

Text Mining

- ❖ Data stored in most text databases are **semistructured** data
- ❖ A document may contain:
 - **structured** fields: title, authors, publication date, category, and,
 - largely **unstructured** text components: abstract and contents...
 - Information retrieval techniques: text indexing methods (to handle **unstructured** documents)

Information Retrieval

the science of searching for documents, for information within documents, and for metadata about documents, as well as that of searching relational databases and the World Wide Web.

Suppose that a text retrieval system has just retrieved a number of documents for me based on my input in the form of a query. How can we assess how accurate or correct the system was?

{Relevant}: the set of documents relevant to a query

{Retrieved}: the set of documents retrieved

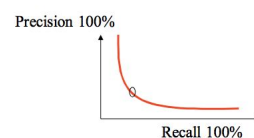
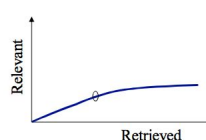
{Relevant} ∩ {Retrieved}: the set of documents that are both relevant and retrieved

- ❖ Measure quality of text retrieval:
 - Precision: This is the percentage of retrieved documents that are in fact relevant to the query (i.e., “correct” responses)
 - Recall: This is the percentage of documents that are relevant to the query and were, in fact, retrieved. It is formally defined as
 - F-score: Trade-Off (harmonic mean of precision & recall)
 - not sensitive to the internal **ranking** of the documents in a retrieved set

$$precision = \frac{|\{Relevant\} \cap \{Retrieved\}|}{|\{Retrieved\}|} \quad recall = \frac{|\{Relevant\} \cap \{Retrieved\}|}{|\{Relevant\}|}$$

+ IR Evaluation: Recall and Precision

	Retrieved	Not-Retrieved	
Relevant	A	B	Precision = A/(A+C)
Non-Relevant	C	D	Recall = A/(A+B)
			Fallout = C/(C+D)
			N _{documents} = A+B+C+D



$$F - score = \frac{recall \times precision}{(recall + precision)/2}$$

- ❖ Measure the quality of a **ranked** list of documents:

- compute an average of **precisions** at all the ranks where a new relevant document is returned
- **Or**, plot a graph of **precisions** at many different levels of **recall**; a **higher** curve represents a **better-quality** information retrieval system.

Nature Language Processing:

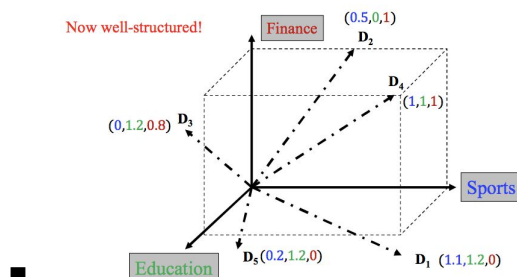
Part-of-speech Tagging	<div>This sentence serves as an example. Det Noun Verb P Det Noun</div>
Entity Recognition	<div>The University of Queensland, St. Lucia Brisbane University Suburb City</div>

Text Mining Challenges

- ❖ Data is not well-organized
- ❖ Natural language text contains ambiguities
- ❖ In traditional data mining: all data are “structured”

How to represent a document “structurally”?

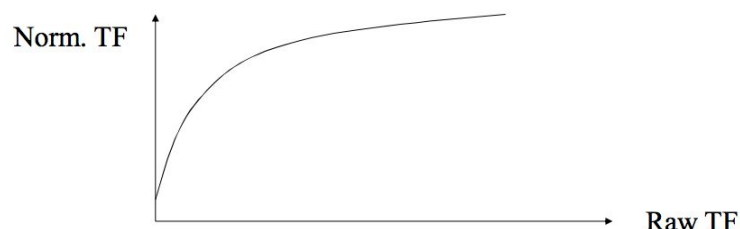
- ❖ Represent by a string?
 - No semantic meaning
- ❖ Represent by a list of sentences?
 - Sentence is just like a short document (recursive definition)
- ❖ Represent documents by concept vectors (**VSM: Vector Space Model**)
 - Each concept defines one dimension
 - k concepts define a k-dimensional space
 - A document is represented as a point in this k-dimensional space
 - E.g., $(1, \dots)$, is the coordinate of dimension i , and also the “importance” of concept in this document.



