

Answering Planning Queries with the Crowd

Haim Kaplan Ilia Lotosh Tova Milo Slava Novgorodov

*School of Computer Science
Tel-Aviv University*

{haimk, ilialoto, milo, slavanov}@post.tau.ac.il

ABSTRACT

Recent research has shown that crowd sourcing can be used effectively to solve problems that are difficult for computers, e.g., optical character recognition and identification of the structural configuration of natural proteins. In this paper we propose to use the power of the crowd to address yet another difficult problem that frequently occurs in a daily life - answering planning queries whose output is a sequence of objects/actions, when the goal, i.e., the notion of “best output”, is hard to formalize. For example, planning the sequence of places/attractions to visit in the course of a vacation, where the goal is to enjoy the resulting vacation the most, or planning the sequence of courses to take in an academic schedule planning, where the goal is to obtain solid knowledge of a given subject domain. Such goals may be easily understandable by humans, but hard or even impossible to formalize for a computer.

We present a novel algorithm for efficiently harnessing the crowd to assist in answering such planning queries. The algorithm builds the desired plans incrementally, choosing at each step the ‘best’ questions so that the overall number of questions that need to be asked is minimized. We prove the algorithm to be optimal within its class and demonstrate experimentally its effectiveness and efficiency.

1. INTRODUCTION

A planning query is a query whose output is a sequence of objects or actions that gets one from some initial state to some ideal goal state. Automated planning is a branch of artificial intelligence that tries to solve this problem using a computer [10]. However, there is a large class of planning queries that we meet in our daily life that is difficult for a computer to solve, not only because of the involved computational complexity, but because the goal state (as well as the consequence of individual actions) is hard or even impossible to formalize. In contrast, in many of these problems, the goal (and the effect of actions) is intuitively understandable by humans, making the planning humanly possible.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Articles from this volume were invited to present their results at The 39th International Conference on Very Large Data Bases, August 26th - 30th 2013, Riva del Garda, Trento, Italy.

Proceedings of the VLDB Endowment, Vol. 6, No. 9

Copyright 2013 VLDB Endowment 2150-8097/13/07... \$ 10.00.

As a simple example, consider a vacation trip planning. A person may have some tentative start and end dates for her vacation, a preference of what she likes to do and a geographic area where she wants to travel. Based on this data she now needs to compile a potential set of places and attractions to visit and, from this set build a vacation schedule (essentially an ordered subset of the original set). A typical goal here may be to enjoy the vacation the most and/or to expand horizons. Such a goal is naturally subjective and hard to formalize (relevant factors may include total travel distances, attractions along the way, price and many more). However, people sharing similar taste/interests are likely to have the same notion of objective function and their experience and opinion can assist in the planning.

In general we are targeting here problems where one has a large set of items from which she needs to choose a subset and then order this subset in a sequence that will give the best value. The “value” definition is domain-specific, hard to formalize but easy to comprehend by humans. The vacation planning example above is one such instance. Another example is academic schedule planning, where the goal for instance is to obtain solid knowledge of a given subject area.

Answering such planning queries requires *expertise* in the domain of the problem, which is often gained by experience, solving instances of the same (or similar) problems. Since many people deal with similar planning problems, it is reasonable to assume that the *crowd* may provide useful insight here. Indeed, several attempts were made in this direction. For example, for academic schedule planning, the *CourseRank* system¹ allows students to rate courses and provides a convenient tool to compile recommended courses into schedule. Another example is the *Cross-Service Travel Engine for Trip Planning* [3] that allows harvesting POIs (points of interest) from various traveling recommendation sites and provides a tool to compile a trip schedule from these POIs. These systems however focus on identifying the *set* of relevant items (courses, POIs), but the non-trivial task of *ordering* them in an ideal way, to form an actual plan, is left to the user.

Assisting the user in this fairly challenging task is the goal of the present work. We refer below to an ordered list of items as a *plan* and present *CrowdPlanr*, a system that employs the crowd to build “good” plans (w.r.t some abstract quality criteria) for specific tasks. It takes as input a set of relevant items (that can be retrieved from the existing systems mentioned above) and intelligently asks users from the crowd series of simple questions (about possible 1-step

¹<http://courserank.stanford.edu/>

continuations of given partial plans), using the answers to identify the plans preferred by the crowd.

Intuitively, the set of all possible plans (ordered lists) that can be built from a given set of items can be modeled as a tree, where each node is an item, its ancestors are the items preceding it in the plan and its children are the items that may follow it. A root-to-leaf path in this tree represents a plan. One may rate (and correspondingly rank) plans by the probability of a person to consider a given plan as the best (w.r.t to the given abstract criteria). As the size of this tree may be extremely large (exponential in the size of the items set), it is clearly impractical to ask the crowd about each possible plan. Instead, we employ in **CrowdPlanr** a novel efficient algorithm that traverses this tree incrementally. It carefully restricts attention to the more promising plans - ones with highest maximum potential score (to be formally defined in the sequel) and optimally chooses at each step the ‘best’ questions (about possibly continuation), so that the overall number of questions that the crowd needs to be asked is minimized.

Note that the problem we are solving here can be viewed as a particular type of *sorting*. Using the crowd for implementing a sort-by operator is a problem that received much attention in recent crowdsourcing research [25, 5]. A key difference is that all these previous works assume the order between two elements to be independent of preceding elements, and thus the developed algorithms are based on the assumption that users can be asked to compare pairs of individual elements (e.g. be asked if $A < B$). This is not the case here: the order in which plan items are selected depend not only on their individual value/properties but also on what precedes them in the plan (e.g. city A may be more attractive than B , but if in a trip a user first visits C , then A (being rather similar to it) may be skipped altogether and B should be visited instead. Consequently a new algorithm that efficiently provides users with the *context* relevant for their choice had to be developed here.

A first prototype of **CrowdPlanr** was demonstrated in [17]. The demonstration gives only a high level overview of the the system capabilities and user interface. The present paper provides a comprehensive description of the formal model and algorithmic solutions underlying **CrowdPlanr**.

The technical contributions of this paper can be summarized as follows:

- We introduce a simple generic model for modeling plans and interpreting crowd’s answers to questions about them. Based on this model, we develop a formal definition of the planning problem and the identification of (approximated) best answer.
- We present an effective algorithm for identifying the (approximately) best answer using the crowd. As the search space may be extremely large, and consequently the number of questions that may be posed to user excessively high, the algorithm builds the desired plan incrementally, choosing at each step the ‘best’ questions so that the overall number of questions that need to be asked is minimized.
- We study formally the efficiency of our algorithm. Following common practice [8], we employ the notion of

instance-optimality, that reflects how well a given algorithm performs compared to all other possible algorithms in its class and show our algorithm to be instance-optimal for a large common class of planning queries and data instances. Moreover, we show that the optimality ratio that our algorithm achieves (to be formally defined in the sequel) is far by at most a factor of two from the lowest possible optimality ratio.

- Finally, we discuss the implementation of the **CrowdPlanr** and demonstrate, by means of an extensive experimental evaluation, on both synthetic and real life data, that our algorithm consistently outperforms alternative baseline algorithms.

Paper organization. In Section 2 we describe our data model and formally define the planning problem. In Section 3 we present the algorithm we developed to solve this problem. In Section 4 we discuss the algorithm performance, define the notion of instance-optimality and prove our algorithm to be instance-optimal for a large class of inputs. In Section 5 presents experimental results on both synthetic and real-world datasets. In Section 6 we survey related work. We conclude and consider future work in Section 7.

2. PRELIMINARIES

We start with an intuitive description of our model, then proceed to the formal definitions.

We assume that we are given an initial finite set \mathbb{S} of potential items to build a plan from. This set already reflects the preferences the user has defined when she requested a plan. There are multiple domain-specific tools that can be used for identifying this initial set \mathbb{S} of items, e.g. *TripAdvisor*² for vacation trip planning, and we assume that one such tool has been employed. We will use this set to suggest to the user possible answers when we ask a question. Some of these items may become irrelevant as we progress, which will be reflected by the users not selecting them as answers.

CrowdPlanr allows users to build plans at different levels of granularity, zooming in and out between levels. For instance, in a trip to Europe, one can start by planning the countries to visit, then the cities in each country and the attractions within/between cities. Different granularity levels are often independent and we thus focus below, for simplicity, on a single level and explain things in this simplified context. The model extends naturally to the nested case, by allowing users, when dependencies do exist, to view the full detailed plan constructed so far, when considering its continuation.

