articles of interest.
- Even if some words used in discussion of the topic were available, documents may fail to use precisely those words. [5]

# Text Clustering

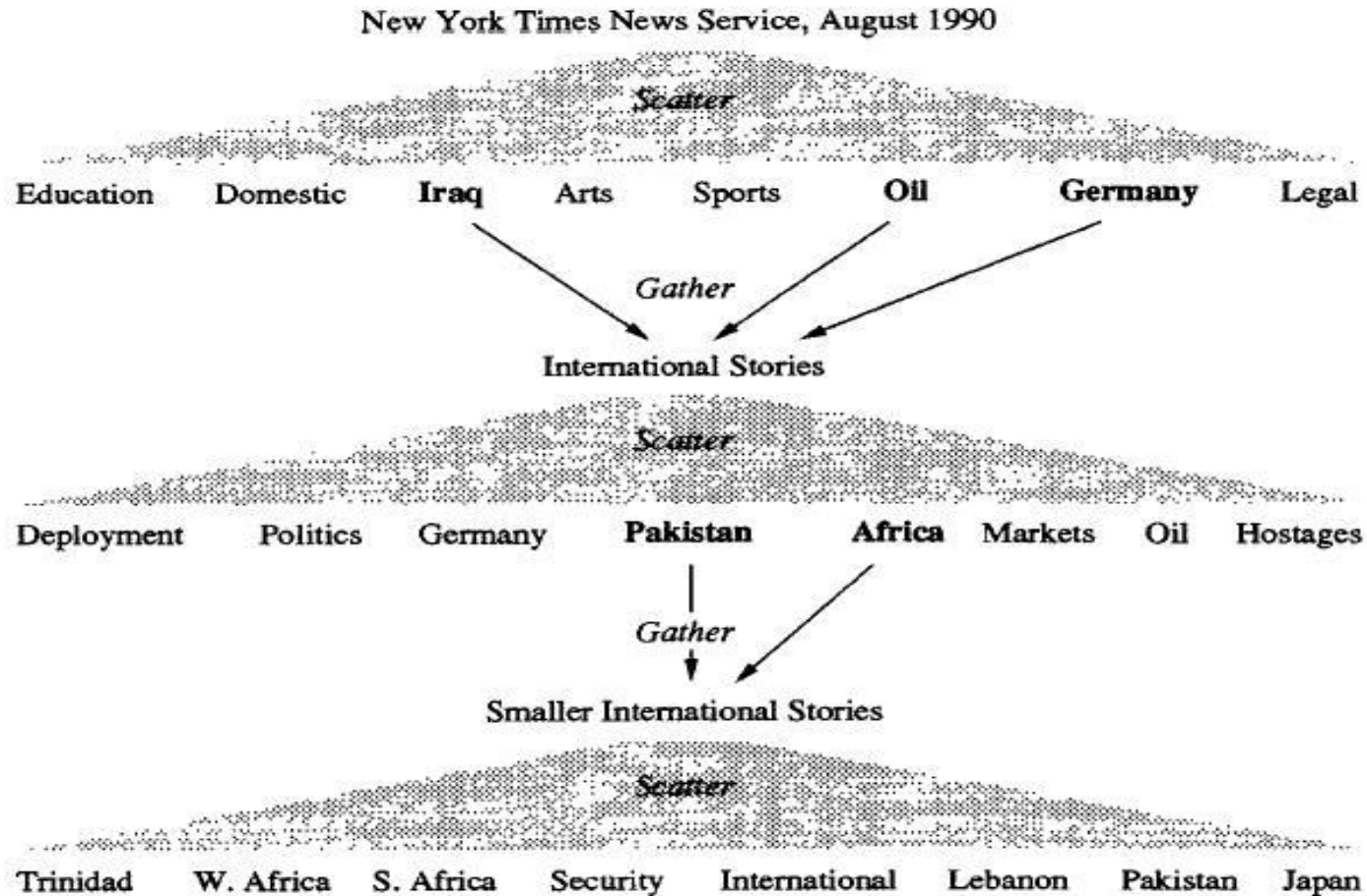


Figure 1: Illustration of Scatter/Gather

# Application II

- ❖ Semi-novel
- ❖ Automatically generating term associations to aid in query expansion
  - Word mismatch
- ❖ Clustering
- ❖ Global / local analysis[7]

# Global vs. Local Analysis

## ❖ Global Analysis

- Expensive
- Clustering based on all documents
- Provides a thesaurus-like resource