VSM doesn't Say:

- ❖ How to define/select the “**basic concept**”: Concepts are assumed to be **orthogonal**
 - Orthogonal:
 - Linearly independent basis vectors
 - “Non-overlapping” in meaning “sports” and “football”?
- ❖ How to assign **weights**: Weights indicate **how well** the concept characterizes the document. (Set automatically and accurately)
 - Document-wise:
 - not all terms are equally important
 - Corpus-wise:
 - some terms carry more information about the document content
 - How to set weight? (2 basic heuristics)
 - **TF** (Term Frequency) = Within-doc-frequency
 - Idea: a term is more important if it occurs more frequently in a document
 - **TF normalization (term frequency)**
 - Two views of document length:
 - ◆ A doc is long because it is verbose / more content
 - Raw TF is inaccurate
 - ◆ Document length variation
 - ◆ “Repeated occurrences” are less informative than the “first occurrence”
 - ◆ Information about semantic does not increase proportionally with number of term occurrence
 - Generally penalize long document, but avoid over-penalizing
 - ◆ Pivoted length normalization
 - Sub-linear TF scaling

$$tf(t, d) = \begin{cases} 1 + \log c(t, d), & \text{if } c(t, d) > 0 \\ 0, & \text{otherwise} \end{cases}$$



- **IDF** (Inverse Document Frequency)

$$IDF(word) = \log \frac{1 + \text{total documents}}{\text{document frequency}}$$

- Assign higher weights to the rare terms

- Formula

$$IDF(t) = \log \left(\frac{N+1}{df(t)} \right)$$

Total number of docs in collection
 Number of docs containing term t

- A corpus-specific property

- Independent of a single document

- TF-IDF weighting

- Combining TF and IDF

- Common in doc \rightarrow high tf \rightarrow high weight
 - Rare in collection \rightarrow high idf \rightarrow high weight
 - $w(t, d) = TF(t, d) \times IDF(t)$

❖ How to define the **distance metric**: What kind of documents are **similar**?

➤ 2 approaches: **Edit Distance** & **Cosine**

- Euclidean distance

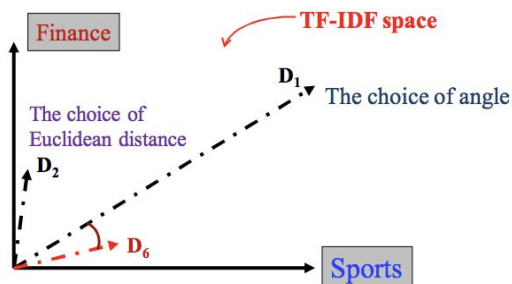
$$dist(d_i, d_j) = \sqrt{\sum_{t \in V} [tf(t, d_i)idf(t) - tf(t, d_j)idf(t)]^2}$$

- Longer documents will be penalized by the extra words

➤ ■ We care more about how these two vectors are overlapped

- Angle: how vectors are overlapped

➤ ■ Cosine similarity – projection of one vector onto another



+ Normalization

This is a data mining course.

(0.602, 0.301, 0, 0.125, 0, 0, 0)

$$w(\text{course}) = \frac{0.602}{\sqrt{0.602^2 + 0.301^2 + 0 + 0.125^2 + 0 + 0 + 0}} = 0.879$$

(0.879, 0.439, 0, 0.220, 0, 0, 0)

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}} = \mathbf{A} \cdot \mathbf{B}$$

ID	term	document frequency	IDF
1	course	1	0.602
2	data	2	0.301
3	interest	1	0.602
4	mine	3	0.125
5	study	1	0.602
6	subfield	1	0.602
7	text	2	0.301

Q: (0, 0.337, 0.879, 0, 0, 0, 0.337)

Document 1: (0.879, 0.439, 0, 0.220, 0, 0, 0)

Document 2: (0, 0.300, 0, 0.188, 0.601, 0.601, 0.390)

Document 3: (0, 0, 0.924, 0.148, 0, 0, 0.355)

ID	term	document frequency	IDF
1	course	1	0.602
2	data	2	0.301
3	interest	1	0.602
4	mine	3	0.125
5	study	1	0.602
6	subfield	1	0.602
7	text	2	0.301

$$\text{Sim}(Q, D1) = 0.337 \times 0.439 = 0.148$$

$$\text{Sim}(Q, D2) = 0.337 \times 0.300 + 0.337 \times 0.390 = 0.334$$

$$\text{Sim}(Q, D3) = 0.879 \times 0.924 + 0.337 \times 0.355 = 0.932$$

Conclusion: Return Document 3

VSM existing Solutions:

- ❖ Topic modelling
- ❖ Bag-of-Words model or N-gram model
 - Bag-of-Words (**BOW**): Each distinct word represents a concept (**most frequently** used **document representation**: Image, speech, gene sequence)
 - **Tokenization**: Break texts into meaningful words
 - Tokens: words, phrases, symbols. (Definition depends on language, corpus, or even context)
 - Rule-based
 - so-called -> 'so', 'called'
 - It's -> 'It's' instead of 'It', 's'
 - Statistical Method
 - 'San Francisco' instead of 'San', 'Francisco'
 - Existing Tool (eg. NLTK)
 - **N-grams** as concepts: a contiguous sequence of N tokens from a given piece of text

Bag-of-Words Advantages Vs. Disadvantages

Advantages	Disadvantages
<ul style="list-style-type: none"> Simple 	<ul style="list-style-type: none"> Words are clearly not linearly independent!

