

Information Retrieval

Venkatesh Vinayakarao

Term: Aug – Sep, 2019
Chennai Mathematical Institute



Computers understand very little of the meaning of human language. This profoundly limits our ability to give instructions to computers, the ability of computers to explain their actions to us, and the ability of computers to analyse and process text. Vector space models (VSMs) of semantics are beginning to address these limits

– **Turney and Pantel, JAIR 2010.**

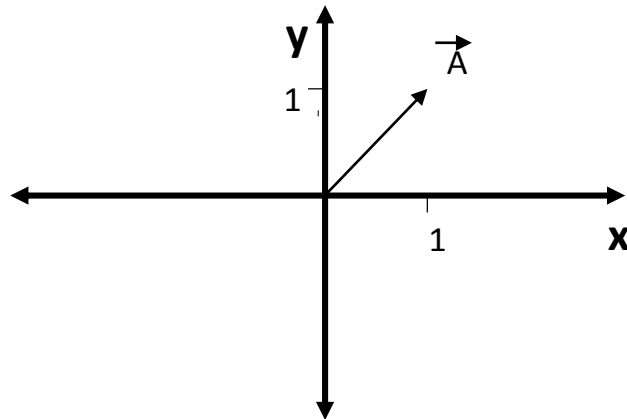


A Better Approach

**Revisiting
Linear Algebra**

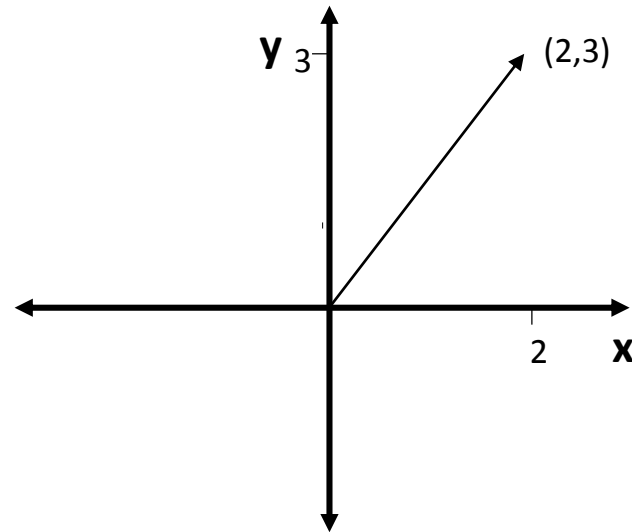
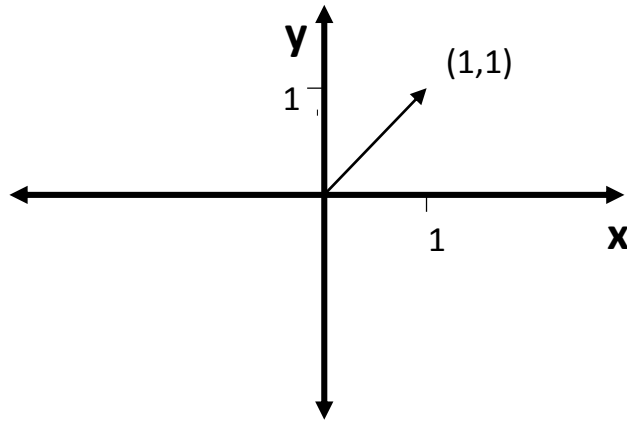
Vectors

