

Context-Aware Top- k Processing Using Views

Silviu Maniu, Bogdan Cautis

University of Hong Kong & Univ. Paris-Sud / INRIA Saclay

CIKM 2013

Location-aware top-k retrieval

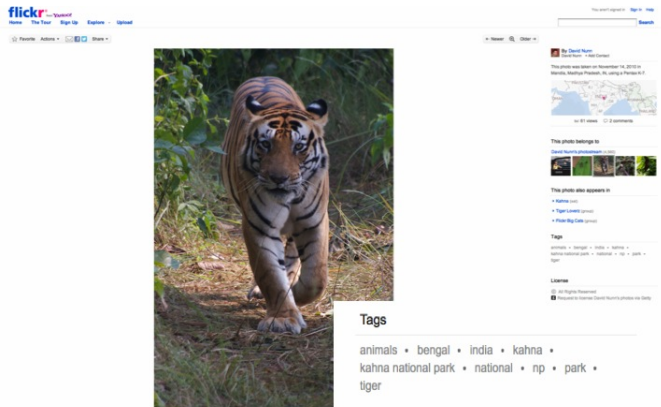
Users search for specific types of restaurants near a given location.

The screenshot shows a web interface for finding restaurants in Auvergne. At the top, there are navigation links for 'Auvergne', 'Créer un compte', and 'S'identifier'. Below this is a breadcrumb trail: 'Restaurant France > Restaurants Auvergne > restaurant Cuisine française Auvergne'. The main search area is titled 'Réserver un restaurant - Auvergne' and includes a search bar with 'Clermont Ferrand' and a dropdown for 'Cuisine française'. There are also fields for 'Jour' (24/10), 'Heure' (20h30), and 'Personne(s)' (2). A 'Trouver une table' button is on the right. Below the search area, there are filters for 'Top restaurants' (Les mieux notés, Les plus réservés, Par type de cuisine, Par ambiance, Par prix) and a '15 restaurants trouvés' section. A note indicates the search is near Clermont-Ferrand (63000) with a 20 km radius. A table lists the results with columns for Nom, Distance, Type, Prix, Avis, Promotions et Menus, and Dispo. The table contains three rows of restaurant data.

| Nom | Distance | Type | Prix | Avis | Promotions et Menus | Dispo |
|---|----------|------------------------------------|---------|--|-----------------------|-------|
| Emmanuel Hodence Clermont-Ferrand | 0,6 km | Gastronomique Cuisine française | € € € € | ★★★★★ 8,1/10 13 avis Gaut Millau 2012 Michelin 2012 | Menu "Saveurs" 78€ * | |
| Campanile Clermont Ferrand - Le Brezet Clermont-Ferrand | 3,0 km | Tendance Cuisine française | € € | ★★★★★ 8,7/10 11 avis | -40% sur la carte ! * | |
| La Table d'Isidore Clermont-Ferrand | 2,0 km | Gastronomique Cuisine française | € € € | ★★★★★ 8,2/10 33 avis Gaut Millau 2012 Michelin 2012 | Formule * | |

Social-aware top-k retrieval

In social tagging applications (Flickr, Delicious, Twitter), users search for photos/pages/items having certain tags.



The screenshot shows a Flickr photo page. The main image is a tiger walking through a forest. To the right of the image is a sidebar with metadata: the photo was taken on November 14, 2010, in Mahuli, Madhya Pradesh, India, using a Canon 40D. It has 81 views and 2 comments. Below this, it lists related photos and tags. The tags listed are: animals, bengal, india, kahna, kahna national park, national, np, park, tiger. At the bottom of the page, there is a 'Tags' section with the following text: animals • bengal • india • kahna • kahna national park • national • np • park • tiger.

flickr — Yahoo!
Home The Year Sign Up Explore Upload

Photo: Animals • [Share](#)

By David Nurn
David Nurn • 1400 Connects
This photo was taken on November 14, 2010 in Mahuli, Madhya Pradesh, IN, using a Canon 40D

81 views 2 comments

This photo belongs to
David Nurn's photostream (100)

This photo also appears in
• [Animals](#) (10)
• [Tiger Land](#) (100)
• [Flickr Big Cats](#) (100)

Tags
animals • bengal • india • kahna • kahna national park • national • np • park • tiger

License
© All Rights Reserved
Respected to licensee David Nurn's photostream via Getty

Tags

animals • bengal • india • kahna • kahna national park • national • np • park • tiger

Outline

Context-aware top-k retrieval

Uncertainty in views

View-based top-k processing

Refinements

Experiments

Context-aware top- k retrieval

- ▶ Collection of **objects** \mathcal{O} , **attributes** \mathcal{T} (e.g., keywords, tags)
- ▶ For a given **context parameter** \mathcal{C} , objects o are associated to certain attributes t , by a function $score(o, t \mid \mathcal{C})$
 - ▶ extended to a set of attributes by monotone aggregation (e.g., sum).

$$score(o, \{t_1, \dots, t_n\} \mid \mathcal{C}) = \sum(score(o, t_1 \mid \mathcal{C}), \dots, score(o, t_n \mid \mathcal{C}))$$

Problem (context-aware top- k retrieval)

Given a query $Q = \{t_1, \dots, t_n\} \subset \mathcal{T}$ and a context \mathcal{C} , retrieve the k objects $o \in \mathcal{O}$ having the highest values $score(o, Q \mid \mathcal{C})$.

