### **Chapter 12: Query Processing**

Computers are useless, they can only give you answers.

-- Pablo Picasso

You have to think anyway, so why not think big?

-- Donald Trump

There are lies, damn lies, and workload assumptions.

-- anonymous







### Outline

- **12.1 Query Processing Algorithms**
- 12.2 Fast Top-k Search
- **12.3 Phrase and Proximity Queries**
- **12.4 Query Result Diversification**



#### loosely following Büttcher/Clarke/Cormack Chapters 5 and 8.6 plus Manning/Raghavan/Schütze Chapters 7 and 9 plus specific literature

# **Query Types**

#### • Conjunctive

(i.e., all query terms are required)

#### • Disjunctive

(i.e., subset of query terms sufficient)

#### • Phrase or proximity

(i.e., query terms must occur in right order or close enough)

# Mixed-mode with negation (e.g., "harry potter" review +movie -book)

• Combined with **ranking of result documents** according to

$$score(q, d) = \sum_{t \in q} score(t, d)$$

with score(t, d) depending on retrieval model (e.g. tf\*idf)

# **Indexing with Document-Ordered Lists**





index-list entries stored in ascending order of document identifiers (document-ordered lists)

process all queries (conjunctive/disjunctive/mixed) by sequential scan and merge of posting lists

# **Document-at-a-Time Query Processing**

**Document-at-a-Time (DAAT)** query processing

- assumes document-ordered posting lists
- scans posting lists for query terms  $t_1, \ldots, t_{|q|}$  concurrently
- maintains an **accumulator** for each candidate result doc:

$$- \ acc(d) = \sum_{i: \ d \ seen \ in \ L(ti)} score(ti, d)$$

$$a \cdots d_{1}, 1.0 \ d_{4}, 2.0 \ d_{7}, 0.2 \ d_{8}, 0.1$$

$$b \cdots d_{4}, 1.0 \ d_{7}, 2.0 \ d_{8}, 0.2 \ d_{9}, 0.1$$

$$c \cdots d_{4}, 3.0 \ d_{7}, 1.0$$
Accumulators
$$d_{1} \ : \ 1.0 \\ d_{4} \ : \ 6.0 \\ d_{7} \ : \ 3.2 \\ d_{8} \ : \ 0.3 \\ d_{9} \ : \ 0.1$$

- always advances posting list with lowest current doc id
- exploit **skip pointers** when applicable
- required memory depends on # results to be returned
- top-k results in priority queue

**IRDM WS 2015** 

# DAAT with Weak And: WAND Method

[Broder et al. 2003]

#### **Disjunctive** (Weak And) query processing

- assumes document-ordered posting lists with known
   maxscore(i) values for each t<sub>i</sub>: max<sub>d</sub> (score (d,t<sub>i</sub>))
- While scanning posting lists keep track of
  - **min-k:** the lowest total score in current top-k results
  - ordered term list: terms sorted by docId at current scan pos
  - **pivot term:** smallest j such that  $min-k \leq \sum_{i \leq j} maxscore(i)$
  - pivot doc: doc id at current scan pos in posting list Lj

Eliminate docs that cannot become top-k results (maxscore pruning)

- **if** pivot term does not exist (min-k >  $\sum_i maxscore(i)$ )
- then stop
- else advance scan positions to pos  $\geq$  id of pivot doc ("big skip")

### **Example: DAAT with WAND Method**

[Broder et al. 2003]

**Key invariant**: For terms i=1..|q| and current scan positions cur<sub>i</sub> assume that cur<sub>1</sub> = min {curi | i=1..|q|} Then for each posting list i there is no docid between cur<sub>1</sub> and cur<sub>i</sub>



Suppose that min-k = 12 then the pivot term is 4  $(\Sigma_{i=1.3} \text{ maxscore}_i > \min\text{-k}, \Sigma_{i=1.4} \text{ maxscore}_i \le \min\text{-k})$ and the pivot docid is 600  $\rightarrow$  can advance all scan positions cur<sub>i</sub> to 600

# **Term-at-a-Time Query Processing**

#### Term-at-a-Time (TAAT) query processing

- assumes document-ordered posting lists
- scans posting lists for query terms  $t_1, \ldots, t_{|q|}$  one at a time, (possibly in decreasing order of idf values)
- maintains an **accumulator** for each candidate result doc
- after processing L(tj):  $acc(d) = \sum_{i \le j} score(ti, d)$

Accumulators

$a \cdots d_{1},$	<b>1.0</b> <i>d</i> <sub>4</sub> , <b>2.0</b>	<i>d</i> 7, 0.2	<i>d</i> 8, <b>0.1</b>	<b>d</b> 1	:	<b>9.</b> 0
$b \cdots b \frac{d_4}{d_4}$	<b>1.0</b> <i>d</i> <sub>7</sub> , <b>2.0</b>	<i>d</i> <sub>8</sub> , 0.2	d9, 0.1	u4 d7	:	9.0 9.9
c ····· d4,	3.0 <i>d</i> <sub>7</sub> , 1.0			ds d9	•	0.9 0.9