- Geometric entity which has magnitude and direction

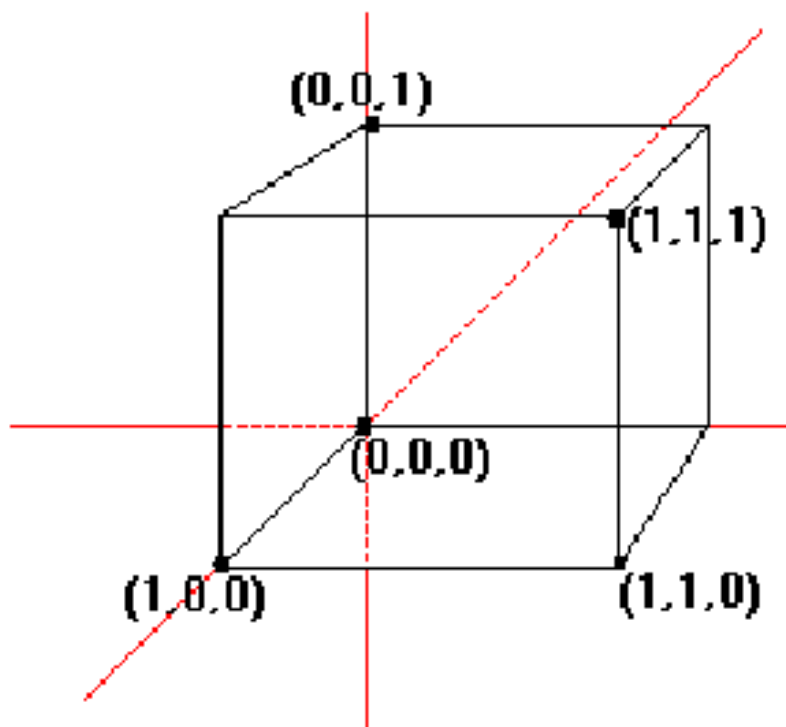


- If (x,y) is our vector of interest, this figure shows \vec{A} vector = $(1,1)$.

How is $(2,3)$ Different?



What is $(1,1,1)$?

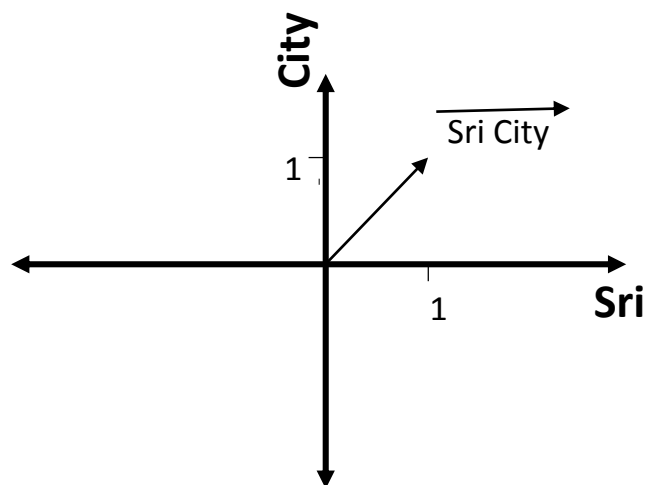


Remember!

**A number is just a mathematical object. We
give meaning to it!**

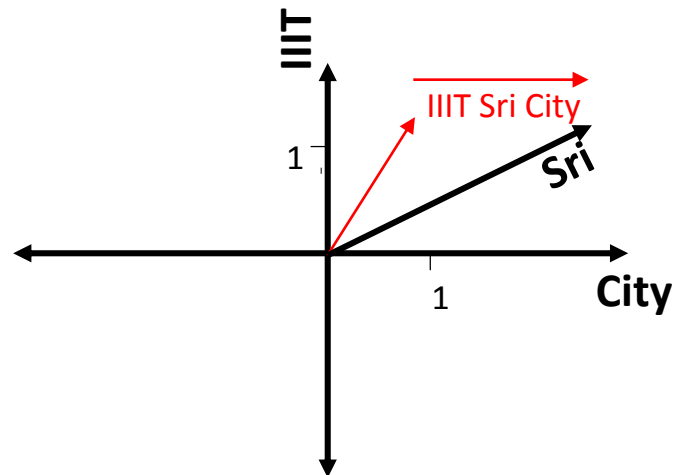
Sentences are Vectors

- “Sri City” as a vector



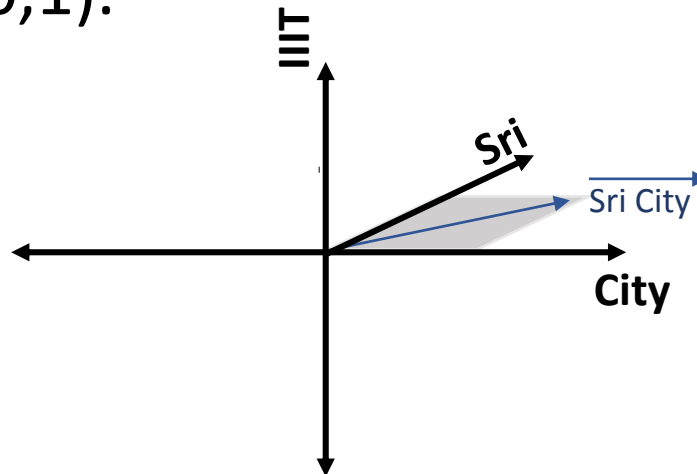
Sentences are Vectors

- “IIT Sri City” is a 3-dimensional vector



Sentences are Vectors

- On this 3D space, “Sri City” vector will lie on the x (City) and z (Sri) plane. If (x,y,z) denotes the vector, “Sri City” is $(1,0,1)$.