As a simple running example we will use below the planning of a vacation in Italy (at the city granularity), starting from Rome. The set of items \mathbb{S} in this case includes commonly visited Italian cities, e.g., $\{\text{Milan}, \text{Venice}, \text{Verona}, \text{Florence}, \text{Pisa}, \text{Trento}, \text{Bologna}, \text{Naples}, \dots\}$. Note that, in general, not every user can answer every question. Indeed users that have never visited/read/heard vacation stories about Italy cannot help much in planning a vacation there. The targeting of questions to relevant users is by itself a challenging problem that may be addressed by a variety of methods (e.g. using semantic knowledge about users [1], employing collaborative-filtering based techniques [1, 2], etc.).

²<http://www.tripadvisor.com/>

In principle, any such black-box algorithm can be plugged into our system and we will assume below that the set of relevant crowd users has been identified.

Model. Given a set \mathbb{S} of items, a *plan* is an ordered subset of \mathbb{S} . We will assume, and use, two special items in \mathbb{S} - \dagger to mark the beginning of the plan and \ddagger to mark an end of the plan. A *complete* plan is an ordered sequence of items $(\dagger, a_1, \dots, a_k, \ddagger)$, with no repetitions, starting with a beginning marker and ending with an end marker. We also consider partial plans - prefixes that can be expanded by adding new items; these do not have an end marker. The set of all possible plans may be represented by a tree, called a *decision tree*, where the root is labeled by the start marker, internal nodes are labeled by items from \mathbb{S} , leaves are labeled by the end marker, and each internal node v_i represents a partial plan $p_i = (\dagger, a_1, \dots, a_i)$, corresponding to the labels of nodes on the path from the root to v_i . We use a tree (and not a graph) to model the dependence of the choice of the next item on the entire history of preceding choices.

More generally, one may also want to consider plans where some items are unordered. For instance, when planning an academic schedule, the set of courses taken in a given semester may be unordered. This may be naturally incorporated into our model by having tree nodes that correspond to sets of items rather than individual ones. We do not describe this generalization here.

The decision tree is built iteratively by asking users questions on its nodes. The question on a node v_k is of the form "Given a sequence (\dagger, \dots, a_k) what should be the next item?", where (\dagger, \dots, a_k) are the labels of the nodes on a path from the root to v_k . To answer the user selects an item from \mathbb{S} . Thus, with each question a user is presented with a context of an existing sequence. Answers to these questions define a probability distribution on the children of every node. We use these distributions to define a score for every node - a *score of a node* is its probability to follow its parent in node's partial plan

Formally we define the decision tree as follows:

DEFINITION 2.1 (DECISION TREE). A Decision tree T is a labeled tree $T(V, E)$ with node labels from \mathbb{S} . The root of the tree is labeled by \dagger , leaves may be labeled by \ddagger , and all other node labels are from $\mathbb{S} \setminus \{\dagger, \ddagger\}$. For every node $v \in V$ the set of its children is denoted as:

$$\text{Children}(v) = \{u | u \in V, (v, u) \in E\}$$

In addition, two functions are defined on the nodes of tree:

- $d_T : V \rightarrow \mathbb{N}$ is a display counter. $d_T(v)$ counts number of questions asked on v .
- $c_T : V \rightarrow \mathbb{N}$ is a choice counter. $c_T(v)$ counts how many times v was chosen as an answer. For every node v it must hold that $\sum_{u \in \text{Children}(v)} c_T(u) = d_T(v)$.

For each node v in the tree the combination of its display counter and the choice counters of its children defines a conditional probability distribution of users choosing a particular child to follow v in a sequence. Thus we can easily define a probability of a sequence to be an optimal one by combining the conditional probabilities of the nodes composing it. Formally it can be defined as follows:

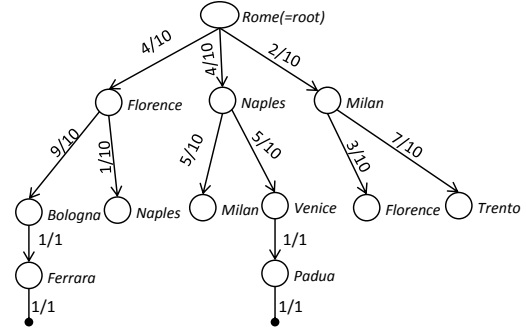


Figure 1: An example of a tree representing a set of plans

DEFINITION 2.2 (NODE SCORE). We define node score in a tree T recursively:

- For the root (node labeled with \dagger): $\text{score}_T(v) = 1$
- For a node u with a parent v , $\text{score}_T(u) = \frac{c_T(u)}{d_T(v)} \text{score}_T(v)$

Example: To continue with our running example, a portion of the tree describing (partial) Italy vacation plans is depicted in Figure 1. The display and choice counts are depicted as labels on the edges incoming the nodes (for example let v, u, w be the nodes labeled with "Florence", "Bologna" and "Naples" respectively, then $d(v) = 10, c(u) = 9$ and $c(w) = 1$). In this figure, 10 questions were asked on most of the nodes, and 1 question on some. Black dots represent leaves marked with \ddagger . The scores of the leaves corresponding to some of the sequences are:

- $(\text{Rome}, \text{Florence}, \text{Bologna}, \text{Ferrara}, \ddagger) - \frac{4}{10} \cdot \frac{9}{10} \cdot \frac{1}{1} \cdot \frac{1}{1} = 0.36$
- $(\text{Rome}, \text{Naples}, \text{Milan}) - \frac{4}{10} \cdot \frac{5}{10} = 0.2$
- $(\text{Rome}, \text{Milan}, \text{Trento}) - \frac{2}{10} \cdot \frac{7}{10} = 0.14$

The previous definitions do not place an upper bound on the number of users that we need to ask in order to compute the probability distribution for a given node. In principle we could ask all available users for each node, but this exhaustive approach can be prohibitively expensive in practice. Instead, we expect applications to place a limit on the number of obtained answers. For this purpose, we define a threshold \mathcal{N} that denotes the desired number of users to be probed at a node. (This may be determined, e.g., based on the desired sampling error bounds [11].) Thus, in principle, by asking \mathcal{N} questions on all of the (incrementally added) nodes (until no more new nodes are added) we can obtain a complete tree.

DEFINITION 2.3 (COMPLETE TREE). A complete tree T is a decision tree in which all leaves are labeled by \ddagger and for each internal node the display counter equals \mathcal{N} .

From the user perspective there is a semantic difference between a complete and partial sequence - a complete sequence cannot be extended further (i.e. the users building it determined that this plan ends here). It makes sense to rank only complete sequences. This difference is naturally reflected in our model where \ddagger markers are used to distinguish complete sequences:

DEFINITION 2.4. A sequence $p = (u_1, \dots, u_k)$ is a complete sequence if and only if u_1 is marked with \dagger and u_k is marked with \ddagger . All other sequences are partial. A set of all complete sequences in a tree T will be denoted as $\mathbb{P}(T)$, and the set of all partial sequences as $\mathbb{P}(T)$. The set of complete sequences in T containing node v will be denoted as $\mathbb{P}_v(T)$.

For example, in Figure 1 the sequence *(Rome, Florence, Bologna, Ferrara)* is a complete sequence, while sequence *(Rome, Naples, Milan)* is a partial one.

Now we can formally define a set of top-k sequences as:

DEFINITION 2.5 (TOP-K SEQUENCES). A set A of complete sequences is a top-k set if $|A| = k$ and for every complete path p' in $\mathbb{P}(T) \setminus A$:

$$\forall p \in A : \text{score}_T(p) \geq \text{score}_T(p')$$

We call the top-1 sequence an optimal sequence.

By nature, the scores computed by sampling a crowd of users are imprecise, in the sense that they only capture general trends: Sequences having similar scores are likely to have a similar “value” for the user. Consequently when two sequences have almost the same score it practically does not matter which one is returned as answer. Namely, it suffices to return a sequence whose score is *approximately* the best. Two types of approximations are common in the literature: relative approximation (i.e. approximation up to a constant *multiplicative* factor) and absolute approximation (i.e. approximation up to an *added* constant). Since we consider here probabilities and when plan scores get very low they become by nature not very interesting, we chose to use additive approximation. Formally we define:

DEFINITION 2.6 (APPROXIMATED TOP-K). A set A of complete sequences is an approximated top-k set if $|A| = k$ and for every complete path p' in $\mathbb{P}(T) \setminus A$:

$$\forall p \in A : \text{score}_T(p) \geq \text{score}_T(p') - \varepsilon$$

The above definitions define a set of optimal (up to a constant) sequences in terms of the complete tree. Note however that, since the size of this tree may be extremely large (exponential in the size of the items set S), it is clearly impractical to build it fully and ask the crowd about each of its nodes. Instead, we employ an efficient algorithm that intelligently traverses the tree and processes only the minimal necessary parts. The algorithm discovers only a partial, as small as possible, decision tree T , that contains sufficient information to guarantee that the set of k highest ranked (up to ε) sequences A in T remains the same in every possible complete tree that can be built by extending T . We call such T a *proof of correctness* for A . We will show in the sequel that the size of the proof of correctness found by our algorithm is $O(\frac{1}{\varepsilon} \cdot |S|)$. Formally, we define a *proof of correctness* as follows.