Social-aware top-k retrieval

[Amer-Yahia et al. VLDB'08, Shenkel et al. SIGIR'08, Maniu et al. CIKM'13]

Top-k retrieval in social tagging applications:

- ▶ Collaborative tagging environment: objects (e.g., photos), users, attributes (tags), a relation $\text{Tagged}(\text{object}, \text{user}, \text{tag})$
- ▶ Social network: associates to pairs of users a social proximity value (σ) (e.g., $[0, 1]$ similarity in tagging)
- ▶ Social score model: a seeker-dependent score (for seeker s)

$$\text{score}(o, t \mid s) = \sum_{u \in \{v \mid \text{Tagged}(o, u, t)\}} \sigma(s, u)$$

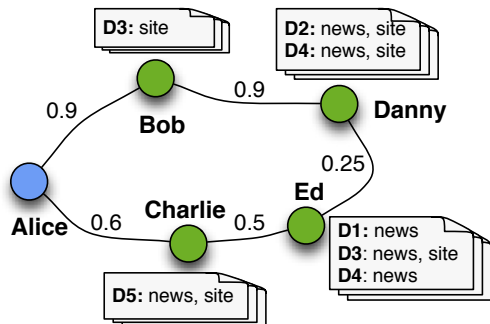
Problem (social-aware top-k retrieval)

Given a query $Q = \{t_1, \dots, t_n\}$ and a context (e.g., the seeker s), retrieve the k objects having the highest scores.

Social-aware top-k retrieval

Alice wants the top two documents for the query $\{news, site\}$

\rightsquigarrow a social-aware result: **D4, D2**



| <i>news</i> | | <i>site</i> | |
|-------------|------|-------------|------|
| doc | tf | doc | tf |
| D4 | 1.11 | D3 | 1.20 |
| D2 | 0.81 | D4 | 0.81 |
| D5 | 0.60 | D2 | 0.81 |
| D3 | 0.30 | D5 | 0.60 |
| D1 | 0.30 | D1 | 0.00 |

| user | prox. |
|---------|-------|
| Bob | 0.90 |
| Danny | 0.81 |
| Charlie | 0.60 |
| Ed | 0.30 |

Location-aware top-k retrieval

[Cong et al. VLDB'09, Christoforaki et al. CIKM'11, Cao et al. SIGMOD'11]

Top-k retrieval in spatial applications:

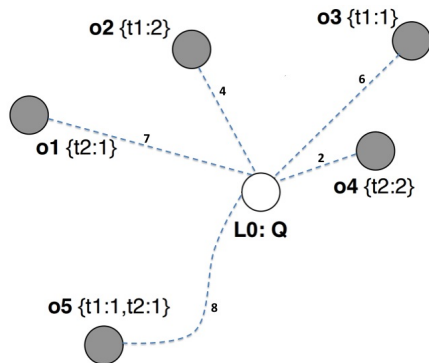
- ▶ Objects (e.g., documents) with attributes and geo-location.
- ▶ Spatial score model: combine textual and location relevance:

$$\text{score}(o, t \mid \text{loc}, \alpha) = \alpha \times \text{tf}(t, o) + (1 - \alpha) \times \text{dist}(o, \text{loc})$$

Problem (location-aware top- k retrieval)

Given a query $Q = \{t_1, \dots, t_n\}$, a context (e.g., location and α), retrieve the k objects having the highest scores.

Location-aware top-k retrieval



Top-2 query $Q = \{t1, t2\}$, $\alpha = 0.7$ at **L0** : **o4:0.92** and **o2:0.85**

Query answering using views

Context-aware retrieval is inherently difficult: joint exploration of the textual and “contextual” (e.g., spatial or social) space.

Our goal: improve **efficiency** by materialization, exploiting results of previous searches (views).

Each view has a context: its usefulness is proportional to distance w.r.t. the new context \rightsquigarrow score **uncertainty**, **approximate** top-k results.