More Linear Algebra...

- So, we learned to represent English phrases on the vector space.
- We need something more!

Revisiting Matrices

Natural Language Phrases as Vectors

Let query $q = \text{"IIIT Sri City"}$.

Let document, $d_1 = \text{"IIIT Sri City"}$ and $d_2 = \text{"IIIT Delhi"}$.

	IIIT	Sri	City	Delhi
q	1	1	1	0
d_1	1	1	1	0
d_2	1	0	0	1

$q = (1,1,1,0)$, $d_1 = (1,1,1,0)$ and $d_2 = (1,0,0,1)$

Quiz

- Considering the following vectors:

	IIIT	Sri	City	Delhi
q	1	1	1	0
d ₁	1	1	1	0
d ₂	1	0	0	1

- What is the Natural Language (NL) equivalent of $(0,1,1,0)$?
- What is the NL equivalent of $(1,0,0,1)$?
- What is the vector for Delhi?
- What is the NL equivalent of q?

Which of the Following are Sets?

- ~~{1, 2, 3, 4, 5, 6, 5, 7, 8, 9, 10, 11, 12, 13}~~
- ~~{A, B, C, D, E, F, G, H, I, I, J, K, L, M, N, O}~~
- ~~{apple, banana, orange, apple, banana, orange}~~

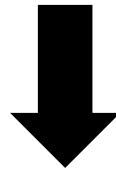


Bag

- {1, 2, 3, 4, 5, 6, 5, 7, 8, 9, 10, 11, 12, 13}
- {A, B, C, D, E, F, G, H, I, I, J, K, L, M, N, O}
- {apple, banana, orange, apple, banana, orange}

Set of Words Representation

- “IIIT Sri City” \rightarrow {IIIT, Sri, City}
- “IIIT Sri City, Sri City” \rightarrow {IIIT, Sri, City}



	IIIT	Sri	City
q	1	1	1

Leads to same term-document matrix

Set Similarity

- Similarity between two sets is easy

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

Jaccard Similarity

Quiz

- What is the Jaccard similarity between
 - $\{1,2,3\}$ and $\{4,5,6\}$?
 - $\{1,2,3\}$ and $\{1,2,4\}$?
 - $\{1,2,3\}$ and $\{1,2,3\}$?

Quiz

- What is the Jaccard similarity between
 - $\{1,2,3\}$ and $\{4,5,6\} = 0$
 - $\{1,2,3\}$ and $\{1,2,4\} = \frac{|\{1,2\}|}{|\{1,2,3,4\}|} = 0.5$
 - $\{1,2,3\}$ and $\{1,2,3\} = 1$

Quiz

- What is the Jaccard similarity between
 - “IIT is Great” and “IITD is Great”?

Quiz

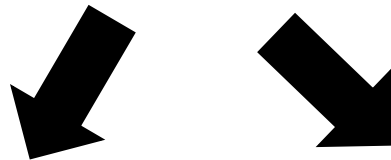
- What is the Jaccard similarity between
 - “IIT is Great” and “IITD is Great”?
- Same as Jaccard Similarity between {IIT, is, Great} and {IITD, is, Great}. Equals 0.5.

Ranked Retrieval

In the **Ranked Retrieval** model, we may want to model documents as **bag of words**.

Bag of Words Representation

- “IIT Sri City” → {IIT, Sri, City}
- “IIT Sri City, Sri City” → [IIT, Sri, Sri, City, City]



	IIT Sri City				IIT Sri City, Sri City		
	IIT	Sri	City		IIT	Sri	City
q	1	1	1	q	1	2	2

Leads to different term-document matrix

Set of Words → Bag of Words

Information Retrieval

Venkatesh Vinayakarao

Term: Aug – Dec, 2018

Indian Institute of Information Technology, Sri City



Searchers may not have a well-developed idea of what information they are searching for, they may not be able to express their conceptual idea of what information they want into a suitable query and they may not have a good idea of what information is available for retrieval. – **Ruthven and Lalmas, The Knowledge Engineering Review, 2003.**



