

# Crowdsourcing for Top-K Query Processing over Uncertain Data (Extended abstract)

Eleonora Ciceri, Piero Fraternali, Davide Martinenghi and Marco Tagliasacchi  
Dipartimento di Elettronica, Informazione e Bioingegneria  
Politecnico di Milano  
Milan, Italy  
Email: first.last@polimi.it

## I. INTRODUCTION

Both social media and sensing infrastructures are producing an unprecedented mass of data characterized by their uncertain nature, due to either the noise inherent in sensors or the imprecision of human contributions. Therefore query processing over uncertain data has become an active research field. In the well-known class of applications commonly referred to as “top- $K$  queries”, the objective is to find the best  $K$  objects matching the user’s information need, formulated as a scoring function over the objects’ attribute values. If both the data and the scoring function are deterministic, the best  $K$  objects can be univocally determined and totally ordered so as to produce a single ranked result set (as long as ties are broken by some deterministic rule). However, in application scenarios involving uncertain data and fuzzy information needs, this does not hold: when either the attribute values or the scoring function are nondeterministic, there may be no consensus on a single ordering, but rather a space of possible orderings. To determine the correct ordering, one needs to acquire additional information so as to reduce the amount of uncertainty associated with the queried data and consequently the number of orderings in such a space. An emerging trend in data processing is *crowdsourcing*, defined as the systematic engagement of humans in the resolution of tasks through online distributed work. Our approach combines human and automatic computation in order to solve complex problems: when data ambiguity can be resolved by human judgment, crowdsourcing becomes a viable tool for converging towards a unique or at least less uncertain query result. The goal of this paper is to define and compare task selection policies for uncertainty reduction via crowdsourcing, with emphasis on the case of top- $K$  queries.

**Problem formulation:** Given a data set with uncertain values and an allowed budget of possible questions, determine the set of questions (to be posed to a crowd) that minimizes the expected residual uncertainty of the result, possibly leading to a unique ordering of the top  $K$  results.

## II. BACKGROUND

We consider the problem of answering a top- $K$  query over a relational database table  $T$  containing  $N$  tuples. The relevance of a tuple to the query is modeled as a score. Let  $t_i \in T$  be a tuple in the database, and  $s(t_i)$  be the score of tuple  $t_i$ , computed by applying a *scoring function* over  $t_i$ ’s attribute values. When the score  $s(t_i)$  is known for each tuple  $t_i$ , the tuples in  $T$  can be totally ordered in

descending order of  $s(t_i)$  by breaking ties deterministically. Instead, if the score  $s(t_i)$  is uncertain and thus modeled as a random variable (with a probability density function (pdf)  $f_i$ ), such an uncertainty induces a partial order over the tuples. Indeed, when the pdf’s of two tuples overlap, their relative order is undefined. Therefore, we define the *space of possible orderings* as the set of all the total orderings compatible with the given score probability functions. This space can be represented by means of a *tree of possible orderings*  $\mathcal{T}$  (henceforth: TPO) [4], in which each node (except the root) represents a tuple  $t_i$ , and an edge from  $t_i$  to  $t_j$  indicates that  $t_i$  is ranked higher than  $t_j$  (denoted  $t_i \prec t_j$ ). Each path  $\omega = t_1 \prec t_2 \prec \dots \prec t_N$  represents a possible ordering of the underlying set of tuples  $T$ , and is associated with a probability  $\Pr(\omega)$  [2]. Since processing a top- $K$  query over uncertain data requires computing the orderings of the first  $K$  tuples, in order to answer such a query it suffices to build the sub-tree  $\mathcal{T}_K$  of possible orderings up to depth  $K$ .

Given a TPO  $\mathcal{T}_K$ , we propose four measures to quantify its level of uncertainty. These measures are based on the idea that the larger the number of orderings in  $\mathcal{T}_K$  and the more similar their probabilities, the higher its uncertainty.