Outline

Context-aware top-k retrieval

Uncertainty in views

View-based top-k processing

Refinements

Experiments

Context transposition

Focus on two applications: **location-aware search**, **social-aware search**

The **context** \mathcal{C}^V of a view V is a pair $(\mathcal{C}^V.l, \mathcal{C}^V.\alpha)$:

- ▶ **location** $\mathcal{C}^V.l$: geo-coordinates or seeker Id in a social network
- ▶ **contextual parameter** $\mathcal{C}^V.\alpha$: the weight of the context in scores

Context transposition

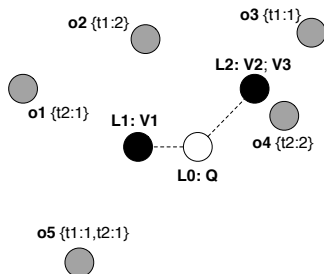
Focus on two applications: **location-aware search**, **social-aware search**

The **context** \mathcal{C}^V of a view V is a pair $(\mathcal{C}^V.l, \mathcal{C}^V.\alpha)$:

- ▶ **location** $\mathcal{C}^V.l$: geo-coordinates or seeker Id in a social network
- ▶ **contextual parameter** $\mathcal{C}^V.\alpha$: the weight of the context in scores

Transposition: adapt results for $(\mathcal{C}^V.l, \mathcal{C}^V.\alpha)$ to a new context $(\mathcal{C}.l, \mathcal{C}.\alpha)$

Context transposition yields uncertainty



$V1=(L1,\{t1,t2\})$

| o | sc |
|----|-------|
| o5 | 1.062 |
| o4 | 1.029 |
| o2 | 1.000 |

$V2=(L2,\{t1\})$

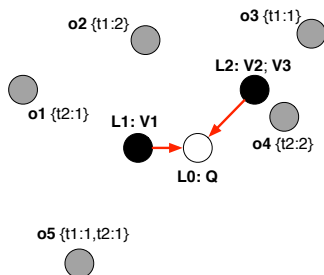
| o | sc |
|----|-------|
| o2 | 0.946 |
| o3 | 0.575 |
| o5 | 0.450 |
| o4 | 0.262 |

$V3=(L2,\{t2\})$

| o | sc |
|----|-------|
| o4 | 0.962 |
| o5 | 0.450 |
| o1 | 0.437 |
| o2 | 0.246 |

Top-2 query $Q=\{t1,t2\}$ at location $L0$

Context transposition yields uncertainty



$V1=(L1,\{t1,t2\})$

| o | sc |
|----|-------|
| o5 | 1.062 |
| o4 | 1.029 |
| o2 | 1.000 |

$V2=(L2,\{t1\})$

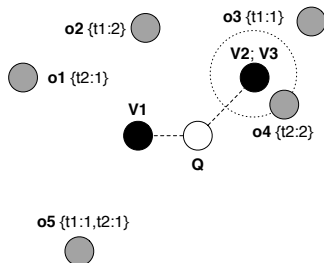
| o | sc |
|----|-------|
| o2 | 0.946 |
| o3 | 0.575 |
| o5 | 0.450 |
| o4 | 0.262 |

$V3=(L2,\{t2\})$

| o | sc |
|----|-------|
| o4 | 0.962 |
| o5 | 0.450 |
| o1 | 0.437 |
| o2 | 0.246 |

Top-2 query $Q=\{t1,t2\}$ at location $L0$

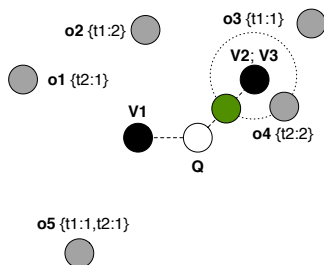
Context transposition yields uncertainty



| V1 =(L1,{t1,t2}) | | V2 =(L2,{t1}) | | V3 =(L2,{t2}) | |
|-------------------------|-------|----------------------|-------|----------------------|-------|
| o | sc | o | sc | o | sc |
| o5 | 1.062 | o2 | 0.946 | o4 | 0.962 |
| o4 | 1.029 | o3 | 0.575 | o5 | 0.450 |
| o2 | 1.000 | o5 | 0.450 | o1 | 0.437 |
| | | o4 | 0.262 | o2 | 0.246 |

distance of **o4** to **Q** unknown, but within [0.987, 1.037] interval

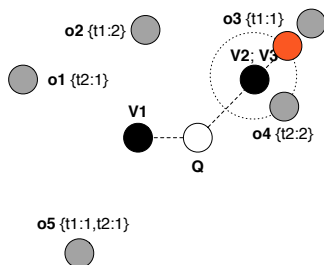
Context transposition yields uncertainty