Relevance Feedback

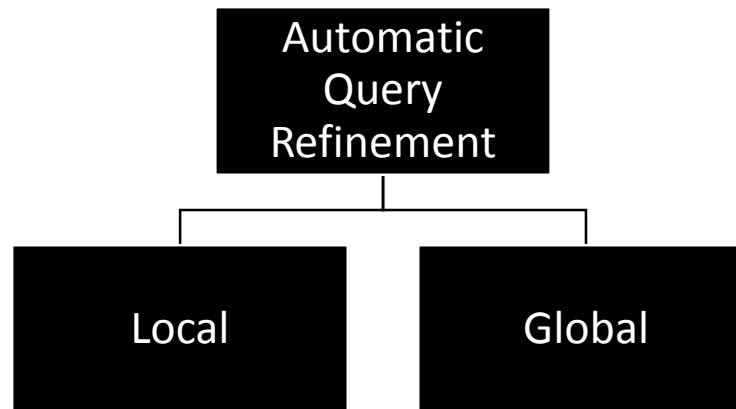
How to improve relevance?

Relevance Feedback
And
Query Expansion

The Problem of Synonymy

- What result do you expect for a query, “plane”?
- What if plane appears in this query, “plane from Delhi to Goa”?
- So many synonyms which will work for web search...
 - Flight
 - Aircraft
 - Airplane
 - Aeroplane
 - By Air
 - Fly
 - Flgt
 - Arcrft

How to ensure good results?



Use the **query** or the **results** for reformulating the query

We will study:

Relevance Feedback

Pseudorelevance

Indirect Relevance Feedback

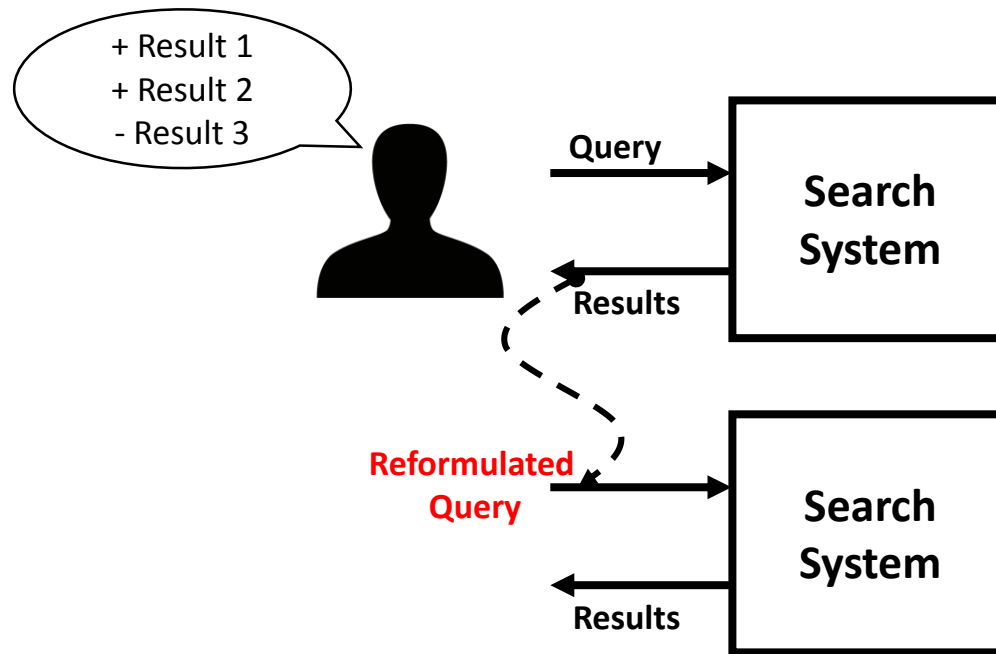
Do not use the **query** or the **results** for reformulating the query.

Eg:

Use Thesaurus.

Do Spelling Correction.

Relevance Feedback



An Example

The image shows a screenshot of the RefMED search engine interface. At the top left is the logo of the University of Hohenheim. The main title is "RefMED: Relevance Feedback Search Engine for PubMed". Below this is a search bar with the text "mrsa" and a "Go" button. To the right of the search bar is a "Push Feedback" button, which is highlighted with a red arrow and a circled "(2)". Below the search bar are several filters: "Display: Summary", "Show: 20", "PMID", "Year: ~", and "Feedback: 3". There are also "Not relevant" and "Relevant" buttons. Below the filters, it says "Items 1 - 20 of 17,656. (0.02662 seconds)". The search results are listed below, with three items shown. Each item has a title, authors, journal information, and PMID. To the right of each item are three radio buttons for feedback. The first item has the third radio button selected, highlighted with a red arrow and a circled "(1)".