DEFINITION 2.7 (POSSIBLE COMPLETION). A complete tree T' is a possible completion of a decision tree T if the following conditions hold:

1. T is a subtree of T'
2. $\forall v \in V_T : d_{T'}(v) \geq d_T(v)$

$$3. \forall v \in V_T : c_{T'}(v) \geq c_T(v)$$

We denote a set of possible completions of T as $\text{Compl}(T)$.

DEFINITION 2.8 (PROOF OF CORRECTNESS). A decision tree T is a proof of a set A being a top-k set if for all $T \in \text{Compl}(T)$:

$$\forall p \in A : \forall p' \in \mathbb{P}(T) \setminus A : \text{score}_T(p) \geq \text{score}_T(p') - \varepsilon$$

Example: In the tree presented in Figure 1 the sequence $p = (\text{Rome}, \text{Florence}, \text{Bologna}, \text{Ferrara})$ is the highest ranked sequence, however if we take $\varepsilon = 0.01$ then this tree is not a proof of correctness for the set $\{p\}$ - indeed, there is a possible continuation of this tree - T' , where additional 9 questions are asked (recall that $\mathcal{N} = 10$) on *Bologna* node, and for all these questions we get “*Trento*” as an answer. Then, in T' , the sequence p will have a score of $\frac{4}{10} \cdot \frac{9}{10} \cdot \frac{1}{10} = 0.036$, while a sequence $p' = (\text{Rome}, \text{Florence}, \text{Bologna}, \text{Trento})$ will have a score of $\frac{4}{10} \cdot \frac{9}{10} \cdot \frac{9}{10} = 0.324$.

3. PLANNING USING CROWD

We are now ready to present our algorithm (Algorithm 2) for finding the optimal, up to ε plan. For brevity we will omit below the words “up to a constant” and whenever refer to an optimal plan we mean optimal up to a constant.

To simplify the presentation we will focus on finding the Top-1 sequence. Our results naturally generalize to the Top-k case and we briefly discuss the extension in Section 7³.

Our algorithm for finding an optimal plan will hold a decision tree (initially containing only the root) and will expand it by asking questions on its nodes. To achieve its goal the algorithm has to solve the two following sub-problems:

- Checking stop condition - i.e. checking whether the current tree is a proof of correctness for the current optimal plan
- Deciding which next question to ask in order to reach stop condition as fast as possible

The algorithm is inspired by the well-known A^* algorithm [13] and the key challenge was to find the appropriate solution for these two points, that guarantee optimality.

3.1 Stopping condition

To solve the first sub-problem we define a notion of uncertainty for a sequence. Uncertainty is the maximum possible difference between a given sequence score and the highest sequence score in all possible completions of the current state of the tree. Formally it is defined as follows.

DEFINITION 3.1 (UNCERTAINTY). In a tree T , an uncertainty for a complete sequence p is given by:

$$U(T, p) = \max_{T' \in \text{Compl}(T)} \left[\max_{p' \in \mathbb{P}(T')} \text{score}_{T'}(p') - \text{score}_T(p) \right]$$

Following definitions 3.1 and 2.8 we can use the uncertainty notion to check whether a decision tree is a proof of correctness for a sequence in it:

LEMMA 3.2. A tree T is a proof of correctness for a complete sequence p iff $U(T, p) < \varepsilon$.

³Full details can be found in the technical report [14]

A naïve approach to calculating the value of uncertainty of a given decision tree would be to enumerate its possible completions. However, this approach is ineffective since every incomplete node of a tree can be extended with an arbitrary sub-tree. Instead, we use an efficient algorithm (Algorithm 1) that traverses the current decision tree only once in order to calculate the uncertainty.

Algorithm 1 Calculating $U(T, p)$

Assuming $p = (u_1, \dots, u_m)$

- 1: $Deltas \leftarrow \emptyset$
 - 2: **for all** $\{v|v \in V, v \neq u_m, d(v) < \mathcal{N}\}$ **do**
 Assuming v is a part of a path
 $p' = (u_1, \dots, u_k, v_1, \dots, v_n = v)$ and
 (u_1, \dots, u_k) is a common prefix of p and p'
 (v_1, \dots, v_n) is the remainder of p'
 - 3: $maxCommon \leftarrow \prod_{i=2}^k \frac{c_T(u_i) + \mathcal{N} - d_T(u_{i-1})}{\mathcal{N}}$
 - 4: $maxPPrime \leftarrow \frac{\mathcal{N} - d_T(v)}{\mathcal{N}} \prod_{i=2}^n \frac{c_T(v_i) + \mathcal{N} - d_T(v_{i-1})}{\mathcal{N}}$
 - 5: $minP \leftarrow \prod_{i=k+1}^m \frac{c_T(u_i)}{\mathcal{N}}$
 - 6: $\delta \leftarrow maxCommon \cdot (maxPPrime - minP)$
 - 7: $Deltas \leftarrow Deltas \cup \{\delta\}$
 - 8: **end for**
 - 9: **return** $\max_{\delta \in Deltas} \delta$
-

This algorithm exploits the fact that the maximum difference in scores is achieved when one of the sequences gets its lowest possible score, while some other sequence gets its highest possible score. The algorithm goes iteratively over all nodes in T that we can ask more questions on and for every node v builds a sequence p' that ends one step after v (i.e. a shortest complete sequence that contains v). The algorithm then calculates maximum possible score difference between p' and p (line 3-6). The maximal common prefix of the two sequences is designated as u_1, \dots, u_k . To achieve the maximum possible difference the algorithm assigns: highest possible score to the common part of p and p' (line 3), highest possible score to the remainder of p' (line 4) and lowest possible score to the remainder of p (line 5). At the end, the algorithm returns the maximum of the calculated differences.

Example: While calculating the uncertainty of a path ending by a node labeled “Padua” in a tree presented in Figure 1, the algorithm will build a sequence p' for a node “Bologna”: $p' = (\text{Rome}, \text{Florence}, \text{Bologna})$. The common prefix contains only the root, thus $maxCommon = score(\text{root}) = 1$, $maxPPrime = \frac{9}{10} \cdot \frac{4}{10} \cdot \frac{9}{10} = \frac{324}{1000}$ and $minP = \frac{4}{10} \cdot \frac{5}{10} \cdot \frac{1}{10} = \frac{20}{1000}$, and finally $\delta = \frac{304}{1000}$. The same calculation will be performed for all other nodes of the tree and the maximum δ will be returned.

Theorem 3.3 formally proves the algorithm correctness.

THEOREM 3.3. *Given a decision tree T and a complete sequence p in it Algorithm 1 calculates $U(T, p)$.*

PROOF. The maximum possible score of a complete sequence p in any $T' \in Compl(T)$ is upper bounded by T ’s current state - indeed there are only two options for p :

1. $p \in \mathbb{P}(T)$, then $p = (u_1, \dots, u_k)$ will get a maximum score if for all remaining questions for every u_i , u_{i+1}

will be chosen as an answer. In this case,

$$\max_{T' \in Compl(T)} score_{T'}(p) = \prod_{i=2}^k \frac{c_T(u_i) + \mathcal{N} - d_T(u_{i-1})}{\mathcal{N}}$$

2. p is a continuation of some partial sequence $p' \in \mathcal{P}(T)$ (i.e. p' is a prefix of p), then the maximum score of p in $T' \in Compl(T)$ is exactly the maximum score of p' in T' (the maximum is achieved if all users select p as the only continuation of p'), and thus it can be calculated as in previous case.

On the other hand, the minimal possible score for a sequence $p = (u_1, \dots, u_k)$ is achieved if for all the remaining questions on node u_i all the answers will be different than u_{i+1} . And its minimal score would be:

$$\min_{T' \in Compl(T)} score_{T'}(p) = \prod_{i=2}^k \frac{c_T(u_i)}{\mathcal{N}}$$

Finally, if we have two sequences $p_1 = (u_1, \dots, u_k, v_1, \dots, v_n)$ and $p_2 = (u_1, \dots, u_k, w_1, \dots, w_m)$ (u_1, \dots, u_k is the common prefix of the two sequences) then the difference in their scores in a possible continuation T' is given by:

$$score_{T'}(p_1) - score_{T'}(p_2) = \left(\prod_{i=2}^k \frac{c_{T'}(u_i)}{\mathcal{N}} \right) \cdot \left(\prod_{i=1}^n \frac{c_{T'}(v_i)}{\mathcal{N}} - \prod_{i=1}^m \frac{c_{T'}(w_i)}{\mathcal{N}} \right)$$

And thus, it is maximized when one of the sequences gets all of the remaining votes (including the common prefix part of the sequence) and the second sequence (except for the common prefix) gets no more votes. \square

3.2 Which questions to ask

The second sub-problem any algorithm for finding an optimal sequence has to solve is deciding what question to ask next. We employ a greedy approach to solve this problem - we ask questions on a sequence with the highest “potential”, i.e. a sequence with a highest potential score. This approach is effective (as we will show in section 5) and can further be extended for identifying a bulk of ‘best questions’, e.g. when multiple questions may be posed to users in parallel. For clarity we explain next in details how to choose a single next question, then briefly consider the selection of multiple questions.

