

LC029 정보검색

2022 11 10

Chapter 9 : Relevance Feedback and Query Expansion

For High Recall

- The same concept may be represented with different words.
 - Searching for *aircraft* doesn't match with *plane*
 - Searching for *thermodynamic* doesn't match with *heat*
- Solutions for High Recall
 - Manual refinement of a query
 - Automatic refinement of a query
- We will discuss the automatic refinement methods.

Automatic Refinement of Queries

- Global methods (static)
 - Expanding query terms with semantically similar terms.
 - Query expansion with a **thesaurus** or **WordNet**
 - Query expansion with **automatically generated thesaurus**
- Local methods (dynamic) 가 가
 - Refining a query **relative to the documents** that initially appear to match the query.
 - **Relevance feedback**
 - **Pseudo relevance feedback**
 - **Indirect relevance feedback**

Thesaurus

th | ≡

- From Wikipedia,
 - A thesaurus or synonym dictionary is a reference work for finding synonyms and sometimes antonyms of words.
 - Often used by writers to help find the best word to express an idea:

... to find the word or words, by which an idea may be most fitly and aptly expressed (by Peter Mark **Roget**, 1852)

The screenshot shows the Thesaurus.com website. The header is orange with the Thesaurus.com logo on the left. In the center, there is a search bar with the word 'doctor' entered. To the left of the search bar, the word 'SYNONYMS' is followed by a downward arrow. To the right of the search bar is a magnifying glass icon. Below the header, the text 'SYNONYMS FOR doctor' is displayed. Underneath this, there are four columns of synonyms, each word in an orange box. The synonyms are: expert, specialist, doc, medico, physician, surgeon, healer, quack, professor, MD, intern, general practitioner, scientist, bones, medic, and medical person. At the bottom right, there is a legend: an orange square followed by the text 'MOST RELEVANT'.

Thesaurus.com

SYNONYMS ▾ | doctor 🔍

SYNONYMS FOR doctor

expert	specialist	doc	medico
physician	surgeon	healer	quack
professor	MD	intern	general practitioner
scientist	bones	medic	medical person

■ MOST RELEVANT

Thesaurus



grand, *adj.* large, impressive, magnificent, stately, majestic; pretentious, ostentatious; elegant, lofty.

graph, *n.* diagram, chart, plot, PLAN; bar, circle, *etc.* graph.

graphic, *adj.* pictorial, descriptive; vivid, diagrammatic; delineative; picturesque. See REPRESENTATION,

graceless

attractive. See BEAUTY, ELEGANCE.
graceless, *adj.* ungracious, tactless; inept, awkward, clumsy, ungainly; inelegant; sinful, corrupt. See UGLINESS, IRRELIGION, IMPENITENCE, INELEGANCE.

gracious, *adj.* gentle, courteous, tactful; kind, thoughtful, benign; affable, obliging; generous; charming, attractive. See BENEVOLENCE, COURTESY.

graduation, *n.* stage, DEGREE.

grade, *n.* level, quality, class; rank; standing; gradation, slope, tilt, slant. See OBLIQUITY, DEGREE.

gradual, *adj.* slow, progressive,

[235]

grapple

moderate, leisurely. See DEGREE, SLOWNESS.

graduate, *n.* measure, beaker. *Slang*, grad. See DEGREE. —*v.* measure, classify, grade; calibrate; pass, commence. See ARRANGEMENT, MEASUREMENT, SCHOOL.

graffiti, *n. pl.* See WRITING.

graft, *v.* inoculate, bud, transplant, implant, join, crossbreed. See MIXTURE, INSERTION. —*n.* corruption, porkbarrel politics. See STEALING, IMPROBITY.

grain, *n.* fruit, cereal, seed, grist, kernel; TEXTURE, temper, TENDENCY; mite, speck, bit (see LITTLENESS).

GRAMMAR

Nouns—1, grammar; accident, syntax, analysis, synopsis, praxis, punctuation, syllab[if]ication; parts of speech; participle; article, noun, substantive, pronoun, verb, adjective, adverb, preposition, postposition, interjection, conjunction; particle; prefix, suffix, combining form, element; inflection, inflexion, case, declension, conjugation.

2, tense: present, past, preterit, future; imperfect, perfect, past perfect, pluperfect; progressive, *etc.*

3, mood, mode: infinitive, indicative, subjunctive, imperative.

4, sentence, paragraph, clause, phrase.

5, style; philology, language; phraseology, wording, rhetoric, diction. See SPEECH; WRITING.

Verbs—parse, analyze, diagram; punctuate; conjugate, decline, inflect.

Antonyms, see ERROR.

grand, *adj.* large, impressive, magnificent, stately, majestic; pretentious, ostentatious; elegant, lofty. See GREATNESS, OSTENTATION, REPUTE.

grandeur, *n.* GREATNESS, magnificence; show, ostentation; splendor, majesty; eminence, stateliness, loftiness. See IMPORTANCE.

grandfather, *n.* grandpa, grandsire; gaffer, old man. See ANCESTRY, AGE.

grandiose, *adj.* grand; stately, pompous, bombastic. See BOASTING, OSTENTATION.

grandmother, *n.* granny, grandma, nana, nanny; old woman. See ANCESTRY, AGE.

grandstand, *n.* bleachers. See VISION.

grant, *n.* gift, allotment, contribution.

—*v.* bestow, give, yield; concede; permit; contribute. See GIVING, PERMISSION, CONSENT, RIGHTNESS, TRANSFER.

granular, *adj.* granulated, grainy, mealy, gritty, sandy. See POWDERINESS.

grapevine, *n., colloq.*, rumor; rumor mill, pipeline, a little bird. See INFORMATION.

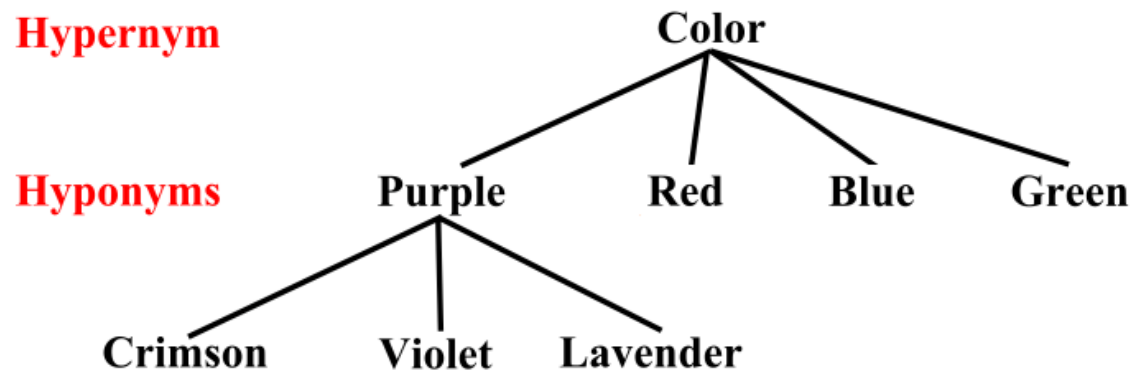
graph, *n.* diagram, chart, plot, PLAN; bar, circle, *etc.* graph.

graphic, *adj.* pictorial, descriptive; vivid, diagrammatic; delineative; picturesque. See REPRESENTATION, DESCRIPTION.

grapple, *v.* seize, grasp, clutch, struggle, contend. See OPPOSITION, CONTENTION.

