

7 Top-k Queries on Web Sources and Structured Data

- 7.1 Top-k Queries over Autonomous Web Sources
- 7.2 Ranking of SQL Query Results

Winter Semester 2003/2004

Selected Topics in Web IR and Mining

7-1

7.1 Computational Model for Top-k Queries over Web Sources

Typical example:

Address = „2590 Broadway“ and Price = \$ 25 and Rating = 30 issued against mapquest.com, nytoday.com, zagat.com

Major complication:

some sources do not allow sorted access

Major opportunity:

sources can be accessed in parallel

→ extension/generalization of TA

distinguish S-sources, R-sources, SR-sources

Winter Semester 2003/2004

Selected Topics in Web IR and Mining

7-2

Source-Type-Aware TA

For each R-source $S_i \in S_{m+1} \dots S_{m+r}$ set $high_i := 1$

Scan SR- or S-sources $S_1 \dots S_m$

Choose SR- or S-source S_i for next sorted access

for object d retrieved from SR- or S-source L_i do {

$E(d) := E(d) \cup \{i\}$; $high_i := si(q,d)$;

$bestscore(d) := \text{aggr}\{x_1, \dots, x_m\}$ with $x_i := si(q,d)$ for $i \in E(d)$, $high_i$ for $i \notin E(d)$;

$worstscore(d) := \text{aggr}\{x_1, \dots, x_m\}$ with $x_i := si(q,d)$ for $i \in E(d)$, 0 for $i \notin E(d)$; }

Choose SR- or R-source S_i for next random access

for object d retrieved from SR- or R-source L_i do {

$E(d) := E(d) \cup \{i\}$;

$bestscore(d) := \text{aggr}\{x_1, \dots, x_m\}$ with $x_i := si(q,d)$ for $i \in E(d)$, $high_i$ for $i \notin E(d)$;

$worstscore(d) := \text{aggr}\{x_1, \dots, x_m\}$ with $x_i := si(q,d)$ for $i \in E(d)$, 0 for $i \notin E(d)$; }

current top-k := k docs with largest worstscore;

worstmin_k := minimum worstscore among current top-k;

Stop when $bestscore(d) \leq worstmin_k$;

Return current top-k;

in contrast to Fagin's TA, keep complete list of candidate objects

Winter Semester 2003/2004

Selected Topics in Web IR and Mining

7-3

Strategies for Choosing the Source for Next Access

for next sorted access:

Escore(S_i) := expected si value for next sorted access to S_i

(e.g.: $high_i$)

$rank(S_i) := w_i * \text{Escore}(S_i) / c_s(S_i)$ // w_i is weight of S_i in aggr

choose SR- or S-source with highest $rank(S_i)$

for next random access (probe):

Escore(S_i) := expected si value for next random access to S_i

(e.g.: $(high_i - low_i) / 2$)

$rank(S_i) := w_i * \text{Escore}(S_i) / c_r(S_i)$

choose SR- or S-source with highest $rank(S_i)$

Winter Semester 2003/2004

Selected Topics in Web IR and Mining

7-4

The Upper Strategy for Choosing Next Object and Source (Marian et al.: TODS 2004)

idea: eagerly prove that candidate objects cannot qualify for top-k

for next random access:

among all objects with $E(d) \neq \emptyset$ and $R(d) \neq \emptyset$

choose d^* with the highest $bestscore(d^*)$;

if $bestscore(d^*) < bestscore(d^{**})$ for object d^{**} with $E(d^{**}) = \emptyset$ then

perform sorted access next (i.e., don't probe d^*)

else {

$\Delta := bestscore(d^*) - worstmin_k$;

if $\Delta > 0$ then {

consider S_i as „redundant“ for d^* if for all $Y \subseteq R(d^*) - \{S_i\}$

$\sum_{j \in Y} w_j * high_j + w_i * high_i \geq \Delta \Rightarrow \sum_{j \in Y} w_j * high_j \geq \Delta$;

choose „non-redundant“ source with highest $rank(S_i)$ }

else choose source with lowest $c_r(S_i)$;

};

Alternative for early stage of algorithm:

could prioritize random access for objects in current top-k to improve $worstmin_k$

Winter Semester 2003/2004

Selected Topics in Web IR and Mining

7-5

The Parallel Strategy pUpper (Marian et al.: TODS 2004)

idea: consider up to $MPL(S_i)$ parallel probes to the same R-source S_i

choose objects to be probed based on

bestscore reduction and expected response time

for next random access:

probe-candidates := m objects d with $E(d) \neq \emptyset$ and $R(d) \neq \emptyset$

such that d is among the m highest values of $bestscore(d)$;

for each object d in probe-candidates do {

$\Delta := bestscore(d) - worstmin_k$;

if $\Delta > 0$ then {

choose subset $Y(d) \subseteq R(d)$ such that $\sum_{j \in Y} w_j * high_j \geq \Delta$

and expected response time

$\sum_{S_j \in Y(d)} (\{ \{ d^* | bestscore(d^*) > bestscore(d) \text{ and } Y(d) \cap Y(d^*) \neq \emptyset \} \}$

$* c_r(S_j) / MPL(S_j))$

is minimum };

