

Optimal Testing of Structured Knowledge

Michael Munie and Yoav Shoham

Stanford University, CA

{munie, shoham}@stanford.edu

Abstract

Adopting a decision-theoretic perspective, we investigate the problem of optimal testing of structured knowledge – the canonical example being a qualifying examination of a graduate student. The setting is characterized by several factors: examinee’s knowledge structured around several inter-dependent topics, a limited “budget” of questions available to the examiner, a decision to be made (pass/fail), and an utility for good and bad decisions. The existence of multiple professors brings up additional issues such as committee formation, and the existence of multiple students brings up issues such as fairness.

1. Introduction

At Stanford University, many different formats of qualifying examinations exist, varying among departments as well as among groups within a given department. In some cases, each student is given an oral examination by a three-member committee for ninety minutes. In other cases, the students are all given identical written exams, lasting about two hours. In yet other cases, each student moves from professor to professor for a one-on-one exam lasting eight minutes each. Which is the best format? We set out to tackle this question scientifically, adopting a decision-theoretic perspective.

The problem is characterized by the following features: (a) structured knowledge being tested (e.g., knowledge of different subfields), (b) limited time of professors (and, less critically, of students), (c) a decision (i.e., whether or not to pass the student) in whose service the exam takes place, and (d) a utility function associated with good and bad decisions. Additional features have to do with the multiagent aspects of the problem. For example, the issues of committee formation and fairness.

Some of these features are quite generic and apply in other domains where value of information plays an important role, including sensor networks, medical diagnosis, and market research. We will stick to the student examination storyline, both for concreteness and because certain features make particular sense here. Nonetheless, work in other areas is rele-

vant, and we will review its similarities and differences after we set up our formal model.

The paper is organized as follows. In Section 2, we present the basic model; intuitively, it corresponds to a single professor examining a single student via a written examination and expecting him to demonstrate knowledge of a certain fraction of the topics. With that we’re in a position to discuss related past work more meaningfully in Section 3. In Section 4, we examine the difficulty of optimal questioning in the basic model (it is NP-hard), and its approximation (one trivially gets a factor-2 approximation, but one can not do any better), and then present a greedy algorithm which is provably optimal in certain restricted models. In Section 5, we look at moving from a written examination to an oral one and show that the latter can be maximally better. In Section 6, we look at what happens in the model if the basic threshold utility model is replaced by a more continuous one. Finally, in Section 7, we look at two sample problems in the multiagent case. In the case of multiple professors, we prove that having committees can be arbitrarily better or worse than having none, and that deciding between the two is NP-hard. In the case of multiple students we show that requiring fairness can lead to the maximally possible degradation of decisions regarding the individual students.

2. The Basic Model

In this section we present the basic formal model; in later sections we will modify and extend it in different ways. It consists of the following elements:

- \mathcal{U} , a set of knowledge of nodes;
- \mathcal{S} , a set of question nodes disjoint from \mathcal{U} ;
- \mathcal{B} , a Bayesian network over $\mathcal{U} \cup \mathcal{S}$ such that for all $S \in \mathcal{S}$, $Children(S) = \emptyset$ and $Parents(S) \subseteq \mathcal{U}$;
- $b \in \mathbb{R}^+$, a budget of questions;
- $c : \mathcal{S} \mapsto \mathbb{R}^+$, the cost function for observing a node;
- $\mathcal{D} = \{0, 1\}$, the set of possible decisions;
- \mathcal{A} , a design space, defined below; and
- $u : (\mathcal{D}, Val(U)) \mapsto \mathbb{R}$, a utility function defined below.

The intended interpretation of these elements is as follows. \mathcal{B} represents the student’s knowledge of the various

subareas and the various questions he can be asked. The nodes – both knowledge and question nodes – are assumed to be binary, so there are no degrees of knowledge, and answers to questions are either correct or incorrect. The design space \mathcal{A} is simply all the possible sets of questions the professor can ask whose combined cost does not exceed b . This models a written exam, which can be thought of as a professor asking a set of questions, and then receiving the set of answers in parallel. We will therefore refer to this as the parallel mechanism. After asking a set of questions and seeing the student’s response, the professor can then make a decision, here assumed to be binary (pass/fail). The utility for making a decision is the threshold function, parameterized by $0 \leq \theta \leq 1$ ¹:

$$u_\theta(d, u) = \begin{cases} 1 & \text{if } d = 1 \ \& \ |\{x \in u : x = 1\}| \geq \theta \ |\mathcal{U}| \\ 1 & \text{if } d = 0 \ \& \ |\{x \in u : x = 1\}| < \theta \ |\mathcal{U}| \\ 0 & \text{otherwise} \end{cases}$$

Given the basic model, there are two possible problems to tackle: what are the questions to ask, and based on the answers what is the decision to take so as to maximize utility? The answer to the second one is well known in the literature. For a binary decision D based on a binary hypothesis H (in this case, if the student knows at least $\theta|\mathcal{U}|$ areas, or not), it is shown by Heckerman *et al.* in (Heckerman, Horvitz, and Middleton 1993) what the decision point p' is such that for $P(H|Observations) \geq p'$ we will decide one way, and for values less than p' we will decide the other way. For our $u_\theta(\cdot)$, $p' = .5$. For a more in-depth background, see (DeGroot 2004). This leaves the problem of what questions to ask. We return to this question after we discuss related work.