WordNet

- From Wikipedia,
 - A lexical database of semantic relations between words
 - WordNet links words into semantic relations including synonyms, hyponyms, meronyms and holonyms
 - It is accessible to human users via a web browser, but most importantly, it is **machine readable** and its primary use is in NLP (Natural Language Processing) applications



WordNet

WordNet Search - 3.1

- [WordNet home page](#) - [Glossary](#) - [Help](#)

Word to search for:

Display Options:

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations

Display options for sense: (gloss) "an example sentence"

Noun

- **S: (n)** **computer**, [computing machine](#), [computing device](#), [data processor](#), [electronic computer](#), [information processing system](#) (a machine for performing calculations automatically)
 - [direct hyponym](#) / [full hyponym](#)
 - [part meronym](#)
 - [domain category](#)
 - [domain term category](#)
 - [direct hypernym](#) / [inherited hypernym](#) / [sister term](#)
 - [part holonym](#)
 - [derivationally related form](#)
- **S: (n)** [calculator](#), [reckoner](#), [figurer](#), [estimator](#), **computer** (an expert at calculation (or at operating calculating machines))

WordNet

- **S: (n) computer**, [computing machine](#), [computing device](#), [data processor](#), [electronic computer](#), [information processing system](#) (a machine for performing calculations automatically)
 - **direct hyponym** / **full hyponym**
 - **S: (n) analog computer**, [analogue computer](#) (a computer that represents information by variable quantities (e.g., positions or voltages))
 - **S: (n) digital computer** (a computer that represents information by numerical (binary) digits)
 - **S: (n) home computer** (a computer intended for use in the home)
 - **S: (n) node**, [client](#), [guest](#) ((computer science) any computer that is hooked up to a computer network)
 - **S: (n) number cruncher** (a computer capable of performing a large number of mathematical operations per second)
 - **S: (n) pari-mutuel machine**, [totalizer](#), [totaliser](#), [totalizator](#), [totalisator](#) (computer that registers bets and divides the total amount bet among those who won)
 - **S: (n) predictor** (a computer for controlling antiaircraft fire that computes the position of an aircraft at the instant of a shell's arrival)
 - **S: (n) server**, [host](#) ((computer science) a computer that provides client stations with access to files and printers as shared resources to a computer network)
 - **S: (n) Turing machine** (a hypothetical computer with an infinitely long memory tape)
 - **S: (n) web site**, [website](#), [internet site](#), [site](#) (a computer connected to the internet that maintains a series of web pages on the World Wide Web)
"the Israeli web site was damaged by hostile hackers"

WordNet

- **S: (n) computer**, [computing machine](#), [computing device](#), [data processor](#), [electronic computer](#), [information processing system](#) (a machine for performing calculations automatically)
 - [direct hyponym](#) / [full hyponym](#)
 - **part meronym**
 - **S: (n) busbar**, [bus](#) (an electrical conductor that makes a common connection between several circuits) *"the busbar in this computer can transmit data either way between any two components of the system"*
 - **S: (n) cathode-ray tube**, [CRT](#) (a vacuum tube in which a hot cathode emits a beam of electrons that pass through a high voltage anode and are focused or deflected before hitting a phosphorescent screen)
 - **S: (n) central processing unit**, [CPU](#), [C.P.U.](#), [central processor](#), [processor](#), [mainframe](#) ((computer science) the part of a computer (a microprocessor chip) that does most of the data processing) *"the CPU and the memory form the central part of a computer to which the peripherals are attached"*
 - **S: (n) chip**, [microchip](#), [micro chip](#), [silicon chip](#), [microprocessor chip](#) (electronic equipment consisting of a small crystal of a silicon semiconductor fabricated to carry out a number of electronic functions in an integrated circuit)
 - **S: (n) computer accessory** (an accessory for a computer) *"when you add in all the computer accessories you are going to need the computer gets pretty expensive"*
 - **S: (n) computer circuit** (a circuit that is part of a computer)
 - **S: (n) data converter** (converter for changing information from one code to another)
 - **S: (n) disk cache** (a cache that stores copies of frequently used disk sectors in random access memory (RAM) so they can be read without accessing the slower disk)

Relevance Feedback

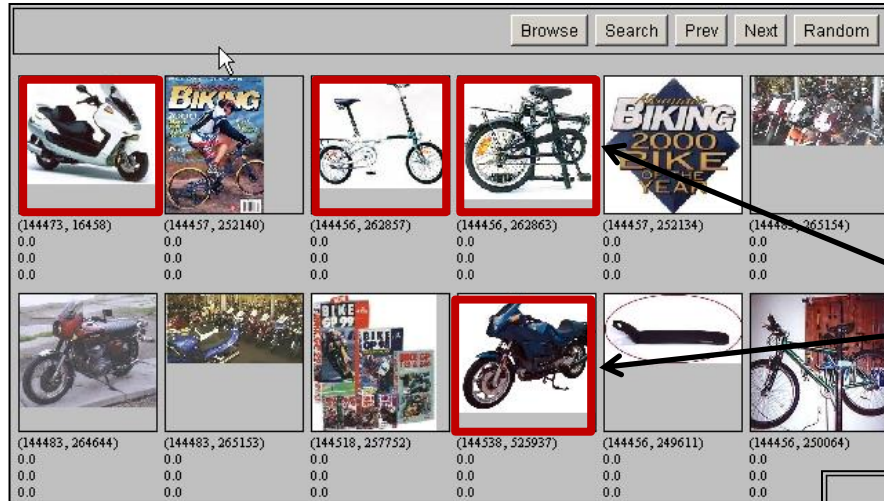
- Idea
 - It may be difficult to formulate a good query when you don't know the collection well.
 - Image search is a good example:
Difficult to formulate a query, but easy to mark retrieved images relevant or nonrelevant.
- Users involved in IR process want to improve the final search results:
 - Users feedback on relevance of docs in initial set of results.
 - IR system refines the initial query to present better results.