Formally, the notion of sequence potential is defined as follows:

DEFINITION 3.4 (POTENTIAL SCORE). *Given a tree T and a sequence p (partial or complete) in it:*

$$M_T(p) = \max_{T' \in Compl(T)} score_{T'}(p)$$

In general, there may be several sequences that have the highest potential score. Since we do not have any additional information that allows us to prefer one over the other, we will consider all of them in a round-robin.

Each iteration of this algorithm finds a node in the current decision tree T and asks a question on it. The algorithm stops (condition on line 3) when the uncertainty of some node p in the tree drops below ε . When this happens, following Lemma 3.2, T is the proof of correctness for p .

On line 4 we find a sequence (partial or complete) with the

Algorithm 2 Finding the optimal sequence

```
1:  $T \leftarrow \text{origin}$ 
2:  $i \leftarrow 0$ 
3: while  $\mathbb{P}(T) = \emptyset$  OR  $\min_{p \in \mathbb{P}(T)} U(p, T) \geq \varepsilon$  do
4:    $Candidates \leftarrow \{\text{argmax}_{p \in \mathbb{P}(T)} M(p)\}$ 
5:    $TopNodes \leftarrow \{tN(\text{argmax}_{p \in \mathbb{P}(T)} M(p))\}$ 
6:    $maxScore \leftarrow \max_{p \in \mathbb{P}(T)} M(p)$ 
7:   if  $maxScore > \varepsilon$  then
8:     for all  $p \in \mathbb{P}(T)$  do
9:       if  $M(p) = maxScore$  then
10:        if  $tN(p) \notin TopNodes$  then
11:           $Candidates \leftarrow Candidates \cup \{p\}$ 
12:           $TopNodes \leftarrow TopNodes \cup \{tN(p)\}$ 
13:        end if
14:      end if
15:    end for
16:     $p \leftarrow Candidates[i \bmod |Candidates|]$ 
17:    Ask a question on  $tN(p)$ 
18:     $i \leftarrow (i + 1) \bmod |Candidates|$ 
19:  else
20:     $p \leftarrow Candidates[0]$ 
21:    Ask a question on lowest node of  $p$ 
22:  end if
23: end while
24: return  $\text{argmin}_{p \in \mathbb{P}(T)} U(p, T)$ 
```

If there is more than one minimum(maximum) item, $\text{argmin}(\text{argmax})$ shall return one arbitrary

maximum potential score. Following the proof of Theorem 3.3, max potential score can be calculated using a simple formula in linear time (in the length of the sequence).

Given a sequence we have also to choose a node in it to ask a question, this is done on line 5. We prefer to ask questions on higher nodes as they affect more paths in the tree. Formally we define:

DEFINITION 3.5 (TOP-NODE). *For a path $p = (v_1, \dots, v_k)$ a top node - $tN(p)$ is a node v_i , s.t. i is the minimum item in the set $\{i | 1 \leq i \leq k, d(v_i) < \mathcal{N}\}$ (i.e. the topmost node at which we have not asked \mathcal{N} questions).*

The loop in lines 8-15 selects all the sequences that have the maximum potential score. We ask questions on all of them in the round-robin manner. Finally, on line 7 we check for a special condition - if all of the sequences in the tree cannot have score greater than ε , then it does not matter which sequence we return - all of them are optimal by definition, thus we just need to discover one complete sequence and return it. The easiest way to do it is by asking questions on the lowest possible node - take any sequence and ask a question on its last node, if the answer terminates the sequence - return the sequence, otherwise ask a question on a newly discovered node.

Example: given a decision tree presented in Figure 1 our algorithm will ask a question on a node labeled “Bologna” since it’s the highest non-exhausted node of a sequence with the highest potential (the potential of a sequence (Rome, Florence, Bologna, Ferrara) is $\frac{4}{10} \cdot \frac{9}{10} \cdot \frac{10}{10} = \frac{36}{100}$).

It is clear that the algorithm eventually halts - the number of questions we can ask is bounded by the size of \mathcal{T} (times \mathcal{N}) which is finite. It is also clear that when it does, there is a sequence p in T for which $U(T, p) < \varepsilon$ and this is the

sequence that is returned (lines 3 and 24). Thus, following Lemma 3.2 the algorithm returns an optimal sequence. In Section 4 we perform a detailed analysis of the algorithm’s efficiency.

Asking questions in bulk. A common problem in crowdsourcing applications is assigning a bulk of questions - we have M users ready to answer questions, so we want to ask them all at once to prevent wasting a valuable human time. Algorithm 2 was constructed to ask one question at time, but it can easily be extended to assign a bulk of questions, by choosing questions from the following sets:

- On every node that enters *Candidates* list we may need to ask up to \mathcal{N} questions. All may be asked in parallel
- The *Candidates* list may contain several nodes, all equivalent from the algorithm’s point of view. Questions on them can thus be asked in parallel
- Finally, the *Candidates* list contains nodes with currently maximum potential. If there are users pending we can also ask questions on nodes that are top-k in potential.

There is a trade-off between utilizing many users in parallel and asking the minimal possible number of questions since every answer we get helps us target next questions better, thus asking questions one by one helps minimizing the total number of questions. On the other hand, asking questions in parallel helps utilizing more users. Detailed analysis of this trade-off is an interesting direction for future work.

4. EFFICIENCY AND OPTIMALITY

In this section we will discuss the efficiency and the optimality of the algorithm presented above and will also provide a lower-bound for the possible optimality ratio. To discuss optimality we need to define a cost measure and a set of inputs, based on which we will compare different algorithms.

We use the number of visited nodes in a tree (i.e. number of nodes we asked questions about) as our *cost measure*. This number is in direct correlation to the actual number of questions asked - indeed we assumed that in order to learn a probability distribution of a continuation from a node, one has to ask \mathcal{N} questions on this node. Furthermore, using number of nodes as a cost measure, rather than the actual number of asked questions, makes reasoning about optimality much simpler.

The class of inputs we consider is the set of all possible complete trees composed from items in \mathbb{S} where the difference between the shortest and the longest sequence is at most \mathbf{k} , for some predefined constant \mathbf{k} . For a given \mathbf{k} we denote the corresponding class of inputs as $\mathbf{I}_{\mathbf{k}}$. As we saw in our experiments, typical real-world inputs fall into $\mathbf{I}_{\mathbf{k}}$ for fairly small value of \mathbf{k} (comparable plans of the same granularity usually contain similar number of items).

We use an instance-optimality notion as it appears in [8]:

DEFINITION 4.1 (C-OPTIMALITY). *For a class of inputs \mathbf{I} and a class of algorithms \mathbf{A} , algorithm $\mathcal{A} \in \mathbf{A}$ is c-optimal if for every input $\mathcal{I} \in \mathbf{I}$ and for every algorithm $\mathcal{B} \in \mathbf{A}$:*

$$\text{cost}(\mathcal{A}, \mathcal{I}) \leq c \cdot \text{cost}(\mathcal{B}, \mathcal{I}) + c'$$

We refer to c as the optimality ratio. c' is also a constant.

We will prove next the following two results. The first proves the instance optimality of our algorithm and the second shows its optimality ratio is far by at most a factor of two from the lower bound.

THEOREM 4.2. *Algorithm 2 (A) is $\frac{1}{\varepsilon} - \text{optimal}$ on trees from \mathbf{I}_k .*

THEOREM 4.3. *Let \mathcal{B} be a deterministic algorithm for finding an optimal (up to ε) sequence which is $c - \text{optimal}$ on trees from \mathbf{I}_k . Then $c \geq \frac{2}{\varepsilon}$.*

To prove Theorem 4.2 we will first analyze the performance of Algorithm 2 and show that it asks questions about at most $\frac{1}{\varepsilon} - 1$ different sequences (later, in Section 5, we show that in real-life cases our algorithm considers even less sequences). To do this we will introduce a concept of “Candidates pool”, a set of nodes that our algorithm considers to ask questions about. We will show that during the run of the algorithm, only a limited number of nodes can enter the pool, thus limiting the total number of different paths considered by our algorithm.