- memory depends on the number of accumulators maintained
- TAAT is attractive when scanning many short lists

# **Indexing with Impact-Ordered Lists**

#### Data items: $d_1, \ldots, d_n$



index-list entries stored in descending order of per-term score impact (impact-ordered lists)

aims to avoid having to read entire lists rather scan only (short) prefixes of lists

### **Greedy Query Processing Framework**

Assume index lists are sorted by  $tf(t_i,d_j)$  or  $tf(t_i,d_j)*idl(d_j)$  values idf values are stored separately

Open scan cursors on all m index lists L(i) **Repeat** 

Find **pos(g)** among current cursor positions pos(i) (i=1..m) with the **largest value of idf(t<sub>i</sub>)\*tf(t<sub>i</sub>,d<sub>j</sub>)** (or idf(t<sub>i</sub>)\*tf(t<sub>i</sub>,d<sub>j</sub>)\*idl(d<sub>j</sub>)); Update the accumulator of the corresponding doc;

Increment pos(g);

Until stopping condition

### **Stopping Criterion: Quit & Continue Heuristics**

m

[Zobel/Moffat 1996]

For scoring of the form score
$$(q, d_j) = \sum_{i=1}^{m} s_i(t_i, d_j)$$
  
with  $s_i(t_i, d_j) \sim tf(t_i, d_j) \cdot idf(t_i) \cdot idl(d_j)$ 

Assume hash array of *accumulators for* summing up score mass of candidate results

*quit* heuristics (with docId-ordered or tf-ordered or tf\*idl-ordered index lists):

- ignore index list L(i) if  $idf(t_i)$  is below tunable threshold or
- stop scanning L(i) if  $idf(t_i)*tf(t_i,d_i)*idl(d_i)$  drops below threshold or
- stop scanning L(i) when the number of accumulators is too high

#### *continue* heuristics:

upon reaching threshold, continue scanning index lists, but do not add any new documents to the accumulator array

# 12.2 Fast Top-k Search

**Top-k aggregation query** over relation *R* (*Item, A1, ..., Am*): *Select Item, s(R1.A1, ..., Rm.Am) As Aggr From Outer Join R1, ..., Rm Order By Aggr Limit k* with monotone s:  $(\forall i: x_i \ge x_i^{\circ}) \Rightarrow s(x_1 \dots x_m) \ge s(x_1^{\circ} \dots x_m^{\circ})$ (example: item is doc, attributes are terms, attr values are scores)

- Precompute per-attr (index) **lists sorted in desc attr-value order** (score-ordered, impact-ordered)
- Scan lists by **sorted access (SA)** in round-robin manner
- Perform random accesses (RA) by Item when convenient
- Compute aggregation s **incrementally** in **accumulators**
- Stop when **threshold test** guarantees correct top-k (or when heuristics indicate ,,good enough" approximation)

#### simple & elegant, adaptable & extensible to distributed system

following R. Fagin: Optimal aggregation algorithms for middleware, JCSS. 66(4), 2003 IRDM WS 2015

### **Threshold Algorithm (TA)**

[Fagin 01,Güntzer 00, Nepal 99, Buckley 85]

simple & DB-style; needs only O(k) memory

Data items:  $d_1, \ldots, d_n$ 



Query: 
$$q = (t_1, t_2, t_3)$$

**Threshold algorithm (TA):** scan index lists; consider d at pos<sub>i</sub> in L<sub>i</sub>;  $high_i := s(t_i,d);$ if  $d \notin \text{top-k}$  then { **look up**  $s_{\nu}(d)$  in all lists  $L_{\nu}$  with  $\nu \neq i$ ; **score**(**d**) := **aggr** { $s_v$ (**d**) | v=1..m}; if score(d) > min-k then add d to top-k and remove min-score d';  $\min-k := \min\{score(d') \mid d' \in top-k\};\$ threshold := aggr {high<sub>v</sub> | v=1..m}; if threshold ≤ min-k then exit;



### TA with Sorted Access only (NRA) [Fagin 01, Güntzer 01]

sequential access (SA) faster than random access (RA) by factor of 20-1000

Data items:  $d_1, \ldots, d_n$ 



Query: 
$$q = (t_1, t_2, t_3)$$

No-random-access algorithm (NRA): scan index lists; consider d at pos<sub>i</sub> in L<sub>i</sub>;  $E(d) := E(d) \cup \{i\}$ ; high<sub>i</sub> :=  $s(t_i,d)$ ; worstscore(d) := aggr{ $s(t_v,d) \mid v \in E(d)$ }; bestscore(d) := aggr{worstscore(d),  $aggr{high_v \mid v \notin E(d)$ }; if worstscore(d) > min-k then add d to top-k min-k := min{ $worstscore(d') \mid d' \in top-k$ }; else if bestscore(d) > min-k then  $cand := cand \cup \{d\}$ ; threshold := max {bestscore(d') | d' \in cand}; if threshold ≤ min-k then exit;