refine

Relevance Feedback

- Basic Procedure of Relevance Feedback
 - User issues a (short, simple) query.
 - The system returns an initial set of retrieved results.
 - The user marks some results as ***relevant*** or ***non-relevant***.
 - The system computes a better representation of the information need based on feedback.
 - The system displays a **revised set of retrieved results**.
 - The procedure can go through one or more iterations.
 - ***ad hoc retrieval*** is regular retrieval without relevance feedback.
- | | | |
|----|---|---|
| 1. | 가 | |
| 2. | | |
| 3. | 가 | 가 |
| 4. | | |
| 5. | | |

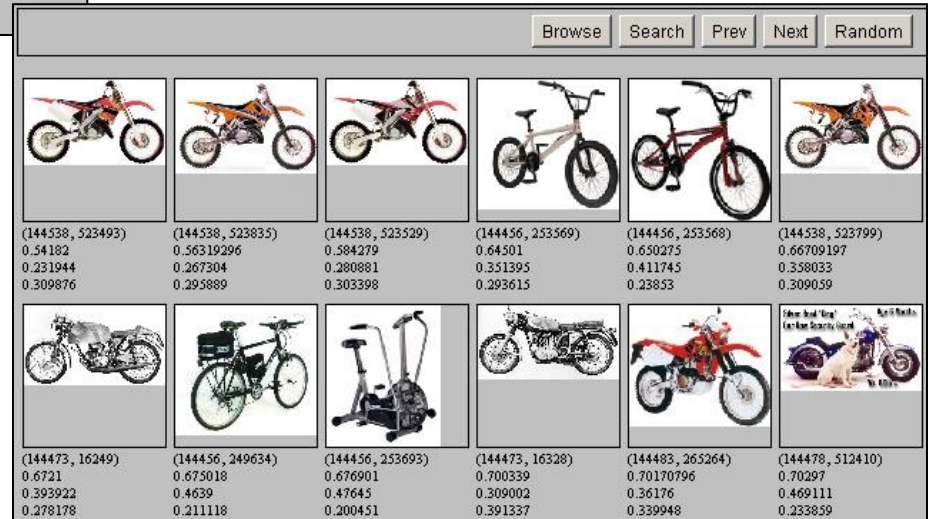
Relevance Feedback : Example 1



← initial search results

user feedback

improved search results →



Relevance Feedback : Example 2

Query: New space satellite applications

**retrieved
documents**

**relevance
feedback**

1. 0.539, 08/13/91, NASA Hasn't Scrapped Imaging Spectrometer
2. 0.533, 07/09/91, NASA Scratches Environment Gear From Satellite Plan
3. 0.528, 04/04/90, Science Panel Backs NASA Satellite Plan, But Urges Launches of Smaller Probes
4. 0.526, 09/09/91, A NASA Satellite Project Accomplishes Incredible Feat: Staying Within Budget

**Improved
results**

1. 0.513, 07/09/91, NASA Scratches Environment Gear From Satellite Plan
2. 0.500, 08/13/91, NASA Hasn't Scrapped Imaging Spectrometer
3. 0.493, 08/07/89, When the Pentagon Launches a Secret Satellite, Space Sleuths Do Some Spy Work of Their Own
4. 0.493, 07/31/89, NASA Uses 'Warm' Superconductors For Fast Circuit

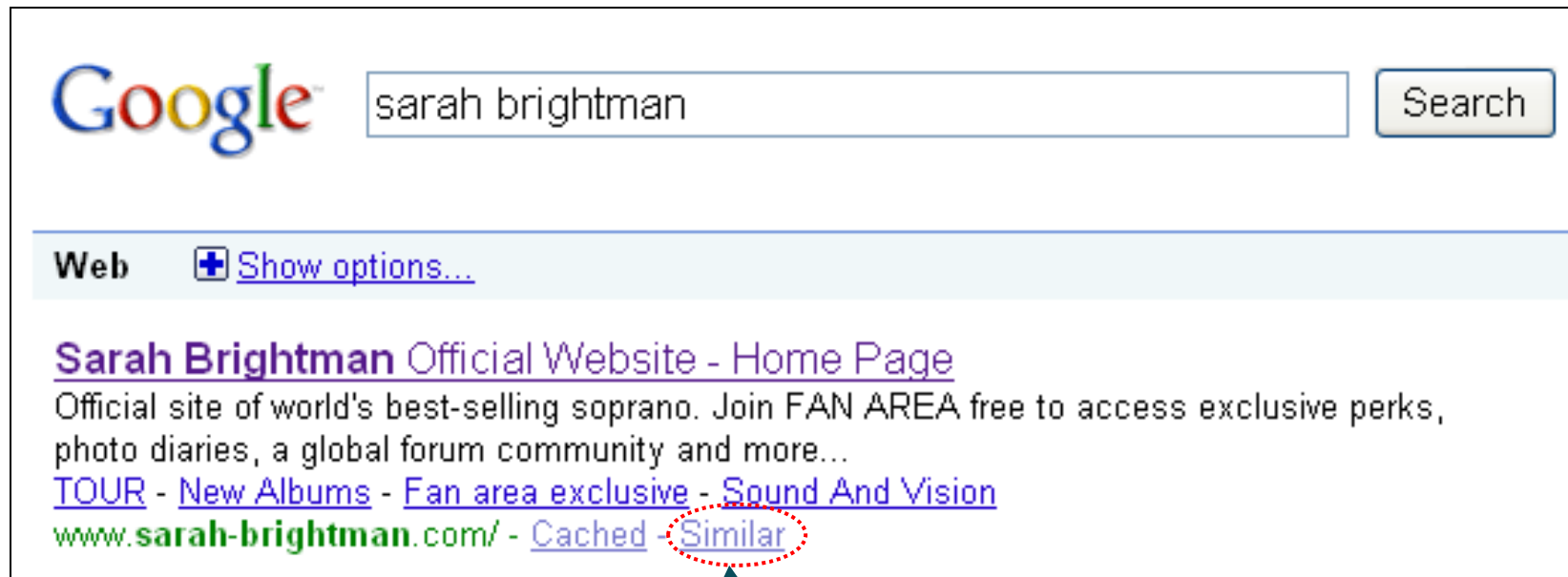
Relevance Feedback : Example 2

- Expanded query after relevance feedback

2.074	new	15.12	space
30.82	satellite	5.660	application
5.991	nasa	5.196	eos
4.196	launch	3.972	aster
3.516	instrument	3.446	arianespace
3.004	bundespost	2.806	ss
2.790	rocket	2.053	scientist
2.003	broadcast	1.172	earth
0.836	oil	0.646	measure

Query: New space satellite applications

Relevance Feedback : Example 3



relevance
feedback

Local Methods

1. Key Concepts
2. Relevance feedback
3. Pseudo relevance feedback
4. Indirect relevance feedback
5. Evaluation of Relevance Feedback