THEOREM 4.4. *Algorithm 2 asks questions about at most $\frac{1}{\varepsilon} - 1$ different sequences during its run.*

PROOF. Let’s call Algorithm 2 A. Only sequences with potential maximum score greater than ε are considered by A, A asks questions only on the top-nodes of the sequences in *Candidates*. Let P_i be the set (pool) of all top-nodes that are part of a path with potential maximum greater than ε after question i . It is clear that the node for question $i + 1$ is chosen only from P_i . If a node v was in P_i but is not in P_{i+1} we say that node v has *left the pool*. Node v can leave the pool only in one of the following cases:

1. v is no longer a top-node, i.e. all possible questions on it were asked, but it still belongs to a path with potential max score greater than ε . In this case one or more children of v will be in P_{i+1} .
2. v no longer belongs to a path with potential max score greater than ε , it means that no descendants of v will be in P_j for $\forall j > i$.

It is clear that once a node has left the pool, it cannot return. Let *Out* be the set of all nodes that left the pool for the reason 2, formally $Out = \{v | \exists 1 \leq i \leq m : v \in P_i \wedge v \notin P_{i+1} \wedge (\forall u \in Children(V) : u \notin P_{i+1})\}$.

After question m every path that was ever considered by A has a node that is a part of it in either P_m or *Out*.

$M(p) \leq score_T(tN(p))$: For every possible completion T' of T , $p \in \mathbb{P}_{tN(p)}(T')$, thus following Lemma 4.6 $score_{T'}(p) \leq score_{T'}(tN(p))$, by the definition of top-node, its score is final (since the display count of its parent equals \mathcal{N}), thus $score_{T'}(tN(p)) = score_T(tN(p))$, hence following the definition of potentially maximum score we get that $M(p) \leq score_T(tN(p))$.

This means that for every $v \in P_m$, $score(v) > \varepsilon$ and also for every $u \in Out$, $score(u) > \varepsilon$ (since every node in *Out* was once in the pool and its score was already final then).

No two nodes in P_m are a part of a same path (by the definition of *top-node*). The same is true for nodes in *Out* (indeed if v has moved to *Out*, its children cannot even enter the pool so they cannot be moved to *Out* either). There are also no

$v \in P_m$ and $u \in Out$ such that u, v are parts of the same path (for the same reason). Thus all the nodes in $P_m \cup Out$ are parts of a different sequences. Hence (by Lemma 4.7) $\sum_{v \in P_m \cup Out} score(v) \leq 1$. And since for every $v \in P_m \cup Out$, $score(v) > \varepsilon$ we have that $\sum_{v \in P_m \cup Out} score(v) > |P_m \cup Out| \cdot \varepsilon$.

From these two facts we get that $|P_m \cup Out| < \frac{1}{\varepsilon}$. Thus algorithm A considers at most $\frac{1}{\varepsilon} - 1$ different sequences during its run. \square

COROLLARY 4.5. *The size of the proof of correctness tree found by Algorithm 2 is $O(\frac{1}{\varepsilon} \cdot |S|)$.*

PROOF. The proof follows from the fact that the tree contains at most $\frac{1}{\varepsilon}$ different paths, and each path contains at most $|S|$ nodes. \square

To complete the proof of Theorem 4.4 we prove the following two lemmas.

LEMMA 4.6. *For every node $v \in V$, $\sum_{p \in \mathbb{P}_v(T)} score(p) = score(v)$. In particular $\sum_{p \in \mathbb{P}(T)} score(p) = 1$.*

PROOF. By induction on the tree structure. For leaves the claim is true since $\mathbb{P}_v(T)$ contains exactly one path - from the root to v , and thus $\sum_{p \in \mathbb{P}_v(T)} score(p) = score(v)$ by definition.

Let v be an internal node and assume the claim is true for its children. Let $\{u_1, u_2, \dots, u_n\}$ be the set of v ’s children. For every u_i there is a single path from the root to u_i - $\{w_1, w_2, \dots, w_k\}$ (where w_k is u_i and w_{k-1} is v). Thus:

$$\begin{aligned} score(u_i) &= \prod_{j=2}^k \frac{c(w_j)}{d(w_{j-1})} = \frac{c(w_k)}{d(w_{k-1})} \prod_{j=2}^{k-1} \frac{c(w_j)}{d(w_{j-1})} \\ &= \frac{c(u_i)}{d(v)} score(v) \end{aligned}$$

Since a path containing u_i cannot contain u_j we can split $\mathbb{P}_v(T)$ into $\bigcup_{i=1}^n \mathbb{P}_{u_i}(T)$, and thus:

$$\begin{aligned} \sum_{p \in \mathbb{P}_v(T)} score(p) &= \sum_{i=1}^n \sum_{p \in \mathbb{P}_{u_i}(T)} score(p) = \\ &= \sum_{i=1}^n score(u_i) = score(v) \sum_{i=1}^n \frac{c(u_i)}{d(v)} \end{aligned}$$

Choice counters of the children sum up to a display counter of the parent, thus $\sum_{i=1}^n \frac{c(u_i)}{d(v)} = 1$, and hence $\sum_{p \in \mathbb{P}_v(T)} score(p) = score(v)$. \square

LEMMA 4.7. *Let A be a set of nodes of a tree T , s.t. no two nodes in A are a part of a same sequence. Then $\sum_{v \in A} score_T(v) \leq 1$.*

PROOF. Let’s build a tree T' from T by terminating a path at each $v \in A$, the rest of the tree remains as is. Now, for each $v \in A$ there is a path p_v in T' (since all the nodes in A are parts of different paths, there is no conflict). Also, $score_{T'}(p_v) = score_T(v)$ (since we left the rest

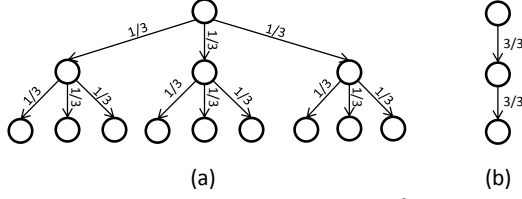


Figure 2: (a) Scattered-tree of depth 2 ($\mathcal{N} = 3$). (b) Chain-tree of length 2 ($\mathcal{N} = 3$).

of the tree as is, all the nodes from the root to v in T and T' have the same counter values). $\{p_v | v \in A\} \subseteq \mathbb{P}(T')$, thus by Lemma 4.6 $\sum_{v \in A} \text{score}_T(v) = \sum_{p \in \{p | v \in A\}} \text{score}_{T'}(p) \leq \sum_{p \in \mathbb{P}(T')} \text{score}_{T'}(p) = 1$. \square

Now we can prove the optimality ratio of Algorithm 2.

PROOF OF THEOREM 4.2. Let $\mathcal{T} \in \mathbf{I}$ be some ground truth tree and let p be the shortest optimal (up to ε) path in it, with length l . \mathcal{A} will consider at most $\frac{1}{\varepsilon}$ different paths during its run, one of them will be p . Each one of these paths is at most $l + k$ nodes long (following the assumption that paths length varies by at most k nodes), thus in total the algorithm will visit at most $\frac{l+k}{\varepsilon}$ nodes. On the other hand any other algorithm \mathcal{B} will have to visit at least l nodes to discover and return p (or any other optimal path, since p is the shortest optimal path). Thus:

$$\begin{aligned} \text{cost}(\mathcal{A}, \mathcal{T}) &\leq \frac{l+k}{\varepsilon} \\ \text{cost}(\mathcal{B}, \mathcal{T}) &\geq l \end{aligned}$$

Combining this we get:

$$\text{cost}(\mathcal{A}, \mathcal{T}) \leq \frac{1}{\varepsilon} \cdot \text{cost}(\mathcal{B}, \mathcal{T}) + \frac{k}{\varepsilon}$$

$\frac{k}{\varepsilon}$ is a constant independent of the input, thus following the definition, Algorithm 2 is $\frac{1}{\varepsilon}$ -optimal. \square

And finally we prove that there is no deterministic algorithm that has optimality ratio better than $\frac{2}{\varepsilon}$. To do so, we will show that for every deterministic algorithm we can construct an input that will require from it to consider $\frac{2}{\varepsilon}$ different paths, while an optimal algorithm will have to consider only one path.