RefMED: Relevance Feedback Search Engine for PubMed

Search PubMed for

Display: Show: Year: Feedback:

Items 1 - 20 of 17,656. (0.02662 seconds)

1: [Methicillin-Resistant Staphylococcus aureus ST9 in Pigs in Thailand.](#)

Skov Robert L , Hinjoy Soawapak , Imanishi Maho , Larsen Jesper , Larsen Anders R , Nelson Kenra E , Davis Meghan F , Duangsong Kwanjit , Tharavichitkul Prasit
PloS one. 2012-00-00;7(2):e31245
PMID: 22363594

2: [Inhibition of Virulence Gene Expression in Staphylococcus aureus by Novel Depsipeptides from a Marine Photobacterium.](#)

Larsen Thomas O , Gram Lone , Ingmer Hanne , Wietz Matthias , Gotfredsen Charlotte H , Kj??rulf Louise , Nielsen Anita , Mansson Maria
Marine drugs. 2011-12-00;9(12):2537-52
PMID: 22363239

3: [The Prevalence, Genotype and Antimicrobial Susceptibility of High- and Low-Level Mupirocin Resistant Methicillin-Resistant Staphylococcus aureus.](#)

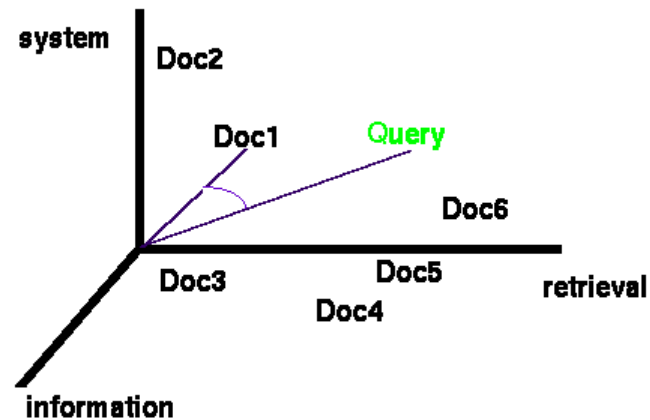
Park Se Young , Kim Shin Moo , Park Seok Don

Image source: <https://sites.google.com/site/postechdm/research>

Interesting Characteristics

- Indexed content is unknown to the user.
- “Information Need” changes after looking at the results.
 - User visits youtube to listen to a specific set of songs.
 - After the first song, he changes his mind and listens to something else!

A Recap of Vector Space Models



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

Image Source: <https://fox.cs.vt.edu/talks/1995/KY95/>

Rocchio Algorithm for Relevance Feedback

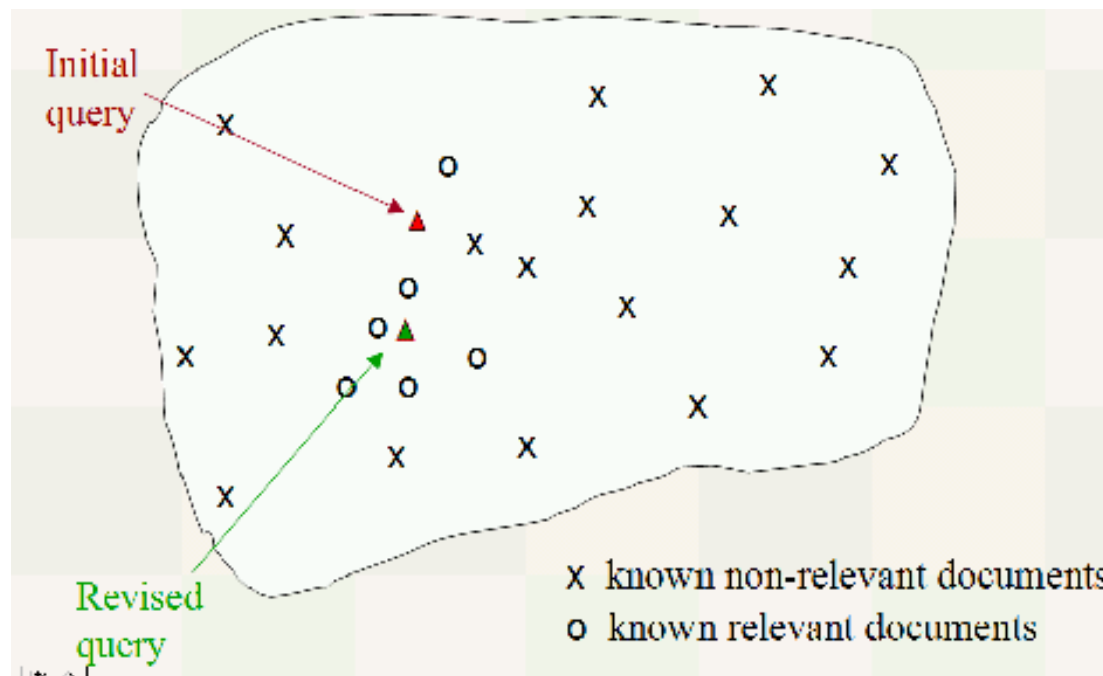


Image Source: <https://nlp.stanford.edu/IR-book/>

Moving the Centroid!