Key Concept : Centroid

Key Concept : Centroid

- Definition

$$\vec{\mu}(C) = \frac{1}{|C|} \sum_{d \in C} \vec{d}$$

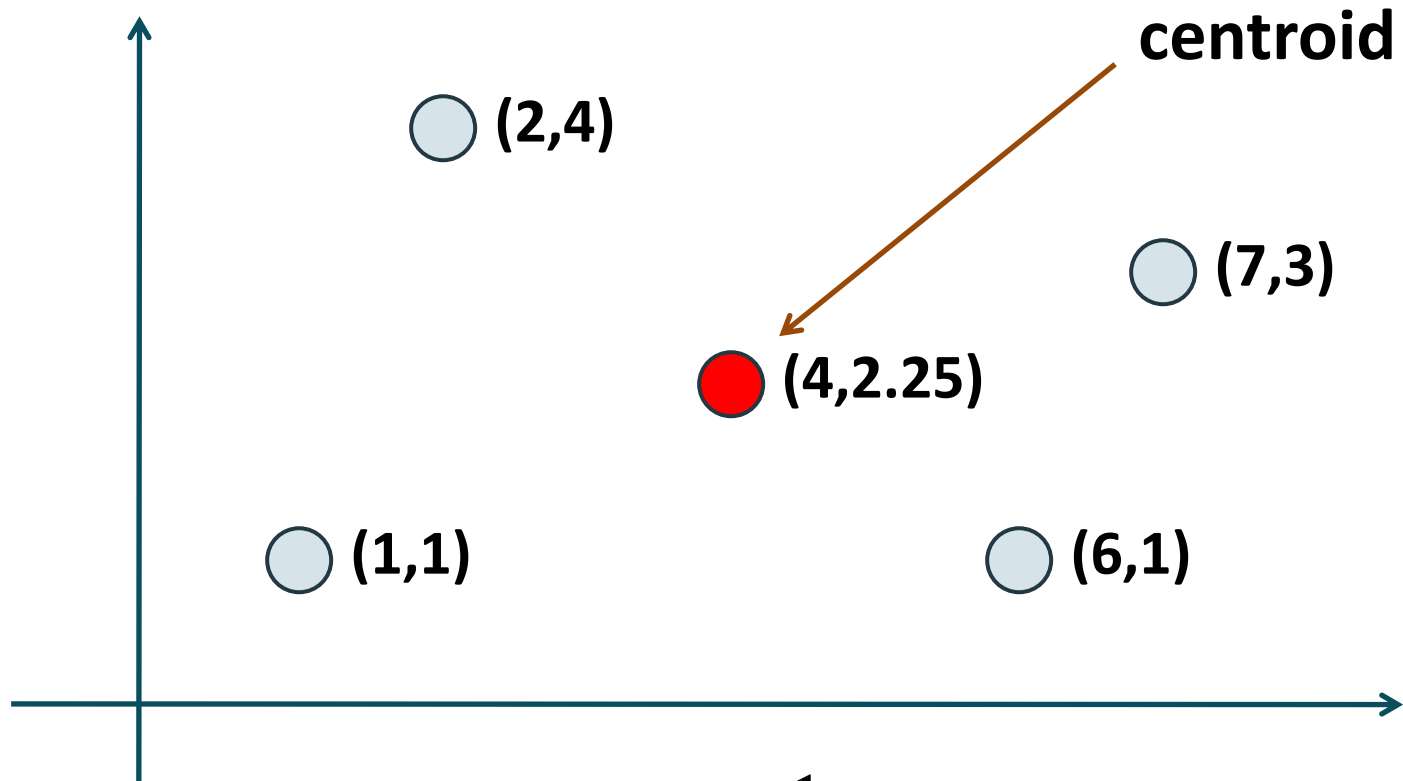
가 Centroid

where C is a set of documents. c

- Recall that we represent documents as points in a high-dimensional vector space.
- The centroid is the center of mass of a set of points.