PROOF OF THEOREM 4.3. For the sake of the proof we will first define 2 special types of subtrees:

- A *Chain-tree* of length k is a subtree consisting of k nodes u_1, \dots, u_k , where u_i is the only child of u_{i-1} for every $1 < i \leq k$. In this subtree $\text{score}(u_k) = \text{score}(u_1)$.
- A *Scattered-tree* of depth k is a subtree rooted in u of depth k where every node has exactly \mathcal{N} children. Every leaf in this subtree has a score of $\frac{\text{score}(u)}{\mathcal{N}^k}$.

Figure 2 illustrates these types of subtrees.

Let's analyze the performance of \mathcal{B} when running on a set of inputs $\{\mathcal{T}_x | x \geq 1\}$. A tree \mathcal{T}_x will be constructed as follows:

- Let k be the largest integer such that $(\frac{1}{2})^k > \varepsilon$

- Starting from the root, which is considered to be on level 0, every node on level $0 \leq i \leq k$ will have 2 children, every child will have choice count of $\lfloor \frac{\mathcal{N}}{2} \rfloor$ (If \mathcal{N} is odd, additional child will be added to every node with choice count of 1 and a scattered-tree underneath it).
- After this, on the level k we will have c leaves, each one of them with a score of $s = (\frac{1}{2})^k$.
- By the way of selection of k we ensured that $s > \varepsilon$ and $\frac{1}{2}s \leq \varepsilon$. Since all the leaves have the same score and their score sum up to 1 (Lemma 4.6) we have that $c \geq \frac{1}{2\varepsilon}$ and since c is an integer, $c \geq \lceil \frac{1}{2\varepsilon} \rceil$. Let's denote these leaves as $u_1 \dots u_c$.
- Under each one of the u_i 's we will put a chain-tree of length M , where M is an arbitrarily large number, we will denote the end of each chain-tree as v_i
- Under each one of the v_i 's except for v_x we will put a scattered-tree of depth T , such that every leaf will have a score less than $s - \varepsilon$, under v_x we will put a chain-tree of length T .
- Every leaf we have now will be labeled with \dagger , making the tree complete

This tree has exactly one correct answer (under v_x). Now, suppose \mathcal{B} does not consider a path under u_i during its run, for some i . Then, \mathcal{T}_i is indistinguishable from \mathcal{T}_x up until the discovery of u_i and the correct answer is under u_i , thus when running against \mathcal{T}_i , \mathcal{B} will not discover a correct answer, contradicting the assumption that \mathcal{B} is a correct algorithm. So \mathcal{B} considers at least c different paths during its run. On the other hand, algorithm \mathcal{B}_x that discovers all of the u_i 's and then proceeds asking questions only about u_x can return after asking \mathcal{N} questions on every node of the u_x chain-tree (since under each u_i all paths have a score of at most s , discovering one path with score s is enough to return a correct answer), thus considering only 1 path during its run. Length of every path in \mathcal{T}_x tree equals $k + M + T$. \mathcal{B} asks questions on at least $k + \frac{1}{2\varepsilon} \cdot M + T$ nodes, while \mathcal{B}_x asks questions on $k + M + T$ nodes. Since M can be arbitrarily large the optimality ratio between \mathcal{B} and \mathcal{B}_x is at least $\frac{1}{2\varepsilon}$. \square

5. EXPERIMENTAL EVALUATION

In this section we will present the results of the experimental evaluation of our algorithm. During the evaluation we explored its behavior on different data sets (both synthetic and real) with different parameters. We also explored the effect of the algorithms parameters (such as allowed error and the required number of answers per node). We compared our algorithm to several baseline algorithms.

5.1 Evaluation setup

To conduct the experiments we have implemented *CrowdPlanr* in C# and PHP while using MySQL as the database engine. Its architecture is presented in Figure 3. For the evaluation purposes the User Interface was replaced with a simulator (called *the Oracle* in the sequel) that returned answers to queries either from a synthetically generated dataset or a dataset recorded from the interaction with real users. Other parts of the system are: *Plan Builder* which executes

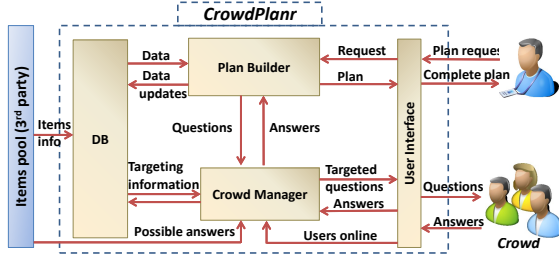


Figure 3: CrowdPlanr architecture

the algorithm for finding the optimal sequence, *Crowd Manager* which formulates the questions for users (by preparing a set of possible answers) and *Database* which holds all the information gathered by the system.

Algorithms. In our experiments we compared the number of visited nodes (and the number of questions asked) by our algorithm (as defined in Chapter 3, from now on it will be called *CrowdPlanr*) to the number of nodes visited by the baseline algorithms running on the same input (since the problem presented in this work has not been studied before, we can't compare our algorithm to other solutions). The baseline algorithms that were considered are:

- *Random* - a naïve algorithm that randomly chooses which question to ask next, from all possible questions (nodes of the tree that are not exhausted yet). This algorithm showed extremely poor performance (asked significantly more questions than all other algorithms), and thus we will not include this algorithm in our comparisons
- *Greedy* - an algorithm that employs a trivial greedy approach: ask a question on a sequence that currently has the maximal score, try to extend it as much as possible (i.e. ask a question on a lowest node possible of the selected sequence).
- *CrowdPlanr⁻* - compared to the greedy, our algorithm is different in two ways, first we choose a sequence with the highest potential score (not the highest current score) and second we ask questions on a highest node possible (not always trying to extend the selected sequence). *CrowdPlanr⁻* is an algorithm that is halfway between the *Greedy* and the *CrowdPlanr* algorithms: it selects a sequence with the highest current score (like the *Greedy*) and asks questions on a highest node possible of that sequence (like the *CrowdPlanr*).

The halting condition for all the algorithms is the same: they can return a sequence only if the tree they had discovered so far forms a proof of correctness (recall Definition 2.8 for that sequence). The algorithms use uncertainty calculation we presented in Section 3 to check this condition.

Synthetic Data. For evaluating the effect of various properties of the input on the number of questions asked by the algorithms we generated a synthetic datasets simulating an input with desired properties. These datasets were represented by complete trees accessible by the *Oracle*. In each experiment we changed only one property of the input while all the others had a default value. We considered the following properties of the input:

- *Tree depth* is ranged from 5 to 10, with default value of 7. Trees that have more than 10 levels are less important for two main reasons: first, human factor reason, in the real world scenarios it is hard for the user from the crowd to hold in her mind more than 10 items as a context and give a good recommendation for continuation of the sequence (usually when one wants to plan a longer sequence, she will do the planning on different granularity levels). Second, since we use probabilities and the final score of the answer is a multiplication of the probabilities of the nodes, for sequences longer than 10 nodes the scores get very small in our setup and hence all the sequences will be optimal up to ϵ (Definition 2.6).
- *Depth difference (k)* is the difference in levels between the highest and the lowest leaf in the tree. This property shows the balance of the tree. The default value is 0, which means that all the sequences have the same length. We also ranged k values up to $\frac{TreeDepth}{3}$. Our experiments showed that this parameter also does not affect the number of questions asked by the algorithms, thus we will not discuss these experiments in detail.
- *Skewness* is a percent of votes that go in favor a specific child of each node in tree (for example if the *Skewness* is 60% then for every node that will be a child that gets 60% of the votes). The default value for skewness is 60%, we also checked skewness of 50% and 70%.

Real Data. To ensure that our algorithm performs well in real life we evaluated the number of questions asked by it (and its baseline competitors) on two datasets coming from different real-world applications:

1. A Large dataset containing a record of 20,000 vacation trips in Europe. The trips included 10 different cities and were approximately of the same length (in terms of visited cities). This dataset was obtained from a traveling agency, we omit its name for privacy reasons.
2. A Medium size dataset containing answers to a question "In which order to watch Star Wars films?". It was obtained by comparing the popularity of the proposed orders on various web sites. This question is asked frequently on the internet and has 100,000,000 results in web search engines. It has $6!$ ($=720$) possible answers (some of them, of course, completely wrong). An important property of this dataset is that there are only a small number of "good" orders, while the rest have very little support.

All the datasets initially contained the ranking of the sequences and were translated into a complete tree that was used by the *Oracle* to answer queries. We believe this is a good approximation of a real-life interaction with the users, because we assume that the users know the rating of a complete sequence and thus their rating of the partial sequence (i.e. the answer to our questions) will be consistent with it.