# **TA Complexity and Instance Optimality**

m-1

TA has worst-case run-time  $O(n^{\overline{m}})$  with high prob. and space O(1)NRA has worst-case run-time O(n) and space O(n)

#### **Definition**:

For class  $\mathcal{A}$  of algorithms and class  $\mathcal{D}$  of datasets, algorithm B is **instance optimal** over  $\mathcal{A}$  and  $\mathcal{D}$  if for every  $A \in \mathcal{A}$  on  $D \in \mathcal{D}$ : cost(B,D)  $\leq c * O(cost(A,D)) + c'$  $(\rightarrow \text{competitiveness c}).$ 

#### Theorem:

- **TA is instance optimal** over all algorithms that are based on sorted and random accesses to m lists (no ,,wild guesses").
- NRA is instance optimal over all algorithms with SA only.

if "wild guesses" are allowed, then no deterministic algorithm is instance-optimal

IRDM WS 2015

# Implementation Issues for TA Family

- Limitation of asymptotic complexity:
  - m (#lists) and k (#results) are important parameters
- Priority queues:
  - straightforward use of Fibonacci heap has high overhead
  - better: periodic rebuild of bounded-size PQs
- Memory management:
  - peak memory use as important for performance as scan depth
  - aim for early candidate pruning even if scan depth stays the same
- Hybrid block index:
  - pack index entries into big blocks in desc score order
  - keep blocks in score order
  - keep entries within a block in item id order
  - after each block read: merge-join first, then PQ update

# **Approximate Top-k Answers**



- IR heuristics for impact-ordered lists [Anh/Moffat: SIGIR'01]: Accumulator Limiting, Accumulator Thresholding
- Approximation TA [Fagin et al.2003] :
  - *θ-approximation* T' for q with  $\theta > 1$  is a set T' of items with:
  - |T'|=k and
  - for each  $d^{*} \in T^{*}$  and each  $d^{*} \notin T^{*}$ :  $\theta * \text{score}(q,d^{*}) \ge \text{score}(q,d^{*})$ <u>Modified TA:</u>
    - ... stop when *min-k*  $\geq$  *aggr* (*high*<sub>1</sub>, ..., *high*<sub>m</sub>) /  $\theta$
- Probabilistic Top-k [Theobald et al. 2004] : guarantee small deviation from exact top-k result with high probability

### **Probabilistic Top-k Answers**



TA family of algorithms based on invariant (with sum as aggr):



- Add *d* to top-k result, if worstscore(d) > min-k
- Drop *d* only if *bestscore(d) < min-k*, otherwise keep in PQ
- → Often overly conservative (deep scans, high memory for PQ)
- → Approximate top-k with probabilistic guarantees:



 $\Rightarrow$  E[rel. precision@k] = 1- $\varepsilon$ 

 $p(d) := P[\sum_{i \in E(d)} s_i(d) + \sum_{i \notin E(d)} S_i > \delta] \quad \text{with } \delta = \min -k$ 

discard candidates d from queue if  $p(d) \le \varepsilon$ 

# Combined Algorithm (CA) for Balanced SA/RA Scheduling [Fagin et al. 03]

cost ratio  $C_{RA}/C_{SA} = r$ 

. . .

```
perform NRA (TA-sorted)
```

after every r rounds of SA (m\*r scan steps) perform RA to look up all missing scores of "best candidate" in Q

cost competitiveness w.r.t. ,,optimal schedule" (scan until  $\Sigma_i$  high<sub>i</sub>  $\leq$  min{bestscore(d) | d  $\in$  final top-k}, then perform RAs for all d' with bestscore(d') > min-k): 4m + k

# Flexible Scheduling of SA's and RA's for Top-k Query Processing

#### Goals:

- 1. decrease **high<sub>i</sub> upper-bounds** quickly
  - $\rightarrow$  decreases bestscore for candidates
  - $\rightarrow$  reduces candidate set
- 2. reduce **worstscore-bestscore gap** for most **promising candidates** 
  - $\rightarrow$  increases min-k threshold
  - $\rightarrow$  more effective threshold test for other candidates

#### Ideas for better scheduling:

- 1. Non-uniform choice of SA steps in different lists
- 2. Careful choice of postponed RA steps for promising candidates when worstscore is high and worstscore-bestscore gap is small



choose  $b_i$  values so as to achieve high score reduction  $\delta$ 

 carefully chosen RAs: score lookups for "interesting" candidates



compute top-1 result using flexible SAs and RAs



Scheduling Example							
	$L_1$	L <sub>2</sub>		L <sub>3</sub>			
	A: 0.8	G: 0.7		Y: 0.9	_		
	B: 0.2	H: 0.5		A: 0.7			
	K: 0.19	R: 0.5		P: 0.3			
	F: 0.17	Y: 0.5		F: 0.25			
	M: 0.16	W: 0.3		S: 0.25			
	Z: 0.15	D: 0.25	5	T: 0.2			
	W: 0.1	W: 0.2		Q: 0.15			
	Q: 0.07	A: 0.2		X: 0.1			
	E	E		÷			
са	indidates:	A: [1.5, 2.0] G: [0.7, 1.6]	Y: [0.9, 1.6] ?: [9.0, 1.4]				