	<ul style="list-style-type: none"> • Grammar and order are missing
--	---

N-grams Advantages Vs. Disadvantages

Advantages	Disadvantages
<ul style="list-style-type: none"> • capture local dependency and order 	<ul style="list-style-type: none"> • A purely statistical view, increase the vocabulary size $O(V^N)$

Document Representation Advantages Vs. Disadvantages

Advantages	Disadvantages
<ul style="list-style-type: none"> • Preserve all information in the text (hopefully) • Fully automatic 	<ul style="list-style-type: none"> • Vocabulary gap: cars v.s., car, talk v.s., talking • Large storage: N-grams needs $O(V^N)$ <ul style="list-style-type: none"> ◦ Solution: Construct controlled vocabulary

❖ Normalisation

- Convert different forms of a word to a normalized form in the vocabulary
 - UQ -> University of Queensland, St. Lucia -> Saint Lucia
- Rule-based
 - Delete periods “.” and hyphens “-”
 - All in **lower cases**
 - Change it to original form
 - interesting -> interest
 - Mining -> mine
- Dictionary-based
 - Construct equivalent class
 - Car -> “automobile, vehicle”
 - UQ -> “University of Queensland”
 - Mobile phone -> “cellphone”

❖ Stemming

- Reduce inflected or derived words to their root form
 - Plurals, adverbs, inflected word forms
 - E.g., ladies -> lady, referring -> refer, forgotten -> forget
 - Reduce the vocabulary size
 - Existing tools (for English)
- Risk: lose precise meaning of the word

❖ Controlled vocabulary

- Remove non-informative & rare words

➤ Zipf's law

- If the word ranked 1st has frequency 10000, what is the frequency of the word ranked 20-th ?
- $f(1)/1 = f(20)/(1/20) \rightarrow f(20) = f(1)/20 = 500$
- Head words take large portion of occurrences, but they are semantically meaningless
 - E.g., the, a, an, we, do, to, in, with, of
- Tail words take major portion of vocabulary, but they rarely occur in documents
 - E.g., dextrosinistral
- The rest is more representative
 - To be included in the controlled vocabulary

Text Mining Tasks:

Step 1 – Tokenization	Step 2 – Stemming /normalization	Step 3 – controlled vocabulary filtering BOW	Step 4 – Compute the term frequencies Term	Step 5 – Normalize the term frequencies Term	Step 6 – Create an indexing file Term	Step 7 – Create the vector space model Term	Step 8 – Compute the IDF Term	Step 9 – Compute TF-IDF Term	Step 10 – Inverted List Term
		Step 3 – construct N-gram N-gram	Step 4 – controlled vocabulary filtering N-gram						Step 11 – Convert Query Term
									Step 12 – Sim(Q,D0)→Sim(Q,Dn) Term
									Step 13 – return doc with highest Sim Term

Query A Document

- ❖ similar to the previous steps:
 - 1. Tokenization/Stemming.
 - 2. Controlled vocabulary filtering.
 - 3. Transform the query string into a vector space model (VSM) by using TD-IDF schema.
 - 4. Normalize the VSM into unit length.
- ❖ Improving Query efficiency:
 - Inverted list
 - Query Q="interesting course"
 - Only calculate the doc exist in the invert list; By Sim(Q,Dn)

Vector Space Model Advantages Vs. Disadvantages

Advantages	Disadvantages
<ul style="list-style-type: none">• Empirically effective!• Intuitive• Easy to implement• Well-studied/mostly evaluate• The Smart system<ul style="list-style-type: none">◦ Developed at Cornell: 1960-1999◦ Still widely used• Warning: many variants of TF-IDF!	<ul style="list-style-type: none">• Assume term independence• Lack of “predictive adequacy”<ul style="list-style-type: none">◦ Arbitrary term weighting◦ Arbitrary similarity measure• Lots of parameter tuning!