| $V1=(L1, \{t1, t2\})$ | | $V2=(L2, \{t1\})$ | | $V3=(L2, \{t2\})$ | |
|-----------------------|-------|-------------------|-------|-------------------|-------|
| o | sc | o | sc | o | sc |
| o5 | 1.062 | o2 | 0.946 | o4 | 0.962 |
| o4 | 1.029 | o3 | 0.575 | o5 | 0.450 |
| o2 | 1.000 | o5 | 0.450 | o1 | 0.437 |
| | | o4 | 0.262 | o2 | 0.246 |

distance of $o4$ to Q unknown, but within $[0.987, 1.037]$ interval

Context transposition yields uncertainty

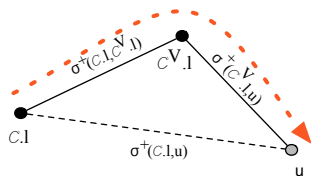


| $V1=(L1,\{t1,t2\})$ | | $V2=(L2,\{t1\})$ | | $V3=(L2,\{t2\})$ | |
|---------------------|-------|------------------|-------|------------------|-------|
| o | sc | o | sc | o | sc |
| o5 | 1.062 | o2 | 0.946 | o4 | 0.962 |
| o4 | 1.029 | o3 | 0.575 | o5 | 0.450 |
| o2 | 1.000 | o5 | 0.450 | o1 | 0.437 |
| | | o4 | 0.262 | o2 | 0.246 |

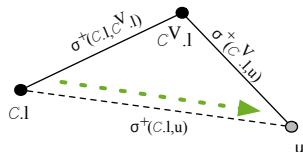
distance of **o4** to **Q** unknown, but within [0.987, 1.037] interval

Context transposition yields uncertainty

Reasoning based on shortest paths, i.e., the optimal is through:



- ▶ a path that has as prefix the $c.l \rightsquigarrow c^v.l$ path - **worstscore**



- ▶ other known paths - **bestscore**

Uncertain views

- ▶ For an input query Q , after context transposition (if necessary),
- ▶ A view V is composed of:
 1. a *definition* $def(V)$: a pair query-context (Q^V, \mathcal{C}^V)
 2. an *answer set* $ans(V)$: triples (o_i, wsc_i, bsc_i) , indicating that object o_i has a score in the range $[wsc_i, bsc_i]$

Outline

Context-aware top-k retrieval

Uncertainty in views

View-based top-k processing

Refinements

Experiments

Using the views for one object's bounds

Given a view set \mathcal{V} and a query Q sharing the same context, compute the tightest worst-score / best-score bounds for some object o .

Via a linear program:

$$\max \sum_{t_i \in Q} \mathbf{sc}(o, t_i \mid \mathcal{C}) \quad (1)$$

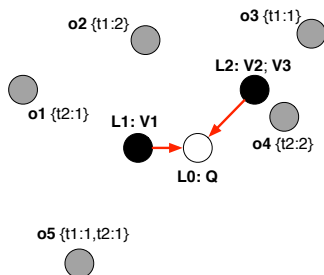
$$\min \sum_{t_j \in Q} \mathbf{sc}(o, t_j \mid \mathcal{C}) \quad (2)$$

$$wsc \leq \sum_{t_j \in Q^V} \mathbf{sc}(o, t_j \mid \mathcal{C}), \quad \forall V \in \mathcal{V} \text{ s.t. } (o, wsc, bsc) \in \mathit{ans}(V)$$

$$\sum_{t_j \in Q^V} \mathbf{sc}(o, t_j \mid \mathcal{C}) \leq bsc, \quad \forall V \in \mathcal{V} \text{ s.t. } (o, wsc, bsc) \in \mathit{ans}(V)$$

$$\mathbf{sc}(o, t_l \mid \mathcal{C}) \geq 0, \quad \forall t_l \in \mathcal{T}$$

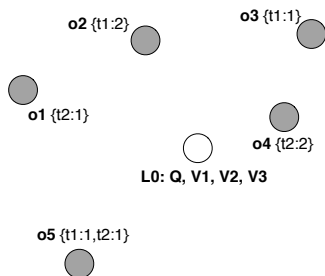
Before context transposition



| $V1=(L1,\{t1,t2\})$ | | $V2=(L2,\{t1\})$ | | $V3=(L2,\{t2\})$ | |
|---------------------|-------|------------------|-------|------------------|-------|
| o | sc | o | sc | o | sc |
| o5 | 1.062 | o2 | 0.946 | o4 | 0.962 |
| o4 | 1.029 | o3 | 0.575 | o5 | 0.450 |
| o2 | 1.000 | o5 | 0.450 | o1 | 0.437 |
| | | o4 | 0.262 | o2 | 0.246 |