Modify the query (and therefore, the query vector from q_0 to q_m):

$$\vec{q}_m = \alpha \vec{q}_0 + \beta \frac{1}{|D_r|} \sum_{\vec{d}_j \in D_r} \vec{d}_j - \gamma \frac{1}{|D_{nr}|} \sum_{\vec{d}_j \in D_{nr}} \vec{d}_j$$

D_r = Set of known relevant documents

D_{nr} = Set of known nonrelevant documents

q_0 = Initial query vector

q_m = Modified query vector

Rocchio relevance feedback - Example

- Given:
 - Initial query = “cheap CDs cheap DVDs extremely cheap CDs”.
 - $d_1 =$ “CDs cheap software cheap CDs” is judged as relevant.
 - $d_2 =$ “cheap thrills DVDs” is judged as nonrelevant
- What would the revised query vector be after relevance feedback?

Let us solve this together

Assume that we are using direct term frequency (with no scaling and no document frequency). There is no need to length-normalize vectors. Assume $\alpha = 1$, $\beta = 0.75$, $\gamma = 0.25$.

Representing Initial Query in Vector Space

Initial query = “cheap CDs cheap DVDs extremely cheap CDs”.

	cheap	CDs	DVDs	extremely	software	thrills
q_0	3	2	1	1	0	0

Rocchio relevance feedback - Example

Quiz: Can you complete the following table?

q_0 = “cheap CDs cheap DVDs extremely cheap CDs”.

d_1 = “CDs cheap software cheap CDs”.

d_2 = “cheap thrills DVDs”.

	cheap	CDs	DVDs	extremely	software	thrills
q_0	3	2	1	1	0	0
d_1						
d_2						

Rocchio relevance feedback - Example

Quiz: Can you complete the following table?

q_0 = “cheap CDs cheap DVDs extremely cheap CDs”.

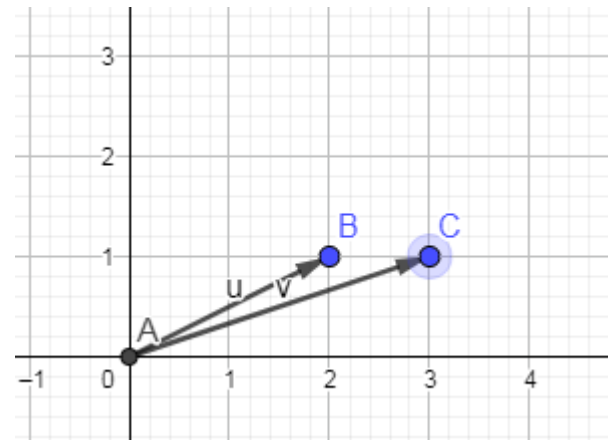
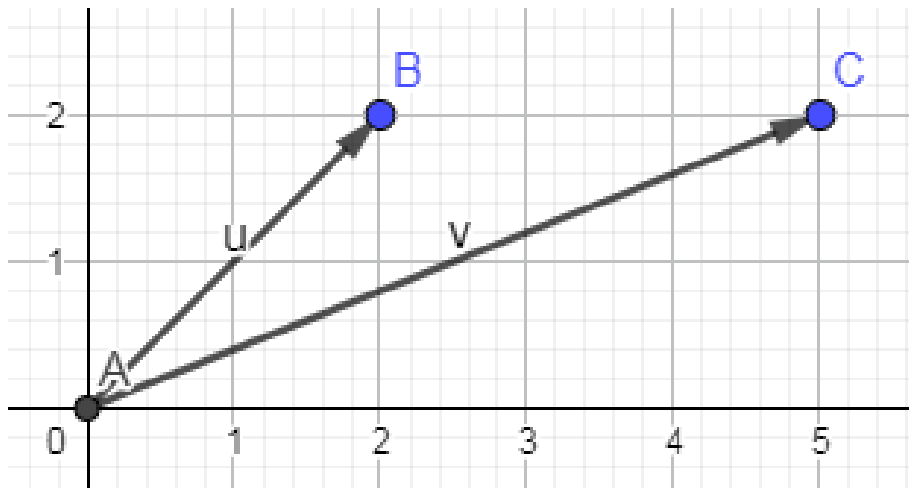
d_1 = “CDs cheap software cheap CDs”.

