# Master of Technology in Knowledge Engineering

## Text Mining

## Advanced Topics in Text Mining

**Rajaraman Kanagasabai**
**Institute for Infocomm Research**
email: kanagasa@i2r.a-star.edu.sg

NUS National University of Singapore | ISS INSTITUTE OF SYSTEMS SCIENCE

# Overview

- Introduction to Advanced Text mining

  » The Semantic Gap in Text mining

  » Traditional Approaches to Text Mining

    ♦ Keywords, NLP, & Ontologies

  » The Deep Learning Approach

- Text Mining with Deep Learning

  » Background

  » Major DL Approaches to Text Mining

- Case studies

- Workshop

  » Hands-on Exercises with Word2Vec

NUS National University of Singapore | ISS INSTITUTE OF SYSTEMS SCIENCE

# What is Text Analytics?

- *"discovery of **new** previously unknown information, by **automatically** extracting information from a usually **large amount** of different unstructured textual resources"*

- Other definitions
  - » Use of **computational techniques** to extract high quality information from text
  - » Extract and discover knowledge hidden in text **automatically**

# Key Text Analytics Problems

- **Analyze Document Collections**

  » Information Retrieval

  » Classification (Supervised Learning)

  » Clustering (Unsupervised Learning)

- **Analyze Document**

  » Summarization

  » Information Extraction: Extract Names, Relations, Facts

- **Analyze Sentence**

  » Sentiment Analysis

  » Co-reference Resolution
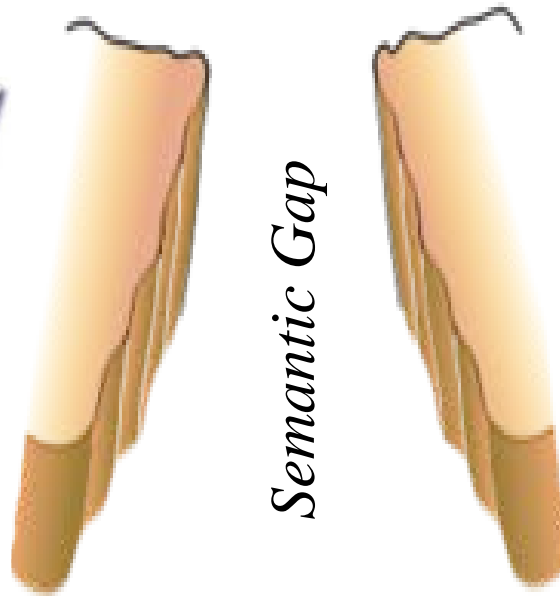
# Why not just apply standard DA?

"Text expresses a vast, rich range of information, but encodes this information in a form that is difficult to decipher automatically."

-- Marti A. Hearst,
"Untangling Text Data Mining," 1999

# Key Issue



Raw Text

*Semantic Gap*

Text Understanding

# Example 1

- *It would be a great <u>help</u> if you can <u>assist</u> me, and I will appreciate the <u>favor</u>*


*=> Known as Synonymy*

# Example 2

- *"I <u>banked</u> on him to meet at the Deutche <u>Bank</u> located near the river <u>bank</u>", said Mr. <u>Banks</u>.*

*=> Known as Polysemy*

# Overview of Approaches

# Approaches

- Keywords-based

  » Statistical

  » Zero semantics approach

- NLP based

  » Language parsing

  » Syntactic approach

- Ontology based

  » Uses formal logic representations

- Deep Learning based

  » Deep neural architectures & ML

Decreasing Semantic Gap
Increasing Accuracy

Increasing Complexity
Decreasing Scalability

# **Keywords-based Approach**

# Bag-of-words document representation

# Bag of words Representation

- Texts are treated as a "bag" of words or terms

- Any document can be represented as a vector: a list of **terms** and their associated **weights**

  » D= {(t_1,w_1),(t_2,w_2),…………

  » $t_i$: i-th term

  » $w_i$: weight for the i-th term



VECTOR SPACE MODEL

# Cosine Similarity

$$sim(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$

# Problems with Keywords

- Bag of words approach

  » Is both the **advantage** and **disadvantage**

- Often results in **huge semantic gap**

  » A positive or negative sentiment word may have opposite orientations in different application domains. ("*This camera sucks.*" -> negative; "*This vacuum cleaner really sucks.*" -> positive)

  » A sentence containing sentiment words may not express any sentiment. (e.g. "*Can you tell me which Sony camera is good*?")