Top-2 query $Q=\{t1,t2\}$ at location $L0$

After context transposition



| $V1=(L0, \{t1, t2\})$ | | | $V2=(L0, \{t1\})$ | | | $V3=(L0, \{t2\})$ | | |
|-----------------------|-------|-------|-------------------|-------|-------|-------------------|-------|-------|
| o | wsc | bsc | o | wsc | bsc | o | wsc | bsc |
| o5 | 0.957 | 1.167 | o2 | 0.871 | 1.000 | o4 | 0.987 | 1.037 |
| o4 | 0.924 | 1.134 | o3 | 0.500 | 0.650 | o5 | 0.500 | 0.525 |
| o2 | 0.895 | 1.105 | o5 | 0.500 | 0.525 | o1 | 0.362 | 0.512 |
| | | | o4 | 0.187 | 0.337 | o2 | 0.171 | 0.321 |

How can we use the views to compute the top-2 for Q?

Using views for one object: example

Top- k using views with **uncertain scores**:

LP formulation to compute tightest bounds - e.g., for **o5**:

$$\max \text{sc}(\mathbf{o5}, \mathbf{t1} \mid \mathcal{C}) + \text{sc}(\mathbf{o5}, \mathbf{t2} \mid \mathcal{C}) \quad (3)$$

$$\min \text{sc}(\mathbf{o5}, \mathbf{t1} \mid \mathcal{C}) + \text{sc}(\mathbf{o5}, \mathbf{t2} \mid \mathcal{C}) \quad (4)$$

$$0.957 \leq \text{sc}(\mathbf{o5}, \mathbf{t1} \mid \mathcal{C}) + \text{sc}(\mathbf{o5}, \mathbf{t2} \mid \mathcal{C}) \leq 1.167 \quad (V1)$$

$$0.500 \leq \text{sc}(\mathbf{o5}, \mathbf{t1} \mid \mathcal{C}) \leq 0.525 \quad (V2)$$

$$0.500 \leq \text{sc}(\mathbf{o5}, \mathbf{t2} \mid \mathcal{C}) \leq 0.525 \quad (V3)$$

Using views for one object: example

Top- k using views with **uncertain scores**:

LP formulation to compute tightest bounds - e.g., for **o5**:

$$\max \quad \text{sc}(\mathbf{o5}, \mathbf{t1} \mid \mathcal{C}) + \quad \text{sc}(\mathbf{o5}, \mathbf{t2} \mid \mathcal{C}) \quad (3)$$

$$\min \quad \text{sc}(\mathbf{o5}, \mathbf{t1} \mid \mathcal{C}) + \quad \text{sc}(\mathbf{o5}, \mathbf{t2} \mid \mathcal{C}) \quad (4)$$

$$0.957 \leq \text{sc}(\mathbf{o5}, \mathbf{t1} \mid \mathcal{C}) + \quad \text{sc}(\mathbf{o5}, \mathbf{t2} \mid \mathcal{C}) \leq 1.167 \quad (V1)$$

$$0.500 \leq \text{sc}(\mathbf{o5}, \mathbf{t1} \mid \mathcal{C}) \quad \leq 0.525 \quad (V2)$$

$$0.500 \leq \quad \text{sc}(\mathbf{o5}, \mathbf{t2} \mid \mathcal{C}) \leq 0.525 \quad (V3)$$

\rightsquigarrow score interval for **o5** between [1.000,1.050]

Our approach for top- k using views

Adapt the TA/NRA early-termination algorithms to the case of uncertain scores \rightsquigarrow the SR-TA and SR-NRA algorithms.

Our approach for top- k using views

Adapt the TA/NRA early-termination algorithms to the case of uncertain scores \rightsquigarrow the SR-TA and SR-NRA algorithms.

Plug the LPs in:

- ▶ the computation of worst-score/ best-score bounds,
- ▶ the computation of the termination threshold.

Most informative answer

In some cases, the exact top- k cannot be extracted with full confidence.

In our running example, at termination:

| Candidates | | |
|------------|-------|-------|
| obj | wsc | bsc |
| o4 | 1.174 | 1.134 |
| o2 | 1.042 | 1.105 |
| o5 | 1.000 | 1.050 |
| o3 | 0.500 | 0.971 |
| * | 0 | 0.849 |

Most informative answer

In some cases, the exact top- k cannot be extracted with full confidence.

In our running example, at termination:

| Candidates | | |
|------------|-------|-------|
| obj | wsc | bsc |
| o4 | 1.174 | 1.134 |
| o2 | 1.042 | 1.105 |
| o5 | 1.000 | 1.050 |
| o3 | 0.500 | 0.971 |
| * | 0 | 0.849 |

► one object guaranteed in the top-2: $G = \{o4\}$

Most informative answer

In some cases, the exact top- k **cannot be extracted with full confidence**.