# Top-k Queries on Internet Sources 🛚 🗮



[Marian et al. 2004]

#### Setting:

- score-ordered lists dynamically produced by Internet sources
- some sources restricted to lookups only (no lists)

#### Example:

preference search for hotel based on *distance*, *price*, *rating* using mapquest.com, booking.com, tripadvisor.com

#### Goal:

good **scheduling** for (parallel) access to restricted sources: **SA-sources, RA-sources, universal sources** with different costs for SA and RA

#### Method (Idea):

- scan all SA-sources in **parallel**
- in each step: choose next SA-source or perform RA on RA-source or universal source with best **benefit/cost** contribution

# **Top-k Rank Joins on Structured Data**

[Ilyas et al. 2008]



Select R.Name, C.Theater, C.Movie From RestaurantsGuide R, CinemasProgram C Where R.City = C.City Order By R.Quality/R.Price + C.Rating Desc

#### RestaurantsGuide

Name	Туре	Quality	Price	City
BlueDragon	Chinese	* * *	€15	SB
Haiku	Japanese	* * *	€30	SB
Mahatma	Indian	* *	€20	IGB
Mescal	Mexican	* *	€10	IGB
BigSchwenk	German	* *	€25	SLS

#### CinemasProgram

Theater	Movie	Rating	City
BlueSmoke Oscar's Holly's GoodNight BigHits 	Tombstone Hero Die Hard Seven Godfather	7.5 8.2 6.6 7.7 9.1	SB SB IGB IGB

# DAAT, TAAT, Top-k: Lessons Learned

- TA family over impact-ordered lists
  - is most elegant and potentially most efficient
  - but depending on score skew, it may degrade badly
- DAAT over document-ordered lists
  - is most versatile and robust
  - has lowest overhead and still allows pruning
  - can be easily scaled out on server farm
- TAAT is of interest for special use-cases
   (e.g. patent search with many keywords in queries)

### **12.3 Phrase Queries and Proximity Queries**

#### phrase queries such as:

*"Star Wars Episode 7", "The Force Awakens", "Obi Wan Kenobi", "dark lord" "Wir schaffen das", "to be or not to be", "roots of cubic polynomials", "evil empire "* 

difficult to anticipate and index all (meaningful) phrases sources could be thesauri/dictionaries or query logs

 $\rightarrow$  standard approach:

combine single-term index with separate position index

4	term	doc	score	term	doc	offset	N
				 emnire	30	191	
<b>EX</b>	empire	77	0.85	empire	e 77	375	
pd /	empire	39	0.82	•••			SO
				evil	12	45	
	evil	49	0.81	evil	39	190	) n
	evil	39	0.78	evil	39 40	194	
	evil	12	0.75	evii	49	190	lde
	evil	77	0.12	 evil	77	190	
IRDM WS 2015		, ,	0.12	•••			

### **Bigram and Phrase Indexing**

#### build index over all word pairs (bigrams):

index lists (term1, term2, doc, score) or

for each term1 nested list (term2, doc, score)

variations:

- treat nearest nouns as pairs, or discount articles, prepositions, conjunctions
- index phrases from query logs, compute correlation statistics

#### query processing by merging posting lists:

- decompose even-numbered phrases into bigrams
- decompose odd-numbered phrases into bigrams with low selectivity (as estimated by df(term1))
- may additionally use standard single-term index if necessary

Examples:

to be or not to be  $\rightarrow$  (to be) (or not) (to be) The Lord of the Rings  $\rightarrow$  (The Lord) (Lord of) (the Rings)

IRDM WS 2015

### **Proximity Search**

#### Example queries: root polynom three, high cholesterol measure, doctor degree defense

<u>Idea:</u> identify positions (pos) of all query-term occurrences and reward short distances

**keyword proximity score** [Büttcher/Clarke: SIGIR'06]: aggregation of per-term scores # + per-term-pair scores attributed to each term

$$score(t_{1}...t_{m}) = \sum_{i=1..m} (score(t_{i}) + \sum_{j\neq i} \left\{ \frac{idf(t_{j})}{(pos(t_{i}) - pos(t_{j}))^{2}} | \neg \exists t_{k} (pos(t_{i}) < pos(t_{k}) < pos(t_{j}) \text{ or } ...) \right\}$$
  
cannot be precomputed count only pairs of query terms

 $\rightarrow$  expensive at query-time

count only pairs of query terms with no other query term in between

### **Example: Proximity Score Computation**

It<sup>1</sup> took<sup>2</sup> the<sup>3</sup> sea<sup>4</sup> a<sup>5</sup> thousand<sup>6</sup> years,<sup>7</sup> A<sup>8</sup> thousand<sup>9</sup> years<sup>10</sup> to<sup>11</sup> trace<sup>12</sup> The<sup>13</sup> granite<sup>14</sup> features<sup>15</sup> of<sup>16</sup> this<sup>17</sup> cliff,<sup>18</sup> In<sup>19</sup> crag<sup>20</sup> and<sup>21</sup> scarp<sup>22</sup> and<sup>23</sup> base.<sup>24</sup>