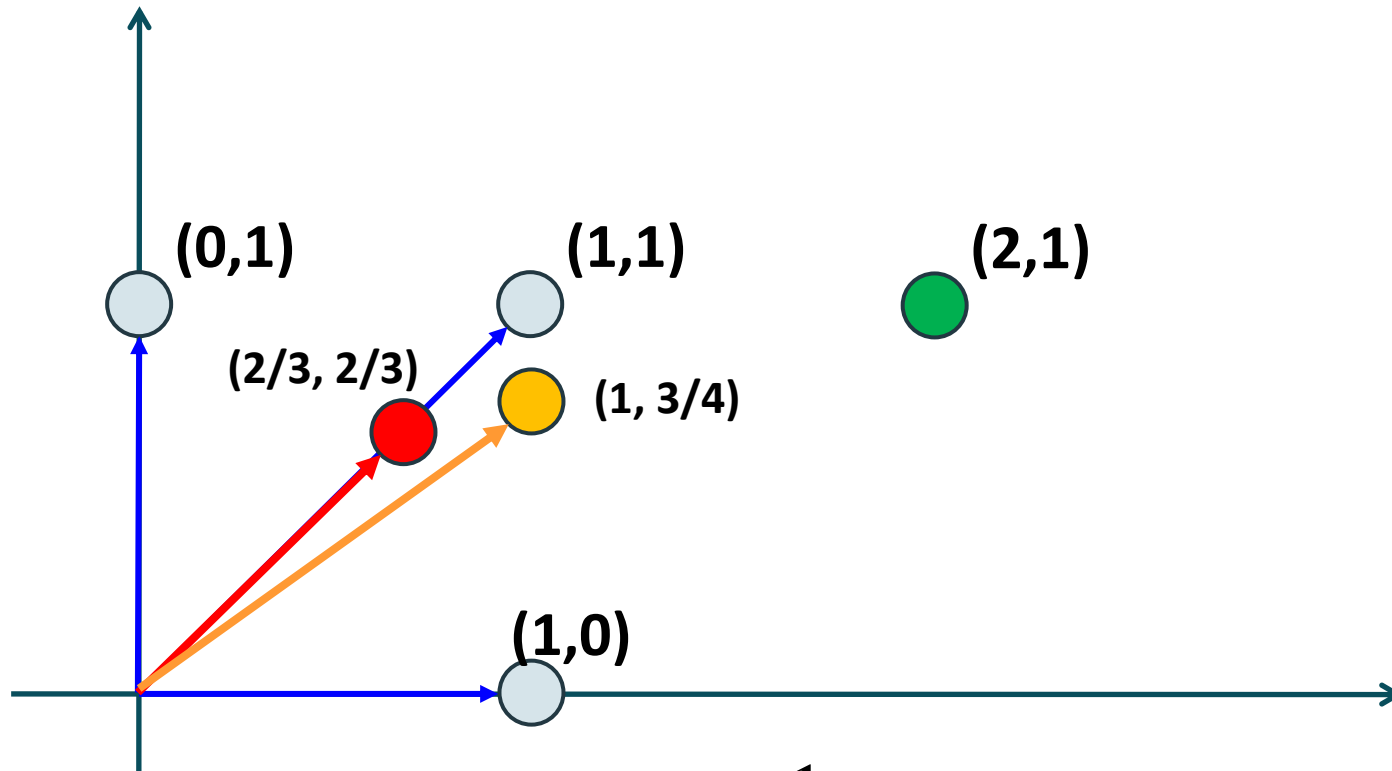
centroid

Key Concept : Centroid



$$\vec{\mu}(C) = \frac{1}{|C|} \sum_{d \in C} \vec{d}$$

Key Concept : Centroid



$$\vec{\mu}(C) = \frac{1}{|C|} \sum_{d \in C} \vec{d}$$

Theoretically Optimal Query

- We want to find a query vector \vec{q} that **maximizes similarity with relevant docs** while **minimizing similarity with non-relevant docs** *in the collection*.

$$\vec{q}_{opt} = \underset{\vec{q}}{\operatorname{argmax}} [sim(\vec{q}, C_r) - sim(\vec{q}, C_{nr})]$$

where C_r is a set of relevant documents and C_{nr} is a set of non-relevant documents in the collection.

$\underset{\vec{q}}{\operatorname{argmax}}$ 가 $\underset{q}{\operatorname{argmax}}$

Theoretically Optimal Query

- In a vector space model, the optimal query is:

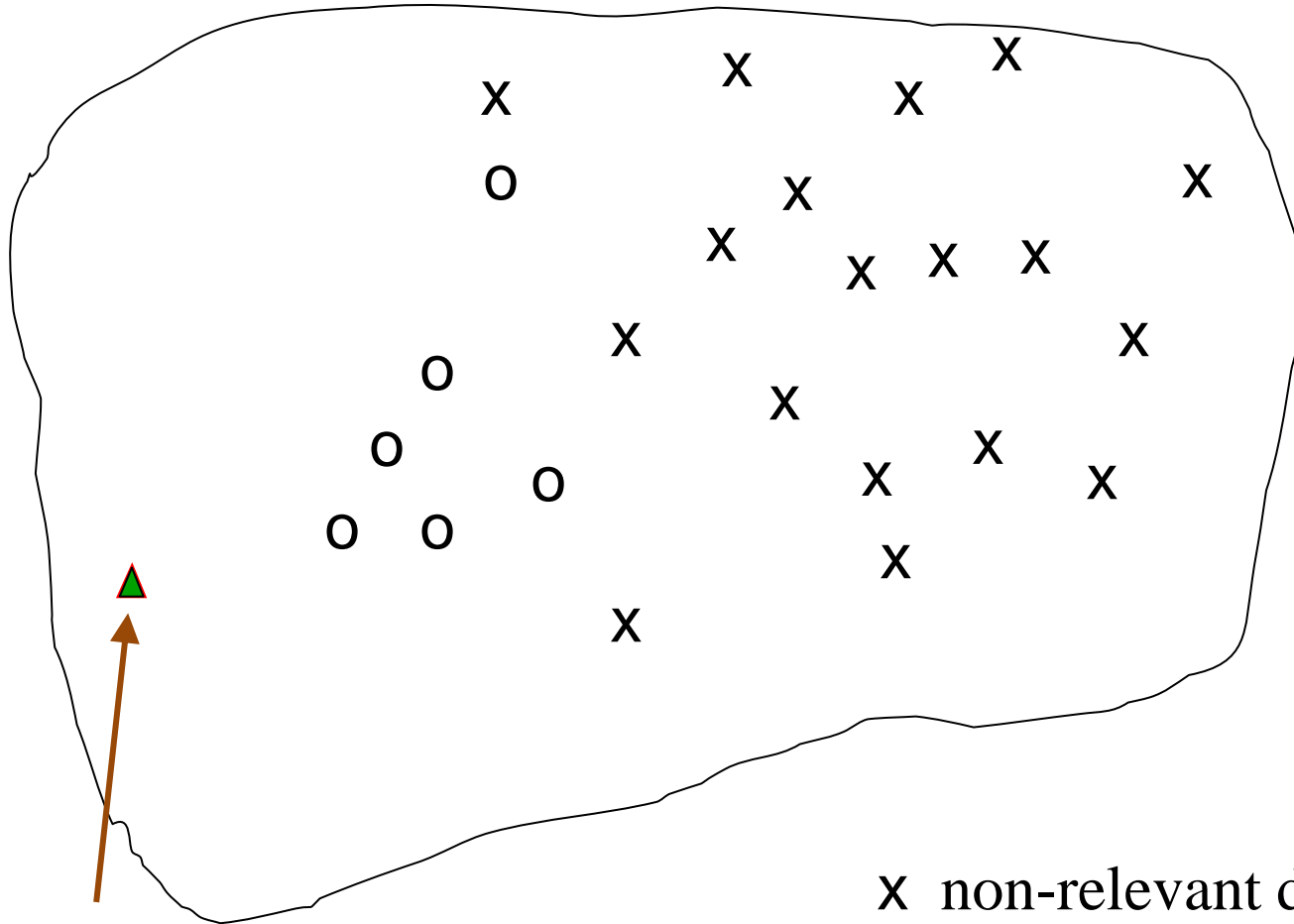
$$\vec{q}_{opt} = \frac{1}{|C_r|} \sum_{\vec{d}_j \in C_r} \vec{d}_j - \frac{1}{|C_{nr}|} \sum_{\vec{d}_j \in C_{nr}} \vec{d}_j \quad \text{centroid} \quad \text{centroid}$$

which is the vector difference between the centroids of the relevant and non-relevant document in the collection.

$$\vec{\mu}(C) = \frac{1}{|C|} \sum_{d \in C} \vec{d}$$

Theoretically Optimal Query

Relevance for all documents in the collection



Optimal Query

x non-relevant documents

o relevant documents

Theoretically Optimal Query

- The problem of theoretically optimal query:
 - The **full set** of relevant and non-relevant documents is **not known in advance**.
 - In fact, it is what we want to find!!! !
- In practice, we have a user query and **partial knowledge** of known relevant and non-relevant documents for the retrieved documents.
- We will use the partial knowledge to modify the query so that we can decide **pseudo optimal query**.
 - This is a relevance feedback mechanism. 가
 - **Rocchio Algorithm : Relevance Feedback Algorithm**

Rocchio Algorithm

Relevance Feedback Algorithm

Rocchio (1971) Algorithm

- Rocchio algorithm uses the vector space model to pick a relevance fed-back query.