d_2 = “cheap thrills DVDs”.

	cheap	CDs	DVDs	extremely	software	thrills
q_0	3	2	1	1	0	0
d_1	2	2	0	0	1	0
d_2	1	0	1	0	0	1

Moving Vectors

- Move $(2,2)$ to $(5,2)$ by adding 3 to x .



Rocchio relevance feedback - Example

Quiz: How to calculate the modified query vector, q_m ?

d_1 is judged as **relevant**. d_2 is judged as **non-relevant**.

Assume $\alpha = 1$, $\beta = 0.75$, $\gamma = 0.25$.

	cheap	CDs	DVDs	extremely	software	thrills
q_0	3	2	1	1	0	0
d_1	2	2	0	0	1	0
d_2	1	0	1	0	0	1
q_m						

$$\vec{q}_m = \alpha \vec{q}_0 + \beta \frac{1}{|D_r|} \sum_{\vec{d}_j \in D_r} \vec{d}_j - \gamma \frac{1}{|D_{nr}|} \sum_{\vec{d}_j \in D_{nr}} \vec{d}_j$$

Rocchio relevance feedback - Example

Quiz: How to calculate the modified query vector, q_m ?

d_1 is judged as relevant. d_2 is judged as nonrelevant.

Assume $\alpha = 1$, $\beta = 0.75$, $\gamma = 0.25$.

	cheap	CDs	DVDs	extremely	software	thrills
q_0	3	2	1	1	0	0
d_1	2	2	0	0	1	0
d_2	1	0	1	0	0	1

$$q_m = q_0 + 0.75 * d_1 - 0.25 * d_2$$

q_m	4.25	3.5	0.75	1	0.75	0
-------	------	-----	------	---	------	---

Negative weight does not make sense. So, leave them as zero.

$$\vec{q}_m = \alpha \vec{q}_0 + \beta \frac{1}{|D_r|} \sum_{\vec{d}_j \in D_r} \vec{d}_j - \gamma \frac{1}{|D_{nr}|} \sum_{\vec{d}_j \in D_{nr}} \vec{d}_j$$

Pseudo (Blind) Relevance Feedback

- No User Judgment.
- Assume that the top-k ranked documents are relevant.

Initial query = “cheap CDs cheap DVDs extremely cheap CDs”.

d_1 = “CDs cheap software cheap CDs”.

d_2 = “cheap thrills DVDs”.

What would the revised query vector be **after pseudo relevance feedback if top-1 document is considered as relevant?**

Assume that we are using direct term frequency (with no scaling and no document frequency). There is no need to length-normalize vectors. Assume $\alpha = 1$, $\beta = 0.75$, $\gamma = 0.25$.

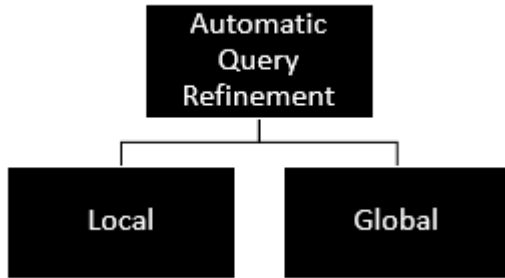
- May lead to query drift.

Indirect (Implicit) Relevance Feedback

- No asking for judgments from users.
- No automatic feedback such as assuming top-k documents as relevant.

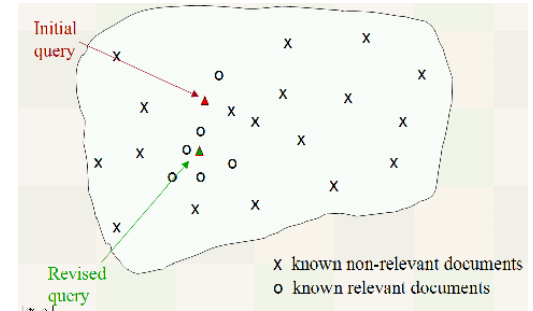
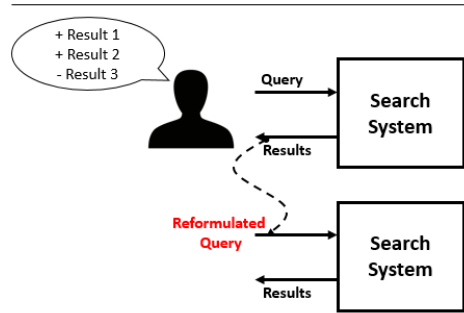
Clickstream Mining