Query: {sea, years, cliff}



### **Efficient Proximity Search**

**Define aggregation function to be distributive** [Broschart et al. 2007] rather than "holistic" [Büttcher/Clarke 2006]:

precompute term-pair distances and sum up at query-time

$$score(t_1...t_m) = \sum_{i=1..m} \left( score(t_i) + \sum_{j \neq i} \left\{ \frac{idf(t_j)}{(p(t_i) - p(t_j))^2} \right\} \right)$$
  
count all pairs of query terms

#### result quality comparable to "holistic" scores

index all pairs within max. window size (or nested list of nearby terms for each term), with precomputed pair-score mass > >

### **Ex.: Efficiently Computable Proximity Score**

It<sup>1</sup> took<sup>2</sup> the<sup>3</sup> sea<sup>4</sup> a<sup>5</sup> thousand<sup>6</sup> years,<sup>7</sup> A<sup>8</sup> thousand<sup>9</sup> years<sup>10</sup> to<sup>11</sup> trace<sup>12</sup> The<sup>13</sup> granite<sup>14</sup> features<sup>15</sup> of<sup>16</sup> this<sup>17</sup> cliff,<sup>18</sup> In<sup>19</sup> crag<sup>20</sup> and<sup>21</sup> scarp<sup>22</sup> and<sup>23</sup> base.<sup>24</sup>

Query: {sea, years, cliff}

acc(d,cliff,sea) = 
$$\frac{1}{(18-4)^2}$$
  
acc(d,cliff,years) =  $\frac{1}{(18-7)^2}$  +  $\frac{1}{(18-10)^2}$   
acc(d,sea,years) =  $\frac{1}{(7-4)^2}$  +  $\frac{1}{(10-4)^2}$ 

### **Relevance Feedback**

Given: a query q, a result set (or ranked list) D, a user's assessment u:  $D \rightarrow \{+, -\}$ yielding positive docs  $D^+ \subseteq D$  and negative docs  $D^- \subseteq D$ 

Goal: derive query q' that better captures the user's intention, by adapting term weights in the query or by query expansion

Classical IR approach: *Rocchio method* (for term vectors)

$$q' = \alpha q + \frac{\beta}{|D^+|} \sum_{d \in D^+} d - \frac{\gamma}{|D^-|} \sum_{d \in D^-} d \qquad \text{with } \alpha, \beta, \gamma \in [0,1]$$
  
and typically  $\alpha > \beta > \gamma$ 

<u>Modern approach</u>: replace explicit feedback by **implicit feedback** derived from **query&click logs** (pos. if clicked, neg. if skipped)

or rely on **pseudo-relevance feedback**: assume that all top-k results are positive

### **Relevance Feedback using Text Classification or Clustering**

Relevant and irrelevant docs (as indicated by user) form two classes or clusters of text-doc-vector distribution

#### **Classifier:**

- train classifier on relevant docs as positive class
- run feature selection to identify best terms for expansion
- pass results of expanded query through classifier

#### **Clustering:**

- refine clusters or compute sub-space clusters:
- user explores the resulting sub-clusters and guides expansion

Search engine examples: <u>http://exalead.com</u> http://yippy.com

# **Query Expansion**

- Query expansion can be beneficial whenever high recall is needed
- **Expansion terms** can come from thesauri/dictionaries/ontologies or personalized profile, regardless of user feedback
- **Term-term similarities** precomputed from co-occurrence statistics

# Example q: traffic tunnel disasters (from TREC benchmark)

. . .



# WordNet: Thesaurus/Ontology of Words and Concepts



• S: (v) traffic (trade or deal a commodity) "They trafficked with us for gold"

# WordNet: Thesaurus/Ontology of Words and Concepts

#### Noun

	<ul> <li>S: (n) traffic (the aggregation of things (pedestrians or vehicles) coming and going in a particular locality during a specified period of time).</li> </ul>
	particular locality during a specified period of time)
	• <u>arect hyponym</u> ( <u>tuli nyponym</u>
	<ul> <li>S: (n) <u>air traffic</u> (traffic created by the movement of aircraft)</li> </ul>
	<ul> <li>S: (n) commuter traffic (traffic created by people going to or returning from</li> </ul>
1	work)
nypor	<ul> <li>S: (n) pedestrian traffic, foot traffic (people coming and going on foot)</li> </ul>
(sub-	concepts) • S: (n) vehicular traffic, vehicle traffic (the aggregation of vehicles coming
	and going in a particular locality)
	<ul> <li>S: (n) automobile traffic, car traffic (cars coming and going)</li> </ul>
	<ul> <li>S: (n) bicycle traffic (bicycles coming and going)</li> </ul>
	<ul> <li>S: (n) bus traffic (buses coming and going)</li> </ul>
	<ul> <li>S: (n) truck traffic (trucks coming and going)</li> </ul>
	<ul> <li><u>direct hypernym / inherited hypernym / sister term</u></li> </ul>
	<ul> <li>S: (n) traffic (buying and selling; especially illicit trade)</li> </ul>
	• S: (n) traffic (the amount of activity over a communication system during a given period
	of time) "heavy traffic overloaded the trunk lines": "traffic on the internet is lightest
	during the pight"
	<ul> <li>S: (n) dealings, traffic (social or verbal interchange (usually followed by `with'))</li> </ul>
	- <u>o.</u> (ii) <u>dealings,</u> <b>name</b> (social or verbal interchange (usually followed by with))