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

where D_r = set of known relevant doc vectors

D_{nr} = set of known non-relevant doc vectors

q_0 = initial query vector

q_m = modified query vector

α, β, γ = weights, hand-chosen or set empirically

D 가 . centroid
 modify 가
 , , 가 :

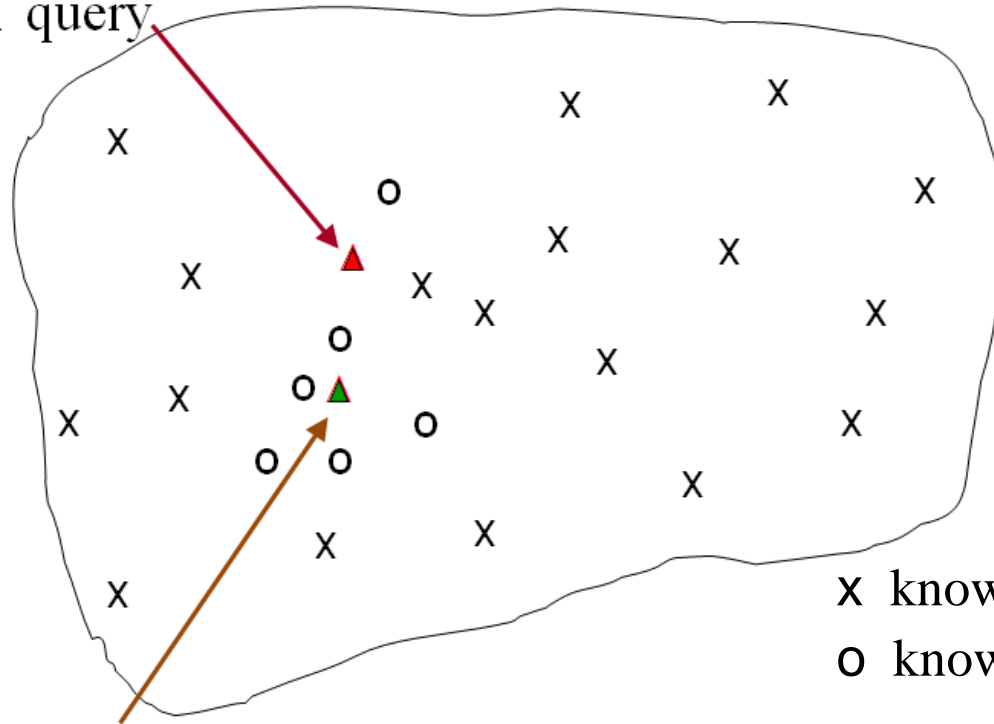
Rocchio (1971) Algorithm

- If we have a lot of judged documents, we want a higher β/γ .
- Positive feedback is more valuable than negative feedback, so, set $\beta > \gamma$.
 - e.g. $\beta = 0.75, \gamma = 0.25$
- Many systems only allow positive feedback ($\gamma=0$).

Rocchio (1971) Algorithm

Modify query with partial knowledge

Initial query



X known non-relevant documents

O known relevant documents

Modified Query

Modified query **moves toward** the centroid of relevant documents and **away from** the centroid of nonrelevant documents.

Optimal Query vs Rocchio Algorithm

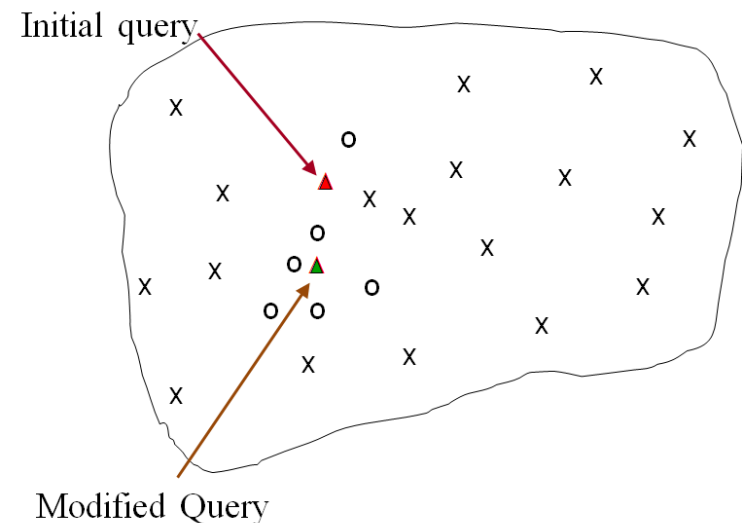
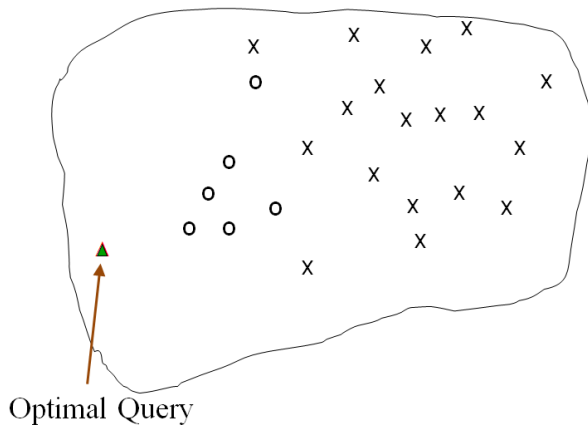
- Optimal query

$$\vec{q}_{opt} = \frac{1}{|C_r|} \sum_{\vec{d}_j \in C_r} \vec{d}_j - \frac{1}{|C_{nr}|} \sum_{\vec{d}_j \in C_{nr}} \vec{d}_j$$

- Rocchio algorithm

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

Why \vec{q}_0 is added here?



Rocchio (1971) Algorithm

- In summary,
 - We can modify the query based on relevance feedback and search again with the modified query.
 - Modify the query only with the docs that were marked relevant or non-relevant.
 - Relevance feedback is most useful for increasing *recall* in situations where recall is important.
 - Users can be expected to review results and to take time to iterate on the search.

Pseudo Relevance Feedback

Pseudo Relevance Feedback

- Also known as ***blind relevance feedback***.
- Pseudo relevance feedback automates the “manual” part of true relevance feedback.
- Pseudo Relevance Feedback algorithm:
 - Retrieve a ranked list of hits for the user’s query.
 - Assume that the top k documents are relevant.
 - Do relevance feedback (e.g. Rocchio).

가 k 가 ' '

Pseudo Relevance Feedback

- Works very well on average. 가
- But can go horribly wrong for some queries. 가
- Several iterations can cause query drift.
 - Assume that the initial query is “copper mines”.
 - The top k documents are all about mines in Chile.
 - Then, the modified query may drift in the direction of documents on Chile.