Algorithm parameters. In addition to the properties of the input we evaluated how the parameters of the algorithm affect the number of nodes visited by it (and the number of questions asked by it). We evaluated the effect of the following parameters:

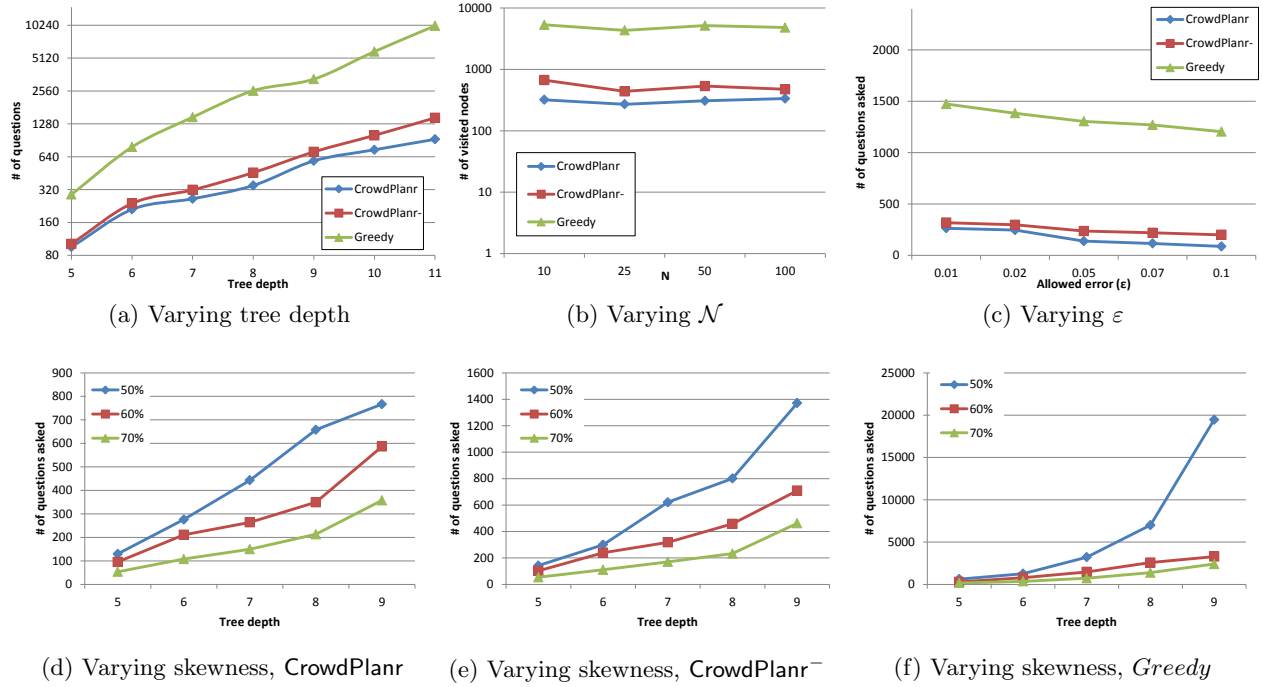


Figure 4: Experiments results

- *Allowed error (ε)* The default value of an allowed error in our experiments was 0.01, in addition we ranged it from 0.002 to 0.1.
- *Number of questions per node (\mathcal{N})* is chosen based on a statistical data, reliability of the crowd and budget constraints. We use a default value of 10, but also run experiments with $\mathcal{N} = 50$ and $\mathcal{N} = 100$.

5.2 Synthetic data evaluation results

In every experiment for each of the algorithms we measured the number of nodes visited (node is considered visited if there was at least one question asked about it) and the total number of questions asked. For each experiment we executed the algorithms on 3 datasets with the same properties and averaged the results.

The experiments were modeled in the following way: generated datasets were read by the *Oracle* and all the questions asked by the algorithms were redirected to the *Oracle* which answered them basing on an input data set in a deterministic way (i.e. several run using the same oracle would yield the same results). The *Oracle* also collected statistics about asked questions from which we derived the results of the experiments.

The results of the experiments are summarized in Figure 4, next we will briefly describe every experiment.

Varying the depth of the tree. In this experiment we evaluated the effect of a tree depth. All other parameters remained default. The results are presented in Figure 4(a). Note that the Y axis has a logarithmic scale and it represents the number of questions that were asked. The X axis represents the depth of the tree. We can see that the *Greedy* algorithm performs significantly worse than both the *CrowdPlanr* and the *CrowdPlanr*⁻ algorithms. Also, we can

see that the difference between the *CrowdPlanr* and the *CrowdPlanr*⁻ grows with the depth of the tree.

Varying the data skewness. In this experiment we examined the effect of data skewness on the number of questions asked by the algorithms. We tested skewness levels of 50%, 60% and 70%. All other parameters had default values. The results of the experiment are presented in Figure 4(d-f). As can be seen in this graph, the more data is skewed the easier it is for the algorithm to find a correct answer. The most prominent effect skewness has on the *Greedy* algorithm.

Varying the number of questions per node. In this experiment we evaluated the impact of different values for \mathcal{N} (10, 25, 50, 100) on the number of nodes visited by the algorithms (the total number of questions is less interesting since it is expected to be linearly dependent on \mathcal{N}). All other parameters had a default value. The results of this experiment are shown in Figure 4(b). Theoretically we have proven that our algorithm is instance-optimal and its optimality ratio does not depend on \mathcal{N} , that means that *CrowdPlanr* algorithm is expected to visit the same number of nodes for any value of \mathcal{N} . In practice it means that *CrowdPlanr* can be used for different approaches with different audience, crowd size and system needs. The graph shows that the reality meets the expectation.

Varying allowed error. In this experiment we evaluated the effect of the allowed error value on the number of questions asked by the algorithms. For this we left all the parameters with default values and ranged ε from 0.01 to 0.1. The results are summarized in Figure 4(c). Here again, the Y axis has a logarithmic scale and represents the number of questions asked by the algorithm. The X axis represents the value of ε . We can see that the effect of the allowed error is much larger on the *Greedy* algorithm than on the other two. One possible explanation to this is that the *Greedy*

strategy causes the algorithm to “jump” all over tree without really reducing the uncertainty. Increasing the allowed error sets the uncertainty bar lower, thus saving the algorithm work. Another interesting property of this graph is that the number of questions asked by the **CrowdPlanr**⁻ algorithm decreases relatively slow with the increase of ε . A possible explanation for this can be that the selection criteria of the sequence to work on is wrong, which causes the algorithm to ask questions on a wrong sequence, and since it asks questions on a highest possible node, the mistake is not revealed fast enough.

Additional experiments. In our last experiment we explored the behavior of all of the algorithms - how many questions were asked on every level of a decision tree. This tells us how “focused” each one of the algorithms is. For this experiment, we fixed the tree depth to be 10 and left all other parameters to have default values. The results are presented in Figure 5.

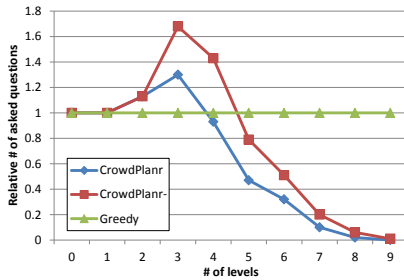


Figure 5: Questions asked on each level

In this graph the X axis represents the level in the tree (0 represents the root, 9 represent leafs) and Y axis represents the normalized number of nodes discovered by the algorithms (the normalization is done by dividing the result of each algorithm by the result of the worst algorithm). This can show how many unnecessary questions were asked. When **CrowdPlanr** and **CrowdPlanr**⁻ are compared to the **Greedy** we can see that on the higher levels **Greedy** performs better than **CrowdPlanr** and **CrowdPlanr**⁻. In this case the greedy strategy to go after the local maximum seems to be a good choice, however starting from the level 4, **Greedy**’s performance gets dramatically worse.

5.3 Real data evaluation

In this experiment we analyzed the number of nodes visited by each one of the algorithms when executed on a real-world data sets. The parameters of the algorithms were set to default values. The experiment was performed on two real-world data sets: trip planning and star wars watching order. The results are shown in Figure 6.

Here the Y axis represents the number of visited nodes. As we can see, the **CrowdPlanr** algorithm is slightly better than the **CrowdPlanr**⁻ algorithm and both are significantly better than the **Greedy** algorithm. We conclude from this experiment that the behavior of all the algorithms seems to be the same as on synthetic data.

6. RELATED WORK

Using the crowd as a source of knowledge, and for solving problems that computers fail to solve, has attracted much research in recent years [7]. The planning problem that we

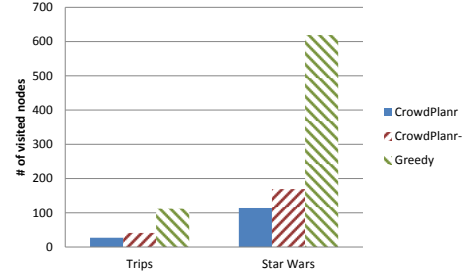


Figure 6: Real-world datasets

consider here is such an instance, as the goal, i.e, the notion of “best output”, is hard to formalize.