# WordNet: Thesaurus/Ontology of Words and Concepts

#### Noun

	• <u>S:</u> (n)	) <b>tunnel</b> (a passageway through or under something, usually underground
	(espe	ecially one for trains or cars)) "the tunnel reduced congestion at that intersection
hyponym	5 °	<u>direct hyponym I full hyponym</u>
(aub aana	anta)	<ul> <li><u>S:</u> (n) <u>catacomb</u> (an underground tunnel with recesses where bodies were</li> </ul>
(sub-conc	epis)	buried (as in ancient Rome))
		<ul> <li>S: (n) railroad tunnel (a tunnel through which the railroad track runs)</li> </ul>
		<ul> <li>S: (n) underpass, subway (an underground tunnel or passage enabling</li> </ul>
meronym	C	pedestrians to cross a road or railway)
meronym	•	part meronym
(part-of)		<ul> <li>S: (n) shaft (a long vertical passage sunk into the earth, as for a mine or</li> </ul>
( <b>T</b> )		tunnel)
	0	domain category
		<ul> <li>S: (n) car, auto, automobile, machine, motorcar (a motor vehicle with four</li> </ul>
		wheels; usually propelled by an internal combustion engine) "he needs a
hvpernvm	IS	car to get to work"
, perny m		direct hypernym / inherited hypernym / sister term
(super-coi	ncepts)	<ul> <li>S: (n) passageway (a passage between rooms or between buildings)</li> </ul>
	0	derivationally related form
	• <u>S:</u> (n)	) <u>burrow,</u> <b>tunnel</b> (a hole made by an animal, usually for shelter)

n

### **Robust Query Expansion**

Threshold-based query expansion:<br/>substitute w by  $exp(w):=\{c_1 \dots c_k\}$  with all  $c_i$  with  $sim(w, c_i) \ge \delta$ Naive scoring:<br/> $s(q,d) = \sum_{w \in q} \sum_{c \in exp(w)} sim(w,c) * s_c(d)$ risk of<br/>topic drift

#### Approach to careful expansion and scoring:

- determine phrases from query or best initial query results (e.g., forming 3-grams and looking up ontology/thesaurus entries)
- if **uniquely mapped** to one concept then expand with synonyms and weighted hyponyms
- avoid **undue score-mass accumulation** by expansion terms  $s(q,d) = \sum w \in q \max_{c \in exp(w)} \{ sim(w,c) * s_c(d) \}$

#### Query Expansion with Incremental Merging [M. Theobald et al.: SIGIR 2005]

relaxable query q: ~*professor research* with expansions  $\exp(t) = \{w \mid sim(t,w) \ge \theta, t \in q\}$ based on ontology relatedness modulating monotonic score aggregation by sim(t,w)

TA/NRA scans of index lists for  $\bigcup_{t \in q} exp(t)$ Better: dynamic query expansion with incremental merging of additional index lists





(ontology / thesuarus) professor lecturer: 0.7 scholar: 0.6 academic: 0.53 scientist: 0.5 ...

meta-index

efficient and robust

# **Query Expansion Example**

#### From TREC 2004 Robust Track Benchmark:

#### **Title:** International Organized Crime

**Description:** Identify organizations that participate in international criminal activity, the activity, and, if possible, collaborating organizations and the countries involved.

Search Word: organized crime	Redisplay Overvia
Searches for organized crime: Noun	Senses:
1 sense of organized crime	
<ul> <li>Sense 1</li> <li>organized crime, gangland, gangdom (underworld organizations)</li> <li>=&gt; yakuza (organized crime in Japan; an alliance of criminal organization enterprises)</li> <li>=&gt; Mafia, Maffia, Sicilian Mafia (a secret terrorist group in Sicily; origination but evolved into a criminal organization in the middle of the 19</li> <li>=&gt; Black Hand (a secret terrorist society in the United States early in the States early in the Camorra (a secret society in Naples notorious for violence and black experimental organizate, mob, family (a loose affiliation of ganged organized criminal activities)</li> </ul>	ions and illegal ginally opposed Pth century) the 20th century) ckmail) gsters in charge of

# **Query Expansion Example**

#### From TREC 2004 Robust Track Benchmark:

#### **Title:** International Organized Crime

**Description:** Identify organizations that participate in international criminal activity, the activity, and, if possible, collaborating organizations and the countries involved.