1. 가 " : "
2. k
3. 가 modify :

Indirect Relevance Feedback

Indirect Relevance Feedback

- Also known as *implicit relevance feedback*.
- Idea:
 - Assumption is that document summaries displayed in results lists are indicative of the relevance of these documents.
 - Then, clicks on links indicate that the page was likely relevant to the query.
 - The data about click rates on pages were gathered globally, rather than being user or query specific.
- Thus, **rank documents higher** that users look at more often.

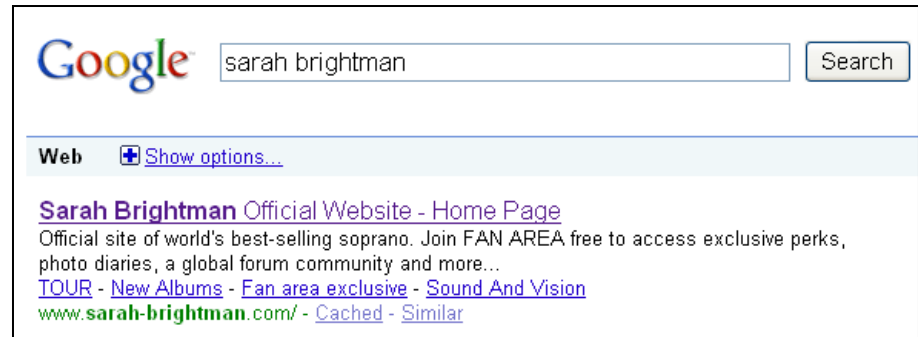
Indirect Relevance Feedback

- This is the general area of **click stream mining**.
- Similar approach is used in ranking the advertisement that match a web search query.
- Less reliable than explicit feedback, but is more useful than pseudo relevance feedback.
- Easy to collect indirect feedback in large quantities, esp. on a web search engine.

and more ...

Relevance Feedback on the Web

- Some search engines offer a *similar/related* pages feature (this is a trivial form of relevance feedback).
 - Google (similar pages)
 - Altavista (more pages)
- Some don't because it's hard to explain to average users.
 - Alltheweb (closed in 2011)
 - Yahoo



Excite Relevance Feedback

- *Excite* initially had true relevance feedback, but dropped it due to lack of use. *Excite*
- Research by Spink et al. 2000
 - Relevance feedback improved results about 2/3 of the time. *2/3*
 - Only about 4% of query sessions from a user used relevance feedback option. *4%*
 - Expressed as “*More like this*” link next to each result.
 - But about 70% of users only looked at first page of results and didn’t pursue things further. *70%* *가*

Problems of Relevance Feedback

- Users are often reluctant to provide explicit feedback.
- It's often harder to understand why a particular document was retrieved after applying relevance feedback.

Problems of Relevance Feedback

- Long (modified) queries are inefficient for typical IR engine.
 - Long response times for user. (가)
 - High computing cost for retrieval system.
 - Partial solution
reweight certain prominent terms only
e.g. reweight only top 20 terms by term frequency

Query: New space satellite applications

2.074	new	15.12	space
30.82	satellite	5.660	application
5.991	nasa	5.196	eos
4.196	launch	3.972	aster
3.516	instrument	3.446	arianespace
3.004	bundespost	2.806	ss
2.790	rocket	2.053	scientist
2.003	broadcast	1.172	earth
0.836	oil	0.646	measure

Evaluation of Relevance Feedback

가

Evaluation of Relevance Feedback

- Empirically, one round of relevance feedback is often very useful. Two rounds is sometimes marginally useful.
가
- **Evaluation Strategy 1**
 - Use q_0 and compute precision and recall graph.
 - Use q_m and compute precision recall graph.
 - Assess on all documents in the collection.
 - Spectacular improvements, but ... it's cheating!
 - Partly due to known relevant documents ranked higher.
 - Must evaluate with respect to documents not seen by user.

가 1 :

가 .

Evaluation of Relevance Feedback

Query: New space satellite applications

**retrieved
documents**

+
+

**relevance
feedback**

1. 0.539, 08/13/91, NASA Hasn't Scrapped Imaging Spectrometer
2. 0.533, 07/09/91, NASA Scratches Environment Gear From Satellite Plan
3. 0.528, 04/04/90, Science Panel Backs NASA Satellite Plan, But Urges Launches of Smaller Probes
4. 0.526, 09/09/91, A NASA Satellite Project Accomplishes Incredible Feat: Staying Within Budget

**improved
results**

1. 0.513, 07/09/91, NASA Scratches Environment Gear From Satellite Plan
2. 0.500, 08/13/91, NASA Hasn't Scrapped Imaging Spectrometer
3. 0.493, 08/07/89, When the Pentagon Launches a Secret Satellite, Space Sleuths Do Some Spy Work of Their Own
4. 0.493, 07/31/89, NASA Uses 'Warm' Superconductors For Fast Circuit

Evaluation of Relevance Feedback

■ Evaluation Strategy 2

- Assess only the documents *not* rated by the user in the first round (*residual collection*).
- Measures are usually lower than for original query because a fair proportion of them have been judged by the user in the first round.
- So, can be used to compare **relative performance** of variant relevance feedback methods.
- Cannot be used to compare performance with and without relevance feedback because the collection size and the number of relevant documents changes.

가 2 : , .

가

Evaluation of Relevance Feedback

■ Evaluation Strategy 3

- Most satisfactory strategy.
- Use two collections each with their own relevance assessments.
 - q_0 and user feedback from first collection
 - q_0 and q_m run on second collection to measure the performance.

가 3 :

Evaluation of Relevance Feedback

- Relevance Feedback vs User's revision of query가
 - User revises and resubmits query.
Users may prefer revision/resubmission to judging relevance of documents.
- True evaluation of usefulness must compare to other methods taking the same amount of time.
- There is no clear evidence that relevance feedback is the “best use” of the user's time.

가

가

Global Methods

1. Query Expansion
2. Automatic Thesaurus Generation

1. Query Expansion

- **Relevance Feedback**

- Users give additional input (relevant/non-relevant) **on documents**, which is used to reweight terms in the documents.