};

enqueue probe(d) to queue(S_i) for all $S_i \in Y(d)$

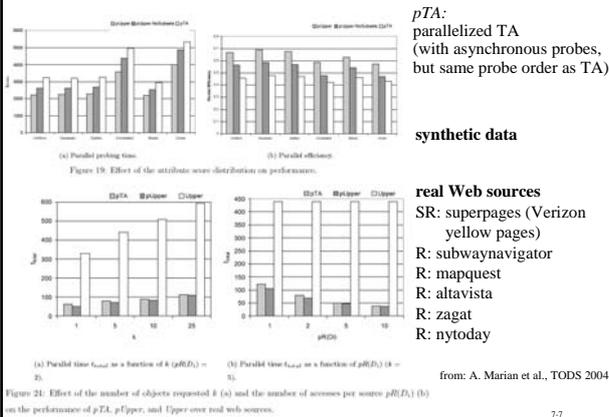
with expected response time as priority;

Winter Semester 2003/2004

Selected Topics in Web IR and Mining

7-6

Experimental Evaluation



7.2 Ranking of SQL Query Results

motivated by applications such as:
 realtor databases (e.g., homeadvisor.msn.com) with attributes
 City, Datebuilt, Price, Sqft, Bedrooms, Bathrooms,
 Deck, Fenced, Culdesac, Pool, Spa, ...

SQL queries such as

```
Select * From Homes Where City In (,Redmond', ,Bellevue') And
Bedrooms >= 3 And Bathrooms >= 2 And Price < 300000
```

often return too many answers or empty answers

Similar examples are car sales databases, customer support data, etc.

Straightforward idea:

apply Fagin's TA algorithm, treating indexed attributes as SR-sources and non-indexed attributes as R-sources, implemented in SQL

Less obvious issue:

similarity measures on numerical and categorical attributes

Winter Semester 2003/2004

Selected Topics in Web IR and Mining

7-8

IDF Based Similarities

for query condition $A_i=q_i$ with categorical attribute A_i
 the score of tuple t with $t.A_i=t_j$ is:

$$si(t_j, q_i) := idf(t_j) \text{ for } q_i=t_j, 0 \text{ else}$$

$$\text{with } idf(t_j) := \log(|R| / |\{t \in R \text{ and } t.A_i=t_j\}|)$$

the total score for a query over multiple categorical attributes then is:

$$s(t, q) := \sum_i si(t.A_i, q_i)$$

Example: q : Type=,Convertible' And Make=,Nissan'
 ranks non-Nissan convertibles higher than Nissan standard cars, 4WDs, etc.

for (practically continuous) numerical attributes A (e.g., price) define:

$$idf(t_j) = \log \left(|R| / \sum_{t \in R} \frac{1}{2} \left(\frac{t.A_i - t_j}{h} \right)^2 \right) \quad si(t_j, q_i) = e^{-\frac{1}{2} \left(\frac{q_i - t_j}{h} \right)^2} \cdot idf(t_j)$$

for conditions A_i In $(q_{i_1}, \dots, q_{i_k})$ set $si(t_j, q_i) := \max_{j=1..k} si(t_j, q_{i_j})$

Winter Semester 2003/2004

Selected Topics in Web IR and Mining

7-9

QF Based Similarities

idf can be misleading about importance

(e.g., old houses are infrequent, new houses are usually looked for)

Query frequency (qf), as reflected in the current workload,
 can indicate importance, too:

$$si(t_j, q_i) := qf(t_j) \text{ for } q_i=t_j, 0 \text{ else}$$

Workloads can also reveal correlations in user preferences
 (e.g., frequent queries with Make In (,Honda Accord', ,Toyota Camry')
 or City In (,Redmond', ,Bellevue', ,Kirkland'))

$W(t)$: set of previous queries with In predicate containing value t

For query value q and tuple with value t define:

$$si(q, t) := J(W(q), W(t)) * qf(t)$$

$$\text{with Jaccard coefficient } J(W(q), W(t)) = \frac{|W(q) \cap W(t)|}{|W(q) \cup W(t)|}$$

Winter Semester 2003/2004

Selected Topics in Web IR and Mining

7-10

Experimental Evaluation

4000 tuples about homes for sale in the Seattle Eastside area
 80 queries each with 2-5 attributes specified
 human relevance assessment by 5 colleagues

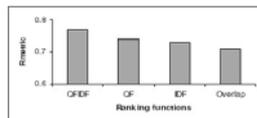


Figure 3: Quality of various ranking functions on Realtor database

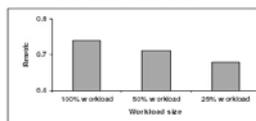


Figure 4: Ranking quality of QF Similarity on Realtor database as workload size varies

from: S. Agrawal et al., CIDR 2003

Winter Semester 2003/2004

Selected Topics in Web IR and Mining

7-11

Ranking quality metric:
 weighted precision R
 for top n ($n=10$)

$$R = \sum_{i=1..n} r_i / 2^{(i-1)/n}$$

with $r_i=1$ if rank is relevant,
 0 otherwise

Literature

- Amelie Marian, Nicolas Bruno, Luis Gravano: Evaluating Top-k Queries over Web-Accessible Databases, to appear in ACM TODS 2004
- Nicolas Bruno, Luis Gravano, Amelie Marian: Evaluating Top-k Queries over Web-Accessible Databases, ICDE Conf. 2002
- Ronald Fagin, Amnon Lotem, Moni Naor: Optimal Aggregation Algorithms for Middleware, Journal of Computer and System Science Vol. 66, 2003
- Sanjay Agrawal, Surajit Chaudhuri, Gautam Das, Aristides Gionis: Automated Ranking of Database Query Results, CIDR Conf. 2003
- Panayiotis Tsaparas, Themistoklis Palpanas, Yannis Kotidis, Nick Koudas, Divesh Srivastava: Ranked Join Indices, ICDE Conf. 2003

Winter Semester 2003/2004

Selected Topics in Web IR and Mining

7-12