**Query** = {international[0.145|1.00],

~META[1.00|1.00][{gangdom[1.00|1.00], gangland[0.742|1.00], ''organ[0.213|1.00] & crime[0.312|1.00]'', camorra[0.254|1.00], maffia[0.318|1.00], mafia[0.154|1.00], ''sicilian[0.201|1.00] & mafia[0.154|1.00]'', ''black[0.066|1.00] & hand[0.053|1.00]'', mob[0.123|1.00], syndicate[0.093|1.00]}], organ[0.213|1.00], crime[0.312|1.00], collabor[0.415|0.20], columbian[0.686|0.20], cartel[0.466|0.20], ...}}

135530 sorted accesses in 11.073s.

#### **Results:**

- 1. Interpol Chief on Fight Against Narcotics
- 2. Economic Counterintelligence Tasks Viewed
- 3. Dresden Conference Views Growth of Organized Crime in Europe
- 4. Report on Drug, Weapons Seizures in Southwest Border Region
- 5. SWITZERLAND CALLED SOFT ON CRIME

### Statistics for Term-Term Similarity or Concept-Concept Relatedness

**Relatedness measures** *sim(c1, c2)* **based on WordNet-like thesaurus:** 

Wu-Palmer distance:|path(c1,lca(c1,c2))| + path(c2,lca(c1,c2))with lowest common ancestor lca(c1,c2) in DAG

Variants with edge weights based on edge type (hyponym, hypernym, ...)

**Relatedness measures** *sim(c1, c2)* **based on co-occurrences in corpus:** 

Dice coefficient: $2 | \{ \text{ docs with } c1 \} \cap \{ \text{ docs with } c2 \} |$  $| \{ \text{ docs with } c1 \} | + | \{ \text{ docs with } c2 \} |$ Jaccard coefficient: $| \{ \text{ docs with } c1 \} \cap \{ \text{ docs with } c2 \} |$  $| \{ \text{ docs with } c1 \} | + | \{ \text{ docs with } c2 \} | -| \{ \text{ docs with } c1 \text{ and } c2 \} |$ PMI (Pointwise<br/>Mutual Information): $\log \frac{\text{freq}(c1 \text{ and } c2)}{\text{freq}(c1) \cdot \text{freq}(c2)}$ Conditional probability:P[ doc has c1 | doc has c2 ]

### **Exploiting Query Logs for Query Expansion**

Given: user sessions of the form (q, D+) with clicked docs D+ (often only a single doc)

We are interested in the correlation between words w in a query and w' in a clicked-on document:

 $P[w'|w] := P[w' \in d \text{ for some } d \in D^+ | w \in q]$   $= \sum_{d \in D^+} P[w' \in d | d \in D^+] \cdot P[d \in D^+ | w \in q]$ Estimate
from query log:
relative frequency
of w' in d
relative frequency of d being clicked on
when w appears in query

Expand query by adding top m words w' in desc. order of  $\prod_{w \in q} P[w'|w]$ 

### Term-Term Similarity Estimation from Query-Click Logs

Use co-occurrences of

- term and term in **same query** (ordered terms)
- term in **query** and term in (title or URL of) **clicked doc**
- term in query without click and term in next query to compute maximum-likelihood estimator for multinomial distribution for ordered term pairs or n-grams and derive P[term u | term w] ~ freq[term u | term w]

Useful for

- Suggestions for alternative queries ("did you mean ...?")
- Suggestions for auto-completion
- Background statistics for geo-localization or user-personalization

### **12.4 Query Result Diversification**

True goal of search engine is to maximize *P[user clicks on at least one of the top-k results]* 

With ambiguity of query terms and uncertainty about user intention (examples: "apple farm", "mpi research", "darwin expedition", "Star Wars 7: The Force Awakens", "physics nobel prize", ...)
we need to diversify the top-10 for risk minimization (portfolio mix)

Given a query q, query results  $D=\{d_1, d_2, ...\}$ , similarity scores for results and the query  $sim(d_i,q)$ and pair-wise similarities among results  $sim(d_i,d_i)$ 

→ Select top-k results  $r_1, ..., r_k \in D$  such that  $\alpha \sum_{i=1..k} sim(r_i, q) - (1 - \alpha) \sum_{i \neq j} sim(r_i, r_j) = max!$ 

### **Alternative Models for Diversification**

Variant 1: Max-Min-Dispersion [Ravi, Rosenkrantz, Tayi 1994] determine results set R={  $r_1, ..., r_k$  } such that  $\alpha \min_{i=1..k} sim(r_i, q) - (1 - \alpha) \max_{\substack{i \neq j}} sim(r_i, r_j) = max!$ 

Variant 2: intention-modulated [Agrawal et al. 2009] assume that q may have m intentions t<sub>1</sub>..t<sub>m</sub> (trained on query-click logs, Wikipedia disambiguation pages, etc.): determine result set R with |R|=k such that

$$P[R \mid q] = \sum_{i=1}^{m} P[t_i \mid q] \cdot \left(1 - \prod_{r \in R} (1 - P[r \mid q, t_i])\right) = \max!$$
  
at least one r clicked  
given intention t<sub>i</sub> for q

More variants in the literature, most are NP-hard But many are **submodular** (have diminishing marginal returns) → **greedy algorithms** with approximation guarantees