In our running example, at termination:

| Candidates | | |
|------------|--------------|--------------|
| obj | wsc | bsc |
| o4 | 1.174 | 1.134 |
| o2 | 1.042 | 1.105 |
| o5 | 1.000 | 1.050 |
| o3 | 0.500 | 0.971 |
| * | 0 | 0.849 |

- ▶ one object **guaranteed** in the top-2: $G = \{o4\}$
- ▶ objects that **may be** in the top-2: $P = \{o2, o5\}$

Most informative answer

In some cases, the exact top- k **cannot be extracted with full confidence**.

In our running example, at termination:

| Candidates | | |
|------------|--------------|--------------|
| obj | wsc | bsc |
| o4 | 1.174 | 1.134 |
| o2 | 1.042 | 1.105 |
| o5 | 1.000 | 1.050 |
| o3 | 0.500 | 0.971 |
| * | 0 | 0.849 |

- ▶ one object **guaranteed** in the top-2: $G = \{o4\}$
- ▶ objects that **may be** in the top-2: $P = \{o2, o5\}$
- ▶ all other objects cannot be in the top-2

Top- k using uncertain views

Problem (Top- k retrieval using uncertain views)

Given a query $Q = \{t_1, \dots, t_n\} \subset \mathcal{T}$ and a context \mathcal{C} , given a set of views \mathcal{V} , retrieve from \mathcal{V} the *most informative answer* (G, P) , with

- ▶ $G \subset \mathcal{O}$ consisting of all *guaranteed objects*; i.e., in any data instance, they are in the top- k for Q and \mathcal{C} .
- ▶ and $P \subset \mathcal{O}$ consisting of all *possible objects* outside G ; i.e., there exist data instances where these are in the top- k for Q and \mathcal{C} .

Outline

Context-aware top-k retrieval

Uncertainty in views

View-based top-k processing

Refinements

Experiments

Beyond the most informative top- k answer

Estimating the **most likely top- k answer**:

Beyond the most informative top- k answer

Estimating the **most likely top- k answer**:

In the example: $P = \{o_2 \in [1.042, 1.105], o_5 \in [1.000, 1.050]\}$.

Beyond the most informative top- k answer

Estimating the **most likely top- k answer**:

In the example: $P = \{o_2 \in [1.042, 1.105], o_5 \in [1.000, 1.050]\}$.

If we assume a **uniform** distribution in the intervals:

$$\mathbf{P}[o_2 \geq o_5] = 0.989$$

$$\mathbf{P}[o_5 > o_2] = 0.011$$

Beyond the most informative top- k answer

Estimating the **most likely top- k answer**:

In the example: $P = \{o_2 \in [1.042, 1.105], o_5 \in [1.000, 1.050]\}$.

If we assume a **uniform** distribution in the intervals:

$$\mathbf{P}[o_2 \geq o_5] = 0.989$$

$$\mathbf{P}[o_5 > o_2] = 0.011$$

\implies the most likely top- k is $G \cup \{o_2\}$: $\mathbf{P}[\{o_4, o_2\}] = 0.989$

Beyond the most informative top- k answer

Estimating the **most likely top- k answer**:

In the example: $P = \{o_2 \in [1.042, 1.105], o_5 \in [1.000, 1.050]\}$.

If we assume a **uniform** distribution in the intervals:

$$\mathbf{P}[o_2 \geq o_5] = 0.989$$

$$\mathbf{P}[o_5 > o_2] = 0.011$$

\implies the most likely top- k is $G \cup \{o_2\}$: $\mathbf{P}[\{o_4, o_2\}] = 0.989$

Ways to evaluate:

- ▶ naive enumeration: good if $|P|$ is small,
- ▶ **sampling** or probabilistic top- k [Soliman et. al, VLDBJ10]

View selection

The P and G sets might be too expensive to compute, if the view set is very large, even using early-termination algorithms.

Solution: select few **most relevant views**, i.e., a subset $\tilde{\mathcal{V}} \subset \mathcal{V}$

- ▶ based on view definition, result statistics (see paper)

View selection

The P and G sets might be too expensive to compute, if the view set is very large, even using early-termination algorithms.

Solution: select few **most relevant views**, i.e., a subset $\tilde{\mathcal{V}} \subset \mathcal{V}$

- ▶ based on view definition, result statistics (see paper)
- ▶ trade-off between size of $\tilde{\mathcal{V}}$ and “quality” of the resulting (\tilde{G}, \tilde{P}) pair, in terms of distance to (G, P) :

$$\Delta = \binom{|\tilde{P}|}{k - |\tilde{G}|} - \binom{|P|}{k - |G|}$$

View selection

The P and G sets might be too expensive to compute, if the view set is very large, even using early-termination algorithms.