- **Entropy.**  $U_H(\mathcal{T}_K)$  measures  $\mathcal{T}_K$ ’s uncertainty via Shannon’s entropy, based *only* on the probabilities of its leaves.
- **Weighted entropy.**  $U_H(\mathcal{T}_K)$  is a weighted combination of entropy values at the first  $K$  levels of the TPO.
- **ORA.**  $U_{ORA}(\mathcal{T}_K)$  is based on the idea of comparing all the orderings in  $\mathcal{T}_K$  with the Optimal Rank Aggregation (ORA) [3], which is a sort of median ordering in  $\mathcal{T}_K$ .
- **MPO.**  $U_{MPO}(\mathcal{T}_K)$  refers to another representative ordering, i.e., the Most Probable Ordering (MPO) [3].

## III. HUMANS FIGHTING UNCERTAINTY

We consider crowd tasks expressed as questions of the form  $q = t_i \succ t_j$ , which compare  $t_i$  and  $t_j$  to determine which one ranks higher. Given a crowd worker’s answer  $ans(q)$  (i.e., either  $t_i \prec t_j$  or  $t_i \not\prec t_j$ ), we can prune from  $\mathcal{T}_K$  all the paths disagreeing with the answer. The goal of *Uncertainty Reduction* (UR) is to find a sequence of questions the answers to which allow pruning the TPO  $\mathcal{T}_K$  so that a single ordering (the “real” ordering  $\omega_r$ ) remains. A UR algorithm is *optimal* if the sequence it finds is always minimal.

*Theorem 3.1:* No deterministic UR algorithm is optimal.

We consider two practical classes of algorithms: *i)* offline algorithms, which determine the questions a priori, before ob-

taining any answer from the crowd, and *ii*) online algorithms, whereby the questions are determined incrementally as the answers to previous questions arrive. These classes reflect two common situations in crowdsourcing markets: one where the interaction is limited to the publication of a batch of tasks, which is evaluated for acceptance as a whole; and one where the employer can inspect the outcome of crowd work as it becomes available and incrementally publish further tasks.

Since optimality is unattainable, we simply aim at maximizing the *expected* uncertainty reduction after pruning  $\mathcal{T}_K$ . The set  $\mathcal{Q}_K$  of relevant questions to ask consists of those that compare tuples with an uncertain relative ordering, i.e., whose pdf's overlap. Let  $B$  denote the maximum number of questions (budget) that can be asked to the crowd workers. Our goal is then to select the best sequence of questions  $\mathcal{Q}^* = \langle q_1^*, \dots, q_B^* \rangle \subseteq \mathcal{Q}_K$  that causes the largest amount of expected uncertainty reduction (with consequent reduction of the number of orderings in  $\mathcal{T}_K$ ). An offline UR algorithm that always does so is said to be *offline-optimal*.

#### A. Offline question selection strategies

We present an offline-optimal one, and two sub-optimal but faster alternatives.

**Best-first search offline algorithm (A\*-off).** This algorithm adapts the well-known A\* search algorithm to explore the space of question sets.

*Theorem 3.2:* A\*-off is offline-optimal.

**Top-B offline algorithm (TB-off).** For each question  $q \in \mathcal{Q}_K$ , we compute the expected residual uncertainty  $\mathcal{R}_q(\mathcal{T}_K)$  (i.e., the uncertainty of  $\mathcal{T}_K$  after  $q$  is asked). Then,  $\mathcal{Q}^*$  is defined as the set of  $B$  questions with the highest  $\mathcal{R}_q(\mathcal{T}_K)$ .

**Conditional offline algorithm (C-off).** This method iteratively selects one question at a time based on the previous selections. Let  $\{q_1^*, \dots, q_i^*\}$  be the first  $i$  selected questions ( $\emptyset$  when  $i = 0$ ). The  $(i+1)$ -th question  $q_{i+1}^*$  is selected by C-off from  $\mathcal{Q}_K \setminus \{q_1^*, \dots, q_i^*\}$  so as to minimize  $\mathcal{R}_{\langle q_1^*, \dots, q_i^*, q_{i+1}^* \rangle}(\mathcal{T}_K)$ , i.e., the residual uncertainty conditioned by the choice of the previously selected questions  $q_1^*, \dots, q_i^*$ . The final output is thus  $\mathcal{Q}^* = \{q_1^*, \dots, q_B^*\}$ .

#### B. Online question selection strategies

An online algorithm can determine the  $i$ -th question based on the answers to all the previously asked  $i - 1$  questions.

