Combining Fuzzy Information -Top-*k* Query Algorithms



Sanjay Kulhari

Outline

- Definitions Objects, Attributes and Scores
- Querying Fuzzy Data
- Top-k query algorithms
 - Naïve Algorithm
 - Fagin's Algorithm (FA)
 - Threshold Algorithm (TA)
 - No Random Access Algorithm (NRA)
- Comparing top-k query algorithms
- References



Objects, Attributes and Scores

- Each object X_i has m scores (r_{i1}, r_{i2}, ..., r_{im}), one for each of m attributes.
- Objects are listed, for each attribute sorted by score.
- Each object is assigned an overall score by combining the attribute score using aggregate function or combining rule.
- Aim: Determine *k* objects with the highest overall score.

	R_1	R_2	R_{3}	
X ₁	1	0.3	0.2	
X ₂	0.8	0.8	0	
X ₃	0.5	0.7	0.6	
X ₄	0.3	0.2	0.8	
X ₅	0.1	0.1	0.1	

	R ₁
X ₁	1
X ₂	0.8
X ₃	0.5
X ₄	0.3
X ₅	0.1

	R_2		R ₃
, 2	0.8	X ₄	0.8
, •3	0.7	X ₃	0.6
, 1	0.3	X ₁	0.2
, •4	0.2	X ₅	0.1
· •5	0.1	X_2	0





Querying Fuzzy Data - Example

Given the following relational structure



Query: Select top-2 for the **sum** aggregate function Monotonicity property: An aggregation function *t* is monotone if $t(x_1, ..., x_m) \le t(x'_1, ..., x'_m)$ whenever $x_i \le x'_i$ for every *i*.

Naïve Algorithm



	R_1	
X ₁	1	
X_2	0.8	
X ₃	0.5	
X ₄	0.3	
X_5	0.1	

	R_2	
X ₂	0.8	
X ₃	0.7	
X ₁	0.3	
X ₄	0.2	
X ₅	0.1	

	R_{3}
X ₄	0.8
X ₃	0.6
X ₁	0.2
X ₅	0.1
X ₂	0





	R_1
X ₁	1
X ₂	0.8
X ₃	0.5
X ₄	0.3
X ₅	0.1

	R_2
X ₂	0.8
X ₃	0.7
X ₁	0.3
X ₄	0.2
X ₅	0.1

	R_3
X ₄	0.8
X ₃	0.6
X ₁	0.2
X ₅	0.1
X ₂	0







	R_1
X ₁	1
X ₂	0.8
X ₃	0.5
X ₄	0.3
X ₅	0.1

	R ₂
X ₂	0.8
X ₃	0.7
X ₁	0.3
X ₄	0.2
X ₅	0.1

	R_3
X ₄	0.8
X ₃	0.6
X ₁	0.2
X ₅	0.1
X_2	0







	R_1
X ₁	1
X ₂	0.8
X ₃	0.5
X ₄	0.3
X ₅	0.1

	R_2
X ₂	0.8
X ₃	0.7
X ₁	0.3
X ₄	0.2
X ₅	0.1

	R_3
X ₄	0.8
X ₃	0.6
X ₁	0.2
X ₅	0.1
X_2	0

X ₁	1.5
X ₂	1.6
X ₃	1.8





	R_1
X ₁	1
X ₂	0.8
X ₃	0.5
X ₄	0.3
X ₅	0.1

	R_2
X ₂	0.8
X ₃	0.7
X ₁	0.3
X ₄	0.2
X ₅	0.1

	R ₃
X_4	0.8
X ₃	0.6
X ₁	0.2
X ₅	0.1
X ₂	0

X ₁	1.5
X ₂	1.6
X ₃	1.8
X ₄	1.3

Naïve Algorithm



	R ₁	
X ₁	1	
X ₂	0.8	
X ₃	0.5	
X ₄	0.3	
X ₅	0.1	

	R_2
X ₂	0.8
X ₃	0.7
X ₁	0.3
X ₄	0.2
X ₅	0.1

	R_3
X ₄	0.8
X ₃	0.6
X ₁	0.2
X ₅	0.1
X ₂	0

X ₁	1.5
X ₂	1.6
X ₃	1.8
X ₄	1.3
X ₅	0.3

Naïve Algorithm



2. Return k objects with the highest overall score.





	R_1	
X ₁	1	
X ₂	0.8	
X ₃	0.5	
X ₄	0.3	
X ₅	0.1	

	R_2
X ₂	0.8
X_3	0.7
X ₁	0.3
X ₄	0.2
X ₅	0.1

	R_3
X ₄	0.8
X ₃	0.6
X ₁	0.2
X ₅	0.1
X ₂	0



	R ₁		R_2		R_3
X ₁	1	X ₂	0.8	X ₄	0.8
X ₂	0.8	X ₃	0.7	X ₃	0.6
X ₃	0.5	X ₁	0.3	X ₁	0.2
X_4	0.3	X ₄	0.2	X ₅	0.1
X_5	0.1	X ₅	0.1	X ₂	0



	R ₁		R_2		R_3
X ₁	1	X ₂	0.8	X ₄	0.8
X ₂	0.8	X ₃	0.7	X ₃	0.6
X ₃	0.5	X ₁	0.3	X ₁	0.2
X_4	0.3	X ₄	0.2	X_5	0.1
X_5	0.1	X ₅	0.1	X ₂	0



	R ₁		R_2		R_3
X ₁	1	X ₂	0.8	X ₄	0.8
X ₂	0.8	X ₃	0.7	X ₃	0.6
X ₃	0.5	X ₁	0.3	X ₁	0.2
X_4	0.3	X ₄	0.2	X ₅	0.1
X_5	0.1	X_5	0.1	X ₂	0