Solution: select few **most relevant views**, i.e., a subset $\tilde{\mathcal{V}} \subset \mathcal{V}$

- ▶ based on view definition, result statistics (see paper)
- ▶ trade-off between size of $\tilde{\mathcal{V}}$ and “quality” of the resulting (\tilde{G}, \tilde{P}) pair, in terms of distance to (G, P) :

$$\Delta = \binom{|\tilde{P}|}{k - |\tilde{G}|} - \binom{|P|}{k - |G|}$$

Final refinement: compute tightest bounds only for objects in $\tilde{G} \cup \tilde{P}$

Formal results

Instance optimality: For $A_1 \in \mathbf{A}$ and $A_2 \in \mathbf{A}$, write $A_1 \preceq A_2$ iff for all sets of views \mathcal{V} and all data instance \mathbf{D} , A_2 costs at least as much as A_1 .

Lemma

$\text{SR-NRA}^{sel} \not\preceq \text{SR-NRA}^{nosel} \not\preceq \text{SR-NRA}^{sel}$.

$\text{SR-TA}^{sel} \not\preceq \text{SR-TA}^{nosel} \not\preceq \text{SR-TA}^{sel}$.

Theorem

When we restrict the class of views to *pairwise disjoint views*:

- ▶ SR-TA^{sel} is instance optimal over \mathbf{A} .
- ▶ SR-NRA^{sel} is instance optimal over \mathbf{A} (when only sequential accesses are allowed).

Putting it all together

ProcessQueryUsingViews(\mathcal{V} , Q , \mathcal{C} , k)

Require: query Q , views \mathcal{V} , context \mathcal{C} , top k required

- 1: **for** $V \in \mathcal{V}$ **do**
- 2: transpose the context \mathcal{C}^V to \mathcal{C}
- 3: **end for**
- 4: $\tilde{\mathcal{V}} \leftarrow$ view selection on \mathcal{V} for Q
- 5: $(\tilde{G}, \tilde{P}) \leftarrow$ SR-TA($Q, k, \tilde{\mathcal{V}}$) or SR-NRA($Q, k, \tilde{\mathcal{V}}$)
- 6: $(G, P) \leftarrow$ REFINE(\tilde{G}, \tilde{P})
- 7: $E =$ ESTIMATE($P, k - |G|$)
- 8: **return** $G \cup E$

Outline

Context-aware top-k retrieval

Uncertainty in views

View-based top-k processing

Refinements

Experiments

Experiments: location-aware search

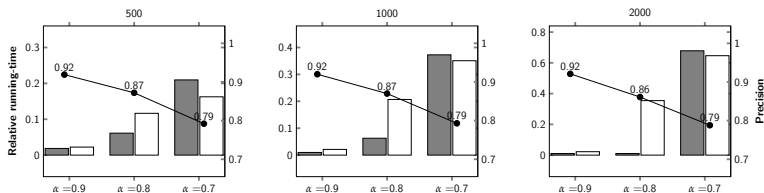


Figure : Performance and precision of SR-TA^{sel} versus exact early-termination algorithm (IR-TREE) (grey=top-10, white=top-20).

- ▶ PolyBot dataset: 6,115,264 objects and 1,876 attributes
- ▶ Views: 20 views of 2-term queries at 5 random locations, various list sizes
- ▶ Test: 10 queries at 5 locations and $\alpha \in \{0.7, 0.8, 0.9\}$

Experiments: social-aware search

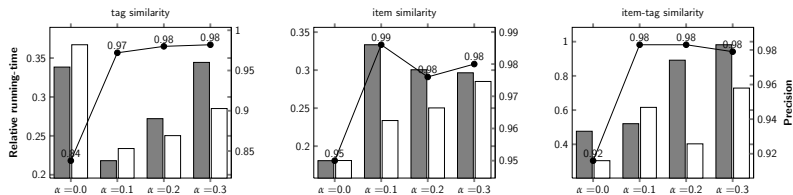


Figure : Social-aware search: performance and precision of SR-TA^{sel} versus CONTEXTMERGE (grey=top-10, white=top-20).

- ▶ Delicious data: 80000 users, 595811 objects, 198080 attributes
- ▶ Social network: 3 similarity networks (tag, item, item-tag)
- ▶ Views: 10 users each having 40 views of 1 and 2 tag queries
- ▶ Test: 10 3-tag queries for 5 seekers and $\alpha \in \{0, 0.1, 0.2, 0.3\}$

Summary

We formalize and study the problem of context-aware top-k processing based on (possibly uncertain) views.

- ▶ Context transposition, exemplified in two application scenarios
- ▶ New semantics based on views: most informative result
- ▶ Sound and complete adaptation of TA / NRA
- ▶ Probabilistic refinement: most likely top-k result
- ▶ Further efficiency: view selection
 - ▶ instance optimality under restrictions

Thank you.