Recap



Dealing with Local Query Refinement

Relevance Feedback



	cheap	CDs	DVDs	extremely	software	thrills
q_0	3	2	1	1	0	0
d_1	2	2	0	0	1	0
d_2	1	0	1	0	0	1
$q_m = q_0 + 0.75 * d_1 - 0.25 * d_2$						
q_m	4.25	3.5	0.75	1	0.75	0

Negative weight does not make sense. So, leave them as zero.

$$\bar{q}_m = \alpha \bar{q}_0 + \beta \frac{1}{|D_r|} \sum_{\bar{d}_j \in D_r} \bar{d}_j - \gamma \frac{1}{|D_{nr}|} \sum_{\bar{d}_j \in D_{nr}} \bar{d}_j$$

Pseudo (Blind) Relevance Feedback
 Assume top-k is relevant.

Indirect (Implicit) Relevance Feedback
 Using clickstreams, query logs, etc.

Global (User/Result-Independent) Query Refinement

- Automatic Thesaurus Generation
 - Fast = rapid
 - Tall = height?
 - Sound = noise?
 - Restaurant = Hotel = Motel?
- How to handle domain specific phrases?
- Slangs!
- ...

How to automate the thesaurus generation?

Co-occurrence Analysis

MAINFRAMES

Mainframes **are primarily** referred to large computers with **rapid**, advanced processing capabilities that **can execute and** perform tasks **equivalent to many** Personal Computers (PCs) machines **networked together**. It is **characterized with high quantity** Random Access Memory (RAM), very large secondary storage devices, and **high-speed** processors to cater for the needs of the computers under its service.

Consisting of advanced components, mainframes have the capability of

MAINFRAMES

Mainframes **usually are** referred those computers with **fast**, advanced processing capabilities that **could perform by itself** tasks **that may require a lot of** Personal Computers (PC) Machines. **Usually mainframes would have lots of** RAMs, very large secondary storage devices, and **very fast** processors to cater for the needs of those computers under its service.

Due to the advanced components mainframes have, **these computers** have the capability of running multiple

Co-occurrence Analysis

- Term-Document Matrix
 - How often does individual terms appear in a document?
- Term-Term Matrix
 - How often terms co-occur?

Quiz: Which books are similar?

	Book1	Book2	Book3	Book4
cricket	400	10	355	3
football	5	5	4	4
hockey	9	330	10	200
tennis	2	6	12	4

Co-occurrence Analysis

- Two documents are similar if the document vectors are similar.

	Book1	Book2	Book3	Book4
cricket	400	10	355	3
football	5	5	4	4
hockey	9	330	10	200
tennis	2	6	12	4

Book1 and Book3 seem to be on cricket.
Book2 and Book4 are about hockey.

Co-occurrence Analysis

- Two terms are similar if the term vectors are similar.

	Book1	Book2	Book3	Book4
boundary	400	310	355	389
four	515	225	390	400
movie	9	4	8	1
film	2	6	9	2

Remember, context is important!

Magic with Matrices

Transpose

- If A is as given below, what is A^T ?

$$A = \begin{pmatrix} 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \end{pmatrix}$$

$$A^T = \begin{pmatrix} 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \end{pmatrix}$$

Co-occurrence Analysis

- Assume a Boolean term-document matrix A .
- What does AA^T mean?

$$\begin{array}{c} t_1 \\ t_2 \\ t_3 \\ t_4 \end{array} \begin{array}{c} d_1 \ d_2 \ d_3 \ d_4 \\ \left(\begin{array}{cccc} 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \end{array} \right) \end{array} \begin{array}{c} d_1 \ d_2 \ d_3 \ d_4 \\ \left(\begin{array}{cccc} 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \end{array} \right) \end{array} = \begin{array}{c} t_1 \ t_2 \ t_3 \ t_4 \\ \left(\begin{array}{cccc} 2 & 2 & 0 & 0 \\ 2 & 2 & 0 & 0 \\ 0 & 0 & 2 & 2 \\ 0 & 0 & 2 & 2 \end{array} \right) \end{array}$$

- Usually, weighted length-normalized tf in a sliding window is used to count co-occurrence.

Thank You