Since k = 2, and X₁ and X₃ have been seen in all the 3 lists



2. Perform random accesses to obtain the scores of all seen objects

	R ₁		R_2		R_3
X ₁	1	X ₂	8.0	X ₄	0.8
X ₂	0.8	X ₃	0.7	X ₃	0.6
X ₃	0.5	X ₁	0.3	X ₁	0.2
X ₄	0.3	X ₄	0.2	X ₅	0.1
X_5	0.1	X ₅	0.1	X_2	0



3. Compute score for all objects and return the top-k



1. Access the elements sequentially

	R ₁
X ₁	1
X ₂	0.8
X ₃	0.5
X ₄	0.3
X_5	0.1

	R_2
X ₂	0.8
X ₃	0.7
X ₁	0.3
X ₄	0.2
X ₅	0.1

	R_{3}
X ₄	0.8
X ₃	0.6
X ₁	0.2
X ₅	0.1
X ₂	0





1. At each sequential access

(a) Set the threshold *t* to be the aggregate of the scores seen in this access.

	R ₁		R ₂		R ₃	
X ₁	1	X ₂	0.8	X ₄	8.0	t = 2.6
X ₂	0.8	X ₃	0.7	X ₃	0.6	
X ₃	0.5	X ₁	0.3	X ₁	0.2	
X_4	0.3	X ₄	0.2	X ₅	0.1	
X_5	0.1	X ₅	0.1	X_2	0	



1. At each sequential access

(b) Do random accesses and compute the scores of the seen objects.



1. At each sequential access

(c) Maintain a list of top-k objects seen so far







1. At each sequential access

(d) Stop, when the scores of the top-k are greater or equal to the threshold.

	R ₁		R ₂		R ₃			
X ₁	1	X ₂	0.8	X ₄	0.8			
X ₂	0.8	X ₃	0.7	X ₃	0.6	t = 2.1	X ₃	1.8
X ₃	0.5	X ₁	0.3	X ₁	0.2		X_2	1.6
X_4	0.3	X ₄	0.2	X ₅	0.1		•	l
X ₅	0.1	X_5	0.1	X ₂	0			



1. At each sequential access

(d) Stop, when the scores of the top-k are greater or equal to the threshold.

	R ₁		R ₂		R_{3}			
X ₁	1	X ₂	0.8	X ₄	0.8			
X ₂	0.8	X ₃	0.7	X ₃	0.6			
X ₃	0.5	X ₁	0.3	X ₁	0.2	t = 1	X ₃	1.8
X ₄	0.3	X ₄	0.2	X ₅	0.1		X ₂	1.6
X ₅	0.1	X ₅	0.1	X ₂	0			

2. Return the top-k seen so far





1. Access sequentially all lists in parallel until there are k objects for which the lower bound is higher than the upper bound of all other objects.

	R ₁		R ₂		R_{3}
X ₁	1	X ₂	0.8	X ₄	0.8
X ₂	0.8	X ₃	0.7	X ₃	0.6
X_3	0.5	X ₁	0.3	X ₁	0.2
X ₄	0.3	X ₄	0.2	X ₅	0.1
X_5	0.1	X ₅	0.1	X ₂	0

	LB	UB
X ₁	1	2.6
X ₂	.8	2.6
X ₄	.8	2.6



1. Access sequentially all lists in parallel until there are k objects for which the lower bound is higher than the upper bound of all other objects.

	R ₁		R_2		R ₃
X ₁	1	X ₂	0.8	X ₄	0.8
X ₂	0.8	X ₃	0.7	X ₃	0.6
X ₃	0.5	X ₁	0.3	X ₁	0.2
X_4	0.3	X ₄	0.2	X ₅	0.1
X_5	0.1	X ₅	0.1	X ₂	0

	LB	UB
X ₂	1.6	2.2
X ₃	1.3	2.1
X ₁	1	2.3
X ₄	0.8	2.3



1. Access sequentially all lists in parallel until there are k objects for which the lower bound is higher than the upper bound of all other objects.

	R ₁		R ₂		R_3
X ₁	1	X ₂	0.8	X ₄	0.8
X ₂	0.8	X ₃	0.7	X ₃	0.6
X ₃	0.5	X ₁	0.3	X ₁	0.2
X ₄	0.3	X ₄	0.2	X ₅	0.1
X ₅	0.1	X ₅	0.1	X ₂	0

	LB	UB
X ₃	1.8	1.8
X ₂	1.6	1.8
X ₁	1.5	1.5
X ₄	0.8	1.6





2. Return top-k objects for which the lower bound is higher than the upper bound of all other objects.





Comparing top-k query algorithms

- Naïve algorithm
 - Buffer space required is equal to the number of database objects.
 - Each entry is looked in the *m* sorted lists. The cost is linear in database size.
 - Not efficient for a large database.
- Fagin's algorithm (FA)
 - Large buffer space is required.
 - Random access is done at the end to get the missing scores.
 - Cost optimal under certain aggregate functions.
- Threshold algorithm (TA)
 - Buffer space required is bounded by *k*.
 - Score of an object not seen during algorithm execution is less than the threshold due to monotonicity property of aggregate function.
 - Less object access is required compared to FA because when k common objects have been seen in FA, their scores are higher or equal to threshold in TA.
 - May perform more random access than FA because in FA random access is done at the end only for the missing scores.
- No Random Access (NRA) algorithm
 - Only sorted access is performed.
 - May not report the exact object scores, since it uses bounds to determine top *k*.

References



- Combining Fuzzy Information: An Overview. Ronald Fagin.
- A Survey of Top-*k* Query Processing Techniques in Relational Database Systems. *Ilyas, Beskales and Soliman.*
- 'Web Information Search' lecture notes. Prof. Leonardi (http://www.dis.uniroma1.it/~leon/didattica/webir/)