## ❖ Local Analysis

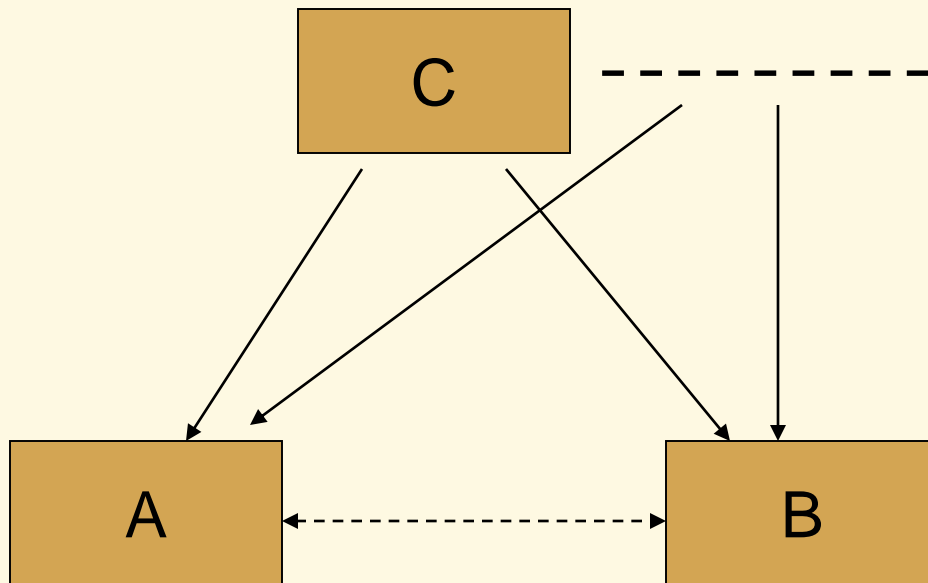
- Cost-efficient
- Clustering based on documents returned from previous query
- Only provides a small test collection

# Application III

- ❖ Semi-novel
- ❖ Using co-citation analysis to find general topics within a collection or identify central web pages

# Co-citation

- ❖ Bibliographic Co-Citation is a popular similarity measure used to establish a subject similarity between two items.
- ❖ E.g. basic idea of Google's algorithm



# Co-citation Analysis Steps

- ❖ Selection of the core set of items for the study.
- ❖ Retrieval of co-citation frequency information for the core set.
- ❖ Compilation of the raw co-citation frequency matrix.
- ❖ Correlation analysis to convert the raw frequencies into correlation coefficients.
- ❖ Multivariate analysis of the correlation matrix, using principle components analysis, cluster analysis or multidimensional scaling techniques.
- ❖ Interpretation of the resulting "map" and validation.

# The Raw Co-citation Frequency Matrix

Raw Cocitation Frequency

ID	
31	
32	98 32
33	74 128 33
34	27 181 25 34
35	41 182 91 123 35
36	4 13 24 1 5 36
37	9 83 19 10 8 8 37
38	8 18 14 12 4 5 5 38
39	127 400 113 120 138 20 700 15 39
310	13 83 20 119 31 1 19 17 75 310
311	88 200 200 71 151 25 130 28 400 97 311
312	3 3 8 2 17 2 0 3 1 2 17 312
313	72 182 91 29 121 19 200 8 1000 23 300 8 313
314	83 185 82 108 109 8 18 11 109 90 173 3 85 314
315	132 181 47 184 25 40 58 28 800 50 154 5 200 38 315
316	5 20 35 15 17 20 33 8 44 18 92 8 45 8 48 316
317	75 180 150 188 110 14 119 32 700 185 200 24 184 85 300 108 317
318	1 28 7 12 8 12 0 4 13 5 22 0 8 4 48 8 11 318
319	27 13 5 88 2 5 8 3 88 4 11 0 18 45 75 5 51 2 319
320	8 44 22 14 5 51 20 10 47 12 40 0 78 10 142 12 23 39 7 320
321	20 53 37 101 15 2 95 3 200 38 43 1 42 81 170 20 200 3 85 19 321
322	1 8 4 0 2 3 1 1 2 4 7 1 4 5 3 4 4 0 0 0 1 322
323	11 89 78 12 15 18 31 21 88 17 75 2 45 23 91 33 200 3 8 23 83 3 323
324	5 14 74 17 31 4 3 8 15 20 82 11 15 15 18 23 73 0 9 5 8 2 10 324
325	5 38 19 8 11 5 128 0 187 27 124 1 130 35 42 31 85 1 1 3 23 2 29 4 325
326	55 85 75 3 38 42 2 5 59 2 108 2 88 80 58 11 18 11 4 70 2 2 8 8 2 326
327	48 131 59 39 31 33 58 9 159 31 132 2 148 48 200 52 123 29 5 51 35 5 59 20 34 37 327
328	80 175 137 33 84 9 54 22 300 72 200 3 101 192 82 11 400 7 8 27 111 5 37 22 24 84 108 328
329	49 112 81 58 88 18 7 10 88 28 118 2 88 117 32 13 28 8 2 23 4 2 9 10 8 58 53 108 329
330	8 39 4 44 28 0 1 4 33 43 54 2 8 80 18 3 78 2 2 3 11 2 8 7 17 1 12 38 13 330
331	138 300 148 185 128 20 50 23 300 80 300 8 200 187 400 48 194 29 8 32 43 5 58 30 47 92 300 200 120 87 331
332	70 200 99 128 83 48 52 58 800 58 200 11 200 93 800 42 300 73 44 132 108 8 100 40 27 100 157 153 98 29 145 332
333	8 40 2 117 24 2 2 2 28 10 12 0 11 18 43 4 21 1 10 4 18 0 10 3 3 3 4 5 4 9 27 21 333
334	88 151 152 25 128 17 13 24 109 59 200 18 95 119 115 28 157 18 8 27 14 2 70 52 18 105 88 127 88 37 175 81 9