Threshold algorithms: SR-TA

Adaptation of TA algorithm[Fagin01], SR-NRA similar.

Require: query Q , size k , views \mathcal{V} (after transposition)

- 1: $D = \emptyset$
- 2: **loop**
- 3: **for** each view $V \in \mathcal{V}$ in turn **do**
- 4: $(o_i, wsc_i, bsc_i) \leftarrow$ next tuple by sequential access in V
- 5: read by random-accesses all other lists $V' \in \mathcal{V}$ for tuples (o_j, wsc_j, bsc_j) s.t.
 $o_i = o_j$
- 6: $wsc \leftarrow$ solution to the MP in Eq. (1) for o_i
- 7: $bsc \leftarrow$ solution to the MP in Eq. (2) for o_i
- 8: add the tuple (o_i, wsc, bsc) to D
- 9: **end for**
- 10: $\tau \leftarrow$ maximal possible score of objects not encountered
- 11: $wsc_t \leftarrow$ lower-bound score of k th candidate in D
- 12: **if** $\tau \leq wsc_t$ **then**
- 13: **break**
- 14: **end if**
- 15: **end loop**
- 16: $(G, P) = \text{PARTITION}(D, k)$
- 17: **return** (G, P)

Threshold algorithms: PARTITION(D, k)

Require: candidate list D , parameter k

- 1: $G \leftarrow \emptyset$ the objects guaranteed to be in the top- k
- 2: $P \leftarrow \emptyset$ the objects that might enter the top- k
- 3: **for** each tuple $(o, bsc, wsc) \in D, o \neq *$ **do**
- 4: $x \leftarrow |\{(o', bsc', wsc') \in D \mid o' \neq o, bsc' > wsc\}|$
- 5: $wsc_t \leftarrow$ lower-bound score of k th candidate in D
- 6: **if** $x \leq k$ and for $(*, wsc_*, bsc_*) \in D, bsc_* \leq wsc$ **then**
- 7: add o to G
- 8: **else if** $bsc > wsc_t$ **then**
- 9: add o to P
- 10: **end if**
- 11: **end for**
- 12: **return** G, P

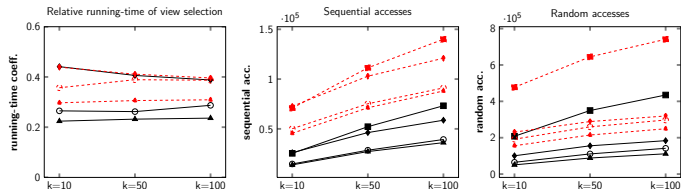
Experiments: context-agnostic setting

Input data:

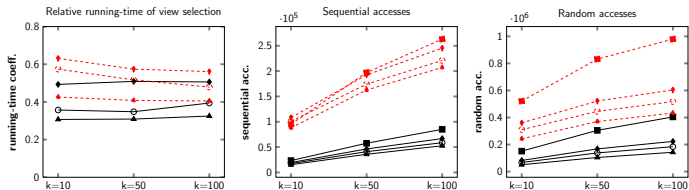
- ▶ synthetic: 100,000 objects and 10 attributes, scores in $[0,100]$
- ▶ views: all possible combinations of 2 and 3 attributes
- ▶ uncertain data: replace each score with a score range (Gaussian distribution, $\sigma \in \{5, 10\}$)

Test: 100 randomly-generated queries of 5 attributes

Experiments: context-agnostic setting



Uniform distribution



Exponential distribution



Experiments: context-agnostic setting

| Sel. + Dist. | Rel. running-time | | | Min. precision | | | P | | |
|--------------|-------------------|-------|-------|----------------|------|------|----|----|-----|
| | 10 | 50 | 100 | 10 | 50 | 100 | 10 | 50 | 100 |
| avg + uni | 0.576 | 0.676 | 0.712 | 0.57 | 0.69 | 0.72 | 10 | 36 | 64 |
| def + uni | 0.350 | 0.446 | 0.544 | 0.57 | 0.69 | 0.72 | 10 | 36 | 64 |
| max + uni | 0.296 | 0.395 | 0.446 | 0.57 | 0.69 | 0.72 | 10 | 36 | 64 |
| avg + exp | 0.732 | 1.128 | 1.287 | 0.60 | 0.63 | 0.64 | 10 | 46 | 86 |
| def + exp | 0.531 | 0.771 | 1.003 | 0.60 | 0.63 | 0.64 | 10 | 46 | 86 |
| max + exp | 0.456 | 0.684 | 0.827 | 0.60 | 0.63 | 0.64 | 10 | 46 | 86 |

Table : Comparison between SR-TA and TA (exact scores), for uniform and exponential distributions, for std 5.