**Best-first search online algorithm (A\*-on).** An online UR algorithm that iteratively applies A\*-off  $B$  times.

**Top-1 online algorithm (T1-on).** The T1-on algorithm builds the sequence of questions  $\mathcal{Q}^*$  iteratively: at each iteration, the algorithm selects the question that minimizes the expected residual uncertainty with budget  $B = 1$ , appends it to  $\mathcal{Q}^*$  and asks it to the crowd. The TPO  $\mathcal{T}_K$  is then updated to the sub-tree that agrees with the received answer. Early termination may occur if all uncertainty is removed with  $|\mathcal{Q}^*| < B$ .

#### C. Handling noisy workers

In a crowdsourcing scenario, the collected answers might be noisy. When a crowd worker's accuracy (i.e., the probability

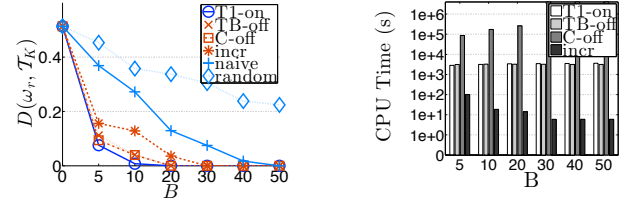


Fig. 1. Performance of faster algorithms as budget  $B$  varies.

that his/her answer is correct) is less than 1, no pruning of  $\mathcal{T}_K$  takes place, but the probabilities of the possible orderings are appropriately adjusted so as to reflect the collected answers.

#### D. Incremental algorithm

The *incr* algorithm, builds the TPO  $\mathcal{T}_K$  incrementally, one level at a time, by alternating tree construction with a round of  $n$  questions and tree pruning, thereby reducing the time needed to materialize trees containing a large number of orderings. The number  $n$  of questions posed at each round is between 1 and  $B$ , thus *incr* can be considered a hybrid between an online and an offline algorithm. Each TPO  $\mathcal{T}_k$ ,  $1 \leq k \leq K$ , is built by adding one level to  $\mathcal{T}_{k-1}$ . We only build new levels if there are not enough questions to ask. Then, we select the best questions, pose them to the crowd, collect the answers and apply the pruning accordingly, until either the budget  $B$  is exhausted or the TPO is entirely built.

## IV. EXPERIMENTAL EVALUATION

The proposed algorithms have been evaluated experimentally against baselines that select questions either randomly or focusing on tuples with an ambiguous order: *i*) the Random algorithm returns a sequence of  $B$  questions chosen at random among all possible tuple comparisons in  $\mathcal{T}_K$ ; *ii*) the Naive algorithm avoids irrelevant questions by returning a sequence of  $B$  questions chosen randomly from  $\mathcal{Q}_K$ . We observed that measures of uncertainty that take into account the structure of the tree in addition to ordering probabilities (i.e.,  $U_{MPO}$ ,  $U_{H_w}$ ,  $U_{ORA}$ ) perform better than state-of-the-art measures (i.e.,  $U_H$ ). The experiments show that the T1-on and C-off algorithms offer a good tradeoff between performance (Figure 1(a)) and costs (Figure 1(b)), with significant reductions of the number of questions wrt. the baselines, and nearly as good as with the A\*-based algorithms, but at a fraction of the cost. The proposed algorithms have been shown to work also with non-uniform tuple score distributions and with noisy crowds [1]. Much lower CPU times are possible with the *incr* algorithm, with slightly lower quality (which makes *incr* suited for large, highly uncertain datasets).

## REFERENCES

- [1] E. Ciceri et al. Crowdsourcing for top-k query processing over uncertain data. *IEEE Trans. Knowl. Data Eng.*, 28(1):41–53, 2016.
- [2] J. Li and A. Deshpande. Ranking continuous probabilistic datasets. *PVLDB*, 3(1-2):638–649, 2010.
- [3] M. Soliman et al. Ranking with uncertain scoring functions: semantics and sensitivity measures. In *SIGMOD '11*, pages 805–816, 2011.
- [4] M. Soliman and I. Ilyas. Ranking with uncertain scores. In *ICDE '09*, pages 317–328, 2009.