  » Sarcastic sentences with or without sentiment words are hard to deal with. (e.g. "*What a great car! It stopped working in two days.*"

# NLP-based Approach

# Natural Language Processing (NLP)

- **NLP**: is a field of computer science, artificial intelligence, and linguistics, concerned with the interactions between computers and human (natural) languages.

- Major NLP applications

  » Part-of-speech tagging (POS tagging)

  » Relationship extraction

  » Sentiment analysis

  » Topic segmentation and recognition

  » Machine translation

# Shallow NLP

- ## Part-of-Speech (POS)

  » Identify Nouns, Verbs, Tenses, Prepositions, etc.

- ## Morphology

  » Do stemming the right way

- ## Syntax

  » Extract sentence structure

# Demonstration: Sentence-level Sentiment – 1/3

- Stanford Sentiment Analyzer

  » http://nlp.stanford.edu:8080/sentiment/rntnDemo.html

# Demonstration: Sentence-level Sentiment – 2/3

- Review 1: This movie doesn't care about cleverness, wit or any other kind of intelligent humor. -> **Negative**

# Demonstration: Sentence-level Sentiment – 3/3

- There are slow and repetitive parts, but it has just enough spice to keep it interesting. -> **Positive**

# Problems with NLP

- Better than Bag-of-words but still not the best

  » synonymy and polysemy not entirely overcome

  » Basically, a syntactic approach

- Examples

  » "I studied in Cambridge" (Which Cambridge?)

  » "I live in Singapore" (Do you live in Asia?)

# Ontology-based Approach

# What Is An Ontology?

- Ontology (Socrates & Aristotle 400-360 BC)
- Word borrowed by computing for the explicit description of the conceptualisation of a domain:
  - » concepts
  - » properties and attributes of concepts
  - » constraints on properties and attributes
  - » Individuals (often, but not always)
- An ontology defines
  - » a common vocabulary
  - » a shared understanding
- E.g.: Wiki/Yahoo categories, WordNet
- Backbone of semantic web

# A simple ontology: Animals

# How is it useful for Text Mining?

"XYZ announced profits in Q3, planning to build a $120M plant in Bulgaria, ….. more and more text..."

# Information Extraction

"**XYZ** announced profits in **Q3**, planning to build a **$120M** plant in **Bulgaria**, ….. more and more text..."

*XYZ*
*Q3*
*$120M*
*Bulgaria*

Loosely structured.. "Semantic Gaps"

OK for simple scenarios, but messy in larger applications!

# With Proper Semantics

# What can we do with ontologies? (1/2)

- **Semantic Annotation**: represent metadata/keywords with proper semantics
  - » to represent 'Cambridge' as a UK location, link it to an ontology instance 'Cambridge, UK' rather than 'Cambridge, Massachusetts'
  - » Link synonyms to the ontology , e.g. link 'heart attack' to 'myocardial infarction'

# Example

# What can we do with ontologies? (2/2)

- **Semantic Search**: Use proper semantics captured in ontologies to retrieve more relevant results

  » search for 'heart attack' retrieves documents containing 'myocardial infarction'

  » search for 'South East Asia' retrieves documents that contain 'Singapore' or 'Thailand'

# Google's Semantic Search
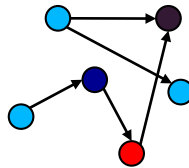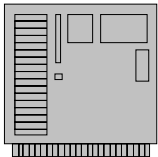
# Semantic Analytics Framework