- **Query Expansion**

:

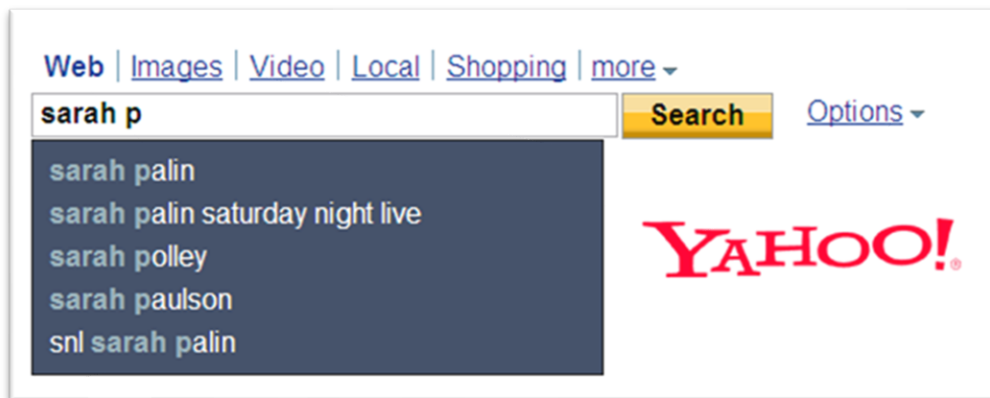
가

- Users give additional input **on query words or phrases**.
- Some search engines suggest related query terms in response to a query.
- **Thesaurus** is often used in expanding queries.

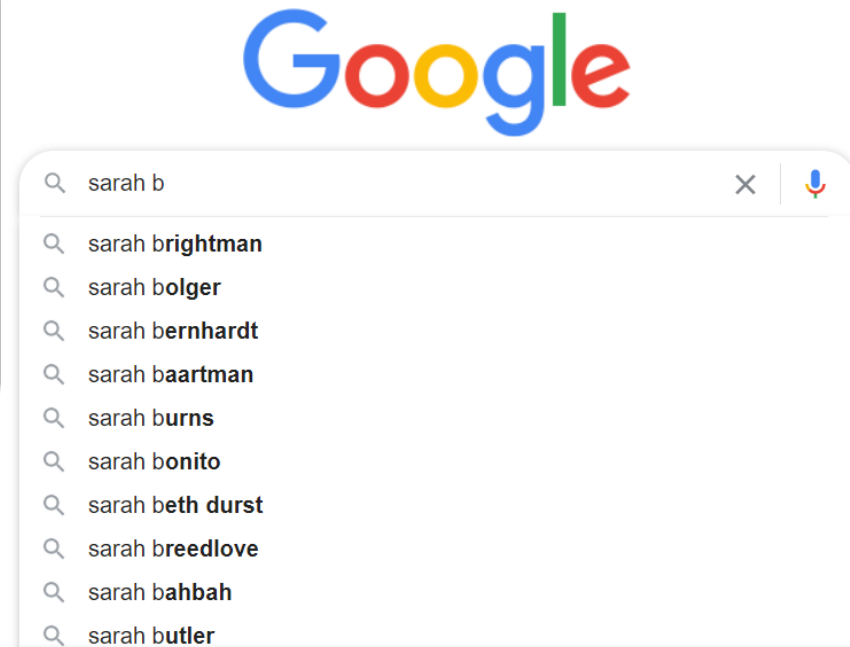
, 가

Query Assist

- Suggest related terms based on **query log mining**.
 - Recommend frequent recent queries that contain partial string typed by user.
 - **Users choose** one from the suggestions.



가



Query Assist

- Suggest related terms based on **query log mining**.
 - Recommend frequent recent queries that contain partial string typed by user.
 - **Users choose** one from the suggestions.

[Yahoo!](#) [My Yahoo!](#) [Mail](#)

[Web](#) | [Images](#) | [Video](#) | [Local](#) | [Shopping](#) | [more](#)

[Options](#)

1 - 10 of about 534,000,000 for

Also try: [palm trees](#), [palm springs](#), [palm centro](#), [palm tree](#), [More...](#)

SPONSOR RESULTS

[Palm - AT&T](#)
[att.com/wireless](#) - Go mobile effortlessly with the **PALM** Treo from AT&T (Cingular).

Thesaurus-based Query Expansion

- For each term t in a query, expand the query with synonyms and related words of t from the thesaurus.
 - feline \rightarrow feline cat
 - Weight may be added to expanded terms, but less than that for original query terms.
- Generally query expansion increases recall.
 - Thus, widely used in many science and engineering fields.
 - However, precision may decrease significantly.

가

가

가

Thesaurus-based Query Expansion

- Use manual thesaurus
 - Thesaurus is a dictionary containing terms such as preferred terms, synonyms, broader terms, narrower terms, ...
 - e.g. MedLine : physician, syn: doc, doctor, MD, medico

User query: cancer

PubMed query: (“neoplasms”[TIAB] NOT Medline[SB]) OR “neoplasms”[MeSH Terms] OR cancer[Text Word]

PubMed does automatic query expansion.

Thesaurus-based Query Expansion

- Drawbacks of manual thesaurus
 - High cost of manually producing a thesaurus.
 - Too general manual thesaurus has little coverage of rich domain-specific vocabularies.
- automatic generation is preferred.

()

()

2. Automatic Thesaurus Generation

- Automatic thesaurus is generated based on the co-occurrence statistics over a collection of documents in a domain.
- Fundamental notion
 - Capture **similarity between two words** from the collection.

가

(co-occurrence)

2. Automatic Thesaurus Generation

- Definition 1 가 / /

Two words are similar if they co-occur in the same document, paragraph or sentence.

Query expansion is often effective in increasing **recall**. However, **query expansion** may also significantly decrease **precision**. In general, a domain-specific **thesaurus** is required.

query expansion	recall precision thesaurus
-----------------	----------------------------

가

가

2. Automatic Thesaurus Generation

: 가

(keeping)

Word	Nearest neighbors
absolutely bottomed captivating doghouse makeup	absurd, whatsoever, totally, exactly, nothing dip, copper, drops, topped, slide, trimmed shimmer, stunningly, superbly, plucky, witty dog, porch, crawling, beside, downstairs repellent, lotion, glossy, sunscreen, skin, gel
<u>mediating</u> <u>keeping</u>	reconciliation, negotiate, case, conciliation hoping, bring, wiping, could, some, would
<u>lithographs</u>	drawings, Picasso, Dali, sculptures, Gauguin
pathogens senses	toxins, bacteria, organisms, bacterial, parasite grasp, psyche, truly, clumsy, naive, innate

lithographs :

2. Automatic Thesaurus Generation

- Definition 2

Two words are similar if they occur in a given grammatical relation with the same words.

They harvest apples. **He peels apples.** She eat apples.
Jane harvests pears. **Tom peels pears.** They eat pears.

Apples and Pears must be similar.

- Co-occurrence based thesaurus is more robust, grammatical relation based thesaurus is more accurate.

2. Automatic Thesaurus Generation

- Quality of word associations is usually a problem.
 - Term ambiguity may introduce statistically correlated but semantically irrelevant terms. 가 가
 - **False positives:** words deemed similar that are not.
 - **False negatives:** words deemed dissimilar that are similar.
- Bad News

Since terms are highly correlated anyway, expansion may not retrieve many additional documents.

()