### **Submodular Set Functions**

Given a set  $\Omega$ , a function f:  $2^{\Omega} \rightarrow \Re$  is **submodular** if for every X,  $Y \subseteq \Omega$  with  $X \subset Y$  and  $z \in \Omega - Y$ the following **diminishing-returns property** holds  $f(X \cup \{z\}) - f(X) \ge f(Y \cup \{z\}) - f(Y)$ 

Typical **optimization** problem aims to choose a subset  $X \subset \Omega$  that minimizes or maximizes f under cardinality constraints for X

- these problems are usually NP-hard but often have polynomial algorithms with very good approximation guarantees
- greedy algorithms often yield very good approximate solutions

### **Maximal Marginal Relevance (MMR): Greedy Reordering for Diversification**

[Carbonell/Goldstein 1998]

Compute a pool of top-m candidate results where m > k(e.g. m=1000 for k=10) Initialize  $S := \emptyset$ 

Choose results in descending order of marginal utility: repeat

 $S \coloneqq S \cup argmax_d (\alpha sim(d,q) - (1 - \alpha) \sum_{r \in S} sim(r,d))$ until |S|=k

### **Summary of Chapter 12**

• **document-ordered** posting lists:

QP based on scan and merge; can optimize order of lists and heuristically control memory for accumulators

• impact-ordered posting lists:

top-k search can be sublinear with Threshold Algorithm family

- additional algorithmic options and optimizations for phrase and proximity queries and for query expansion
- with **ambiguity** of query terms and **uncertainty** of user intention, query result **diversification** is crucial

### **Additional Literature for Chapter 12**

- J. Zobel, A. Moffat: Self-Indexing Inverted Files for Fast Text Retrieval, ACM TOIS 1996
- A. Broder, D. Carmel, M. Herscovici, A. Soffer, J. Zien: Efficient query evaluation using a two-level retrieval process, CIKM 2003
- R. Fagin, A. Lotem, M. Naor: Optimal Aggregation Algorithms for Middleware, Journal of Computer and System Sciences 2003
- C. Buckley, A. Lewit: Optimization of Inverted Vector Searches, SIGIR 1985
- I.F. Ilyas, G. Beskales, M.A. Soliman: A Survey of Top-k Query Processing Techniques in Relational Database Systems, ACM Comp. Surveys 40(4), 2008
- A. Marian, N. Bruno, L. Gravano: Evaluating Top-k Queries over Web-accessible Databases, ACM TODS 2004
- M. Theobald et al.: Top-k Query Processing with Probabilistic Guarantees, VLDB 2004
- H. Bast et al.: IO-Top-k: Index-access Optimized Top-k Query Processing. VLDB 2006
- H.E. Williams, J. Zobel, D. Bahle: Fast Phrase Querying with Combined Indexes, TOIS 2004
- A. Broschart, R. Schenkel: High-performance processing of text queries with tunable pruned term and term pair indexes. ACM TOIS 2012
- S. Büttcher, C. Clarke, B. Lushman: Term proximity scoring for ad-hoc retrieval on very large text collections. SIGIR 2006
- R. Schenkel et al.: Efficient Text Proximity Search, SPIRE 2007

### **Additional Literature for Chapter 12**

- G. Miller, C. Fellbaum: WordNet: An Electronic Lexical Database. MIT Press 1998
- B. Billerbeck, F. Scholer, H.E. Williams, J. Zobel: Query expansion using associated queries. CIKM 2003
- S. Liu, F. Liu, C.T. Yu, W. Meng: An effective approach to document retrieval via utilizing WordNet and recognizing phrases, SIGIR 2004
- M. Theobald, R. Schenkel, G. Weikum: Efficient and Self-Tuning Incremental Query Expansion for Top-k Query Processing, SIGIR 2005
- H. Bast, I. Weber: Type Less, Find More: Fast Autocompletion Search with a Succinct Index, SIGIR 2006
- Z. Bar-Yossef, M. Gurevich: Mining Search Engine Query Logs via Suggestion Sampling. PVLDB 2008
- Z. Bar-Yossef, N. Kraus: Context-sensitive query auto-completion. WWW 2011
- T. Joachims et al.: Evaluating the accuracy of implicit feedback from clicks and query reformulations in Web search, TOIS 2007
- SP.Chirita, C.Firan, W.Nejdl: Personalized Query Expansion for the Web. WWW 2007
- J. Carbonell, J. Goldstein: The Use of MMR: Diversity-based Reranking, SIGIR 1998
- R. Agrawal, S. Gollapudi, A. Halverson, S. Ieong: Diversifying search results. WSDM 2009
- S. Gollapudi, A. Sharma: An Axiomatic Approach for Result Diversification, WWW 2009
   IRDM WS 2015