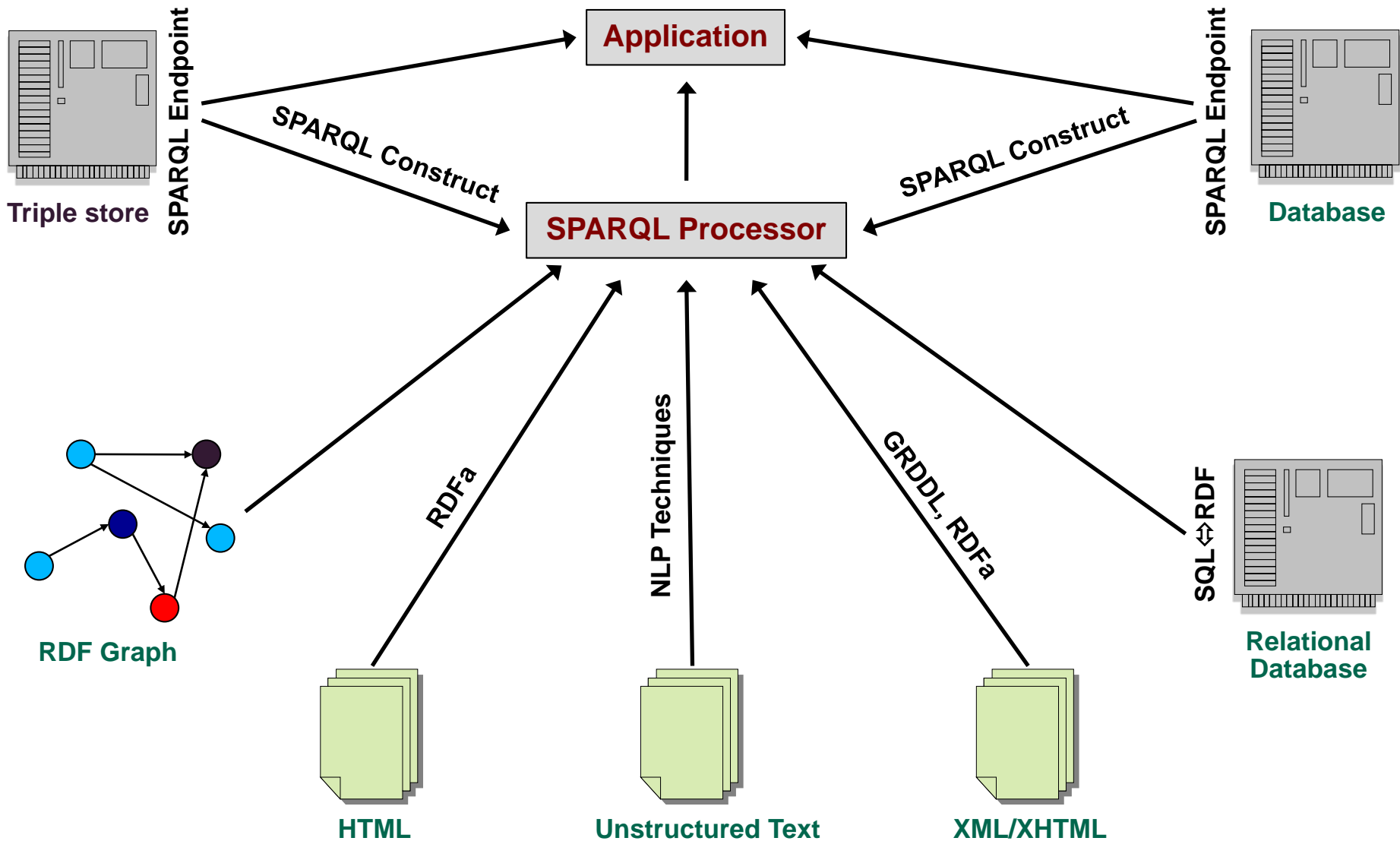**Applications**

**Manipulate Query …**

**Data represented in abstract format**

**Map, Expose, …**

**Data in various formats**

# Data Linking & Integration

# Summary

- Addressing the Semantic Gap is the key to effective Text Mining

- Keyword approach leads to a huge semantic gap, and hence often a big limitation in real life

- NLP addresses it somewhat but is still a syntactic approach

- Knowledge representations, as modeled by Ontologies, can bridge the Semantic gap but require a lot of manual effort in ontology construction