Much research was directed to find ways to effectively collect data from crowd, for example using games (e.g. [6], [26], [18]) or via payment (e.g. [20]). Others considered the development of unified model to allow uniform data collection from both humans and machines [22]. In particular, research has been directed in the Databases community to development of DB systems that allow to specify which parts of the data should be crowdsourced (e.g. CrowdDB [9], Deco [23], Qurk [19]). Crowdsourcing was also suggested as method for data cleaning, integration and analytics, entity resolution, schema expansion (e.g. [12], [24], [27], [16]). Crowdsourcing attracted also interest from the AI community with research aiming at dynamic workflow executions that optimally use the crowd for accomplishing a given complex task (e.g. [4], [15]). This is complementary to our work where users are used to identify and order the items (potentially the to-be-executed workflow components) needed to best accomplish an informally specified goal.

Minimization of the cost (measured in terms of the number of questions that are posed to the crowd) and of the expected error are important goals in crowd-based query processing ([2], [21]). Closest to our work are works that consider max and top-k query processing with the crowds, that involve ordering of query results using the crowd. For example, the problem of finding maximum element has been investigated in [12], considering how, given set of comparison results, one determines an element which is most likely to be the maximum, and which future comparisons will be most effective. Other examples are [25] that provides efficient tunable heuristics and [5] that studies complexity lower/upper bounds. As explained in the Introduction, a key difference between these previous works on max/top-k (and sorting, in general) processing and ours is the *inherent dependency that exists between items in the plan*. Unlike max/top-k processing where users can be asked to compare pairs of individual elements, planning requires a global view of the (preceding sub-)plan, and its possible/ideal completions. These previous algorithms are thus not applicable here. The optimal choice of questions to pose to the crowd has also been considered in [2] to reduce the uncertainty/error in aggregation functions over crowd answers. Here again a key difference from our work is the independence assumption among the aggregated data items. The dependency exhibited in planning problems requires developing corresponding uncertainty measures, and consequently different algorithms.

7. CONCLUSION

In this paper we propose to use the power of the crowd for answering planning queries, when the goal, i.e, the no-

tion of best plan, is hard to formalize. We introduce a simple generic model for modeling plans and for interpreting crowd's answers to questions about them. Based on this model, we present an effective algorithm for identifying the (approximately) best answer using the crowd. The algorithm builds the desired plans incrementally, choosing at each step the best questions so that the overall number of questions that need to be asked is minimized. We prove the algorithm to be instance-optimal for a large common class of planning queries and data instances, showing that the optimality ratio that it achieves is the best possible, and demonstrate experimentally the algorithm's effectiveness and efficiency.

We focused here on identifying the (approximated) best plan. More generally, one may want to identify top-k best answers. Our algorithm naturally generalizes to this context by continuing the execution after a top-1 sequence is found. Intuitively, nodes that are part of the returned sequence should be marked in the tree and ignored when candidates are considered. We omit the details for space constraints and only note that all of the results presented here for the Top-1 case (including optimality) extend for Top-K (full details can be found in the technical report [14]). An interesting challenge for future research is identifying heuristics that can be applied when some prior knowledge about the expected answer distribution, the tree structure, or the specific users expertise is available. How to obtain such information is also an interesting questions. Another possible extension to our algorithm, to be considered in the future research, could be to allow creating plans not necessarily in a successive order (for example when parts of the plan are known and one wants to use the crowd to fill in the gaps).

Acknowledgment

This work has been partially funded by the European Research Council under the European Community's Seventh Framework Programme (FP7/2007-2013) / ERC grant MoDaS, agreement 291071, by the Israel Ministry of Science, and by the US-Israel Bi national Science foundation.

8. REFERENCES

- [1] G. Adomavicius and A. Tuzhilin. Toward the next generation of recommender systems. *IEEE Trans. Knowl. Data Eng.*, 17(6):734–749, 2005.
- [2] R. Boim, O. Greenshpan, T. Milo, S. Novgorodov, N. Polyzotis, and W.-C. Tan. Asking the right questions in crowd data sourcing. In *ICDE*, pages 1261–1264, 2012.
- [3] G. Chen, C. Liu, M. Lu, B. C. Ooi, S. Ying, A. Tung, D. Zhang, and M. Zhang. A cross-service travel engine for trip planning. In *SIGMOD*, pages 1251–1254, 2011.
- [4] P. Dai, Mausam, and D. S. Weld. Decision-theoretic control of crowd-sourced workflows. In *AAAI*, 2010.
- [5] S. Davidson, S. Khanna, T. Milo, and S. Roy. Using the crowd for top-k and group-by queries. In *ICDT*, pages 225–236, 2013.
- [6] D. Deutch, O. Greenshpan, B. Kostenko, and T. Milo. Declarative platform for data sourcing games. In *WWW*, pages 779–788, 2012.
- [7] A. Doan, R. Ramakrishnan, and A. Y. Halevy. Crowdsourcing systems on the world-wide web. *Commun. ACM*, 54(4):86–96, 2011.
- [8] R. Fagin, A. Lotem, and M. Naor. Optimal aggregation algorithms for middleware. *J. Comput. Syst. Sci.*, 66(4):614–656, 2003.
- [9] M. J. Franklin, D. Kossmann, T. Kraska, S. Ramesh, and R. Xin. Crowddb: answering queries with crowdsourcing. In *SIGMOD*, 2011.
- [10] M. Ghallab, D. S. Nau, and P. Traverso. *Automated planning - theory and practice*. Elsevier, 2004.
- [11] R. Groves, F. J. Fowler, M. Couper, J. Lepkowski, E. Singer, and R. Tourangeau. *Survey Methodology*. John Wiley and Sons, 2009.
- [12] S. Guo, A. G. Parameswaran, and H. Garcia-Molina. So who won?: dynamic max discovery with the crowd. In *SIGMOD Conference*, pages 385–396, 2012.
- [13] P. Hart, N. Nilsson, and B. Raphael. A formal basis for the heuristic determination of minimum cost paths. *Systems Science and Cybernetics, IEEE Transactions on*, 1968.
- [14] H. Kaplan, I. Lotosh, T. Milo, and S. Novgorodov. Answering planning queries with the crowd (technical report) <http://slavanov.com/research/crowdplanr-tr.pdf>.
- [15] C. H. Lin, Mausam, and D. S. Weld. Crowdsourcing control: Moving beyond multiple choice. In *UAI*, pages 491–500, 2012.
- [16] X. Liu, M. Lu, B. C. Ooi, Y. Shen, S. Wu, and M. Zhang. Cdas: A crowdsourcing data analytics system. *PVLDB*, 5(10):1040–1051, 2012.
- [17] I. Lotosh, T. Milo, and S. Novgorodov. Crowdplanr: Planning made easy with crowd. In *ICDE*, 2013.
- [18] H. Ma, R. Chandrasekar, C. Quirk, and A. Gupta. Improving search engines using human computation games. In *CIKM*, pages 275–284, 2009.
- [19] A. Marcus, E. Wu, S. Madden, and R. C. Miller. Crowdsourced databases: Query processing with people. In *CIDR*, pages 211–214, 2011.
- [20] Amazon's mechanical turk. <https://www.mturk.com/>.
- [21] A. G. Parameswaran, H. Garcia-Molina, H. Park, N. Polyzotis, A. Ramesh, and J. Widom. Crowdscreen: algorithms for filtering data with humans. In *SIGMOD*, pages 361–372, 2012.
- [22] A. G. Parameswaran and N. Polyzotis. Answering queries using humans, algorithms and databases. In *CIDR*, pages 160–166, 2011.
- [23] H. Park, R. Pang, A. G. Parameswaran, H. Garcia-Molina, N. Polyzotis, and J. Widom. Deco: A system for declarative crowdsourcing. *PVLDB*, 5(12):1990–1993, 2012.
- [24] J. Selke, C. Lofi, and W.-T. Balke. Pushing the boundaries of crowd-enabled databases with query-driven schema expansion. *PVLDB*, 5(6):538–549, 2012.
- [25] P. Venetis, H. Garcia-Molina, K. Huang, and N. Polyzotis. Max algorithms in crowdsourcing environments. In *WWW*, pages 989–998, 2012.
- [26] L. von Ahn and L. Dabbish. Designing games with a purpose. *Commun. ACM*, 51(8):58–67, 2008.
- [27] J. Wang, T. Kraska, M. J. Franklin, and J. Feng. Crowder: Crowdsourcing entity resolution. *PVLDB*, 5(11):1483–1494, 2012.