# Artificial Intelligence

- ❖ **Artificial intelligence (AI)** is a branch of computer science and engineering that deals with intelligent behavior, learning, and adaptation in machines.

# Self-Organizing Map (SOM)

- ❖ One category of neural network models.
- ❖ Neighboring cells in a neural network compete in their activities by means of mutual lateral interactions, and develop adaptively into specific detectors of different signal patterns.
- ❖ Here, learning is called competitive, unsupervised, or self-organizing.[6]

# Self-Organizing Map (SOM)

## ❖ An example of how SOM works?

- Each node has two vectors: input vector, weight vector (location)
- Input vector (Red, Green, Blue)
- Red (255, 0, 0) Green (0, 255,0) Blue (0, 0, 255)
- Randomize the weight vectors of nodes in the map
- Calculate the Euclidean distance formula to find out the smallest difference between the input vector and the weight vector (Best Matching Unit or Winner)
- Pulling neighbors closer to the input vector
- Repeat a large number of cycle.

# Self-Organizing Map (SOM)

- ❖ In attempting to devise neural network models for linguistic representation, the difficulty is how to find metric distance relations between symbolic items.

# SOM Application

## ❖ Web Analysis

- Problem: directory-based search engines such as Yahoo! analyze, index and categorize web content manually.
- Solution:
  - high-precision noun phrase indexing was performed to each page
  - A vector space model of noun phrases and their associated weights were used to present each page
  - All pages were categorized by a SOM clustering program[4]

# Folksonomy

- ❖ **folksonomy** is an Internet-based information retrieval methodology consisting of collaboratively generated, open-ended labels that categorize content such as Web pages, online photographs, and Web links.

# Folksonomy

## ❖ Benefits:

- Lower content categorization costs
- respond quickly to changes and innovations
- the capacity of its tags to describe the "aboutness" of an Internet resource

# Folksonomy

- ❖ Lack of standard: polysemy, synonym
- ❖ Meta noise: inaccurate or irrelevant metadata

# Natural Language Processing

- ❖ **Natural language processing (NLP)** is a subfield of artificial intelligence and linguistics. It studies the problems of automated generation and understanding of natural human languages.
- ❖ Statistical natural language processing uses stochastic, probabilistic and statistical methods to resolve some of the difficulties : e.g. text segmentation, word sense disambiguation

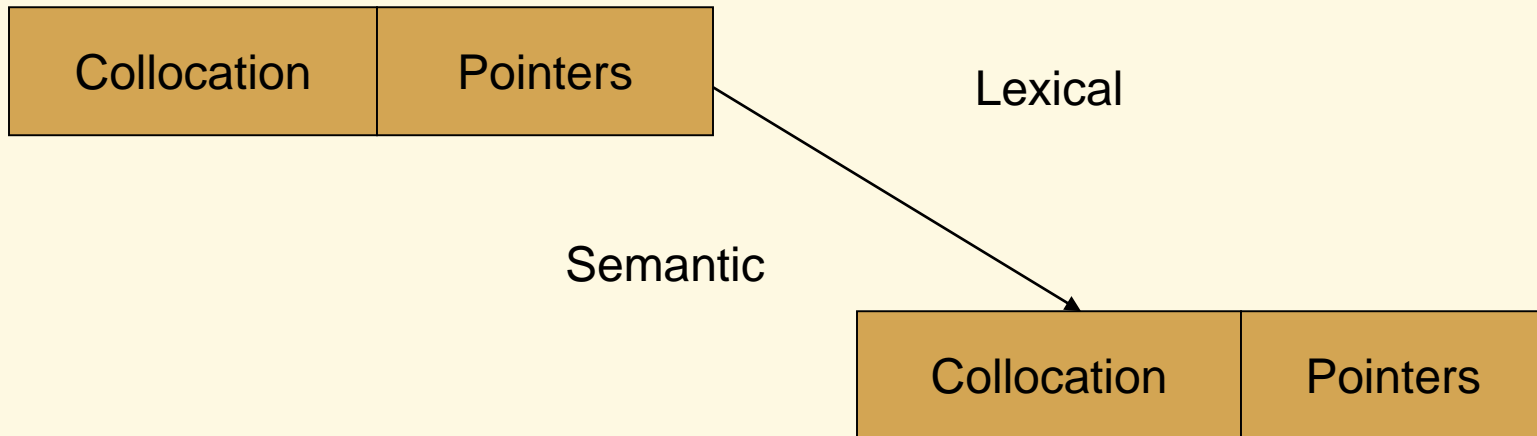
# Application I

❖ Semi-novel

❖ WordNet

<http://wordnet.princeton.edu/>

❖ Structure:



# WordNet

- ❖ Nouns and verbs are organized into hierarchies based on the hypernymy or hyponymy.
- ❖ Adjectives are arranged in clusters containing head synsets and satellite synsets:
  - Head synsets: clusters of antonymous pairs/triplets
  - Satellite synsets: a concept similar in meaning to that of the head synset

# WordNet

- ❖ A synset for a pertainyms contains only one word or collocation and a lexical pointer to the noun that the adjective is "pertaining to".
- ❖ A synset for an adverb contains a lexical pointer to the adjective from which it is derived.

# Glossary

- ❖ Hypernymy: **Y** is a hypernym of **X** if **X** is a (kind of) **Y**
- ❖ Hyponymy: **X** is a hyponym of **Y** if **X** is a (kind of) **Y**

# Application II

- ❖ Novel
- ❖ Bioscience: deduce hypotheses (National Centre for Text Mining in UK)

# An Example

- ❖ When investigate causes of migraine headaches, we extracted various pieces of evidence from titles of articles in the biomedical literature:
  - stress is associated with migraines
  - stress can lead to loss of magnesium
  - calcium channel blockers prevent some migraines
  - magnesium is a natural calcium channel blocker
  - spreading cortical depression (SCD) is implicated in some migraines
  - high levels of magnesium inhibit SCD
  - migraine patients have high platelet aggregability
  - magnesium can suppress platelet aggregability

# Application III

## ❖ Novel

## ❖ Using text to uncover social impact

- the effects of publicly financed research on industrial advances
- 1987 ~ 1988, 1993 ~ 1994; 397,660 patents issued
- found 242,000 identifiable science references and zeroed in on those published in the preceding 11 years (80%)
- 109,000 of these references to known journals and authors' addresses in computer database
- a core collection of 45,000 papers, after eliminating redundant citations to the same paper, and articles with no known American author
- look up the papers and examine their closing lines, which often say who financed the research

# Reference

- ❖ [1] Ronen Feldman and Ido Dagan. “Knowledge discovery in textual databases (KDT)”, *Knowledge Discovery and Data Mining* (1995), 112-117.
- ❖ [2] M. Hearst. “Untangling text data mining”. *Proceedings of ACL'99: the 37th Annual Meeting of the Association for Computational Linguistics*, 1999.
- ❖ [3] chen 2002

# Reference

- ❖ [4] Hsinchun Chen, "Knowledge Management Systems: A Text Mining Perspective", *Knowledge Computing Corporation*, 2001.
- ❖ [5] Douglass R. Cutting, Jan O. Pedersen, David Karger, and John W. Tukey. "Scatter/Gather: A cluster-based approach to browsing large document collections". *In Proceedings of the 15th Annual International ACM/SIGIR Conference 1992*, pages 318-329, Copenhagen, Denmark.

# Reference

- ❖ [6] Teuvo Kohonen. "The Self-Organizing Map". Proceedings of the IEEE (1990), 78, 9, 1464-1480.
- ❖ [7] J. Xu and W. B. Croft. 1996. Query expansion using local and global document analysis. In *SIGIR '96: Proceedings of the 19th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 4-11, Zurich.

# Reference

- ❖ [8] Ray R. Larson. Bibliometrics of the world wide web: An exploratory analysis of the intellectual structure of cyberspace. In *ASIS '96: Proceedings of the 1996 Annual ASIS Meeting*.



# Thank You !

**Munawar, PhD**