3. Related Work

The literature that is potentially relevant to our problem is vast. It can be crudely divided into three classes. The literature on computer adaptive testing is closest in terms of the motivating application but rather different in terms of the models and problems studied. The literature in sensor nets (and related work in social networks) on optimal subset selection is closest technically, but its model makes a key assumption that does not apply in our domain. In addition, it does not ask as wide a set of questions as we do here. Finally, there is the much broader literature on value of information. It provides valuable insights on which we build but does not address the specific questions we ask. We discuss each of these in more detail below.

Most directly related to our model and application is the field of computer adaptive testing. Almond and Mislevy have proposed models based on Bayesian networks for educational testing in (Madigan and Almond 1993), (Robert J. Mislevy and Steinberg 1999), (Almond *et al.* 2002), and (Mislevy 2003) to make approximate inferences about hidden proficiency variables using Markov Chain Monte Carlo and also to provide feedback to the student. They focus on maximizing the expected weight of evidence instead of

¹We note that when $\theta = 1$, we have positive utility for passing the student only when the AND of the \mathcal{U} nodes equals 1.

working directly with the utility function as we do in this paper and don’t consider the complexity of the problem.

Recent work in AI has looked at finding optimal subsets to target in sensor and social networks. Guestrin has shown a $(1 - 1/e - \epsilon)$ approximation algorithm for observing subsets of sensors in (Krause and Guestrin 2005a) and optimal selection algorithms for restricted networks in (Krause and Guestrin 2005b). Unfortunately, this work leverages the fact that they are maximizing a submodular function (information gain of the \mathcal{U} nodes) to produce their results. Our utility functions are not submodular so we are not able to use their approximation. Kempe studies related subset selection problems in (D. Kempe and Tardos 2003) and argues for explicitly modeling utility for sensor networks in (Bian, Kempe, and Govindan 2006), though their model is substantially different from ours.

On a more basic level, we are solving a value of information problem. These problems have been studied for some time with early work being done by (Quinlan 1986) on greedily (myopically) building decision trees. Sequentiality was discussed in (Pearl 1988) but he approached the problem by using heuristics to maximize quasi-utility functions like entropy and didn’t focus on the computability of the problem. (Gaag 1993) outlined a framework covering selection of observations for diagnostic applications. In (Heckerman, Horvitz, and Middleton 1993), the problem of non-myopically selecting optimal observations was studied and a heuristic to select such observations was proposed, but without any guarantees. Jensen presented a means of calculating myopic value of information in influence diagrams in (Dittmer and Jensen 1997). In our case, influence diagrams quickly become unwieldy even for small networks. (Zubek and Dietterich 2005) takes a completely different approach and studies learning diagnostic policies from training data.

Also related is the field of probabilistic diagnosis studied in (Rish *et al.* 2004) and (Rish *et al.* 2005) (among others) where they attempt to find the location of a fault in, for example, a distributed system of computers with noisy observations. However, their focus is on fault detection and localization and not directly applicable to our binary decision problem. (Zheng, Rish, and Beygelzimer 2005) uses entropy as a cost function and shows how to efficiently compute approximate marginal conditional entropy using an algorithm based on loopy belief propagation. Our paper focuses on the orthogonal direction of selecting a set of observations when inference is tractable, instead of on the inference itself and approximately computing the next best myopic observation.

4. Analyzing the Basic Model

Even when modeling a single professor writing a test for a single student, the model exhibits many complexities. We will examine the hardness of selecting a set in the basic model, and how to solve the problem in a restricted setting with a greedy algorithm when approximation turns out to be difficult.

Complexity

As we noted in Section 2, making an optimal decision after making our observation in the parallel mechanism is not dif-

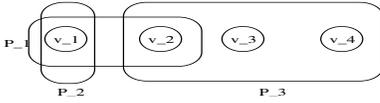


Figure 1: Set Cover

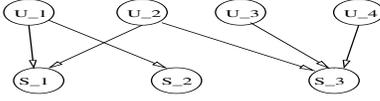


Figure 2: Network for reduction from SET-COVER

ficult. Now we can ask how hard is it to compute the optimal observation set. The answer is that, in general, it is hard:

Theorem 1. *Choosing an action under the parallel mechanism that has at least expected utility w , with $\theta = 1$, is NP-Hard even if inference can be performed efficiently through a belief update oracle.*

Proof. We prove this by a reduction from SET-COVER. In SET-COVER, you are given a set V , subsets P_1, P_2, \dots, P_m , and a budget k . The goal is to find at most k sets P_i such that their union equals V . We will convert an instance of SET-COVER as shown in Figure 1 to an instance of our problem as show in Figure 2. For each element $v_i \in V$ create a node U_i whose value is drawn uniformly at random from $\{0, 1\}$. We create nodes S_1, S_2, \dots, S_m that correspond to the sets P_1, P_2, \dots, P_m . If the set $P_i = \{v_{a_1}, v_{a_2}, \dots, v_{a_{|P_i|}}\}$, then the node $S_i = U_{a_1} \wedge U_{a_2} \wedge \dots \wedge U_{a_{|P_i|}}$.

Suppose we have an efficient algorithm to compute if a subset A exists with expected utility of w . Let R be defined as $R = U_1 \wedge U_2 \wedge \dots \wedge U_{|U|}$. If we let $w = 1$, then our algorithm will select a set with expected utility equal to $\sum_{a \in \text{Val}(A)} P(a) \text{Max}(P(R|a), 1 - P(R|a)) = 1$. If we observed any zeros at all, then $1 - P(R|a) = 0$ no matter which nodes we observed. Let α be the number of U nodes where we didn't observe any of their children. If we observe all ones, $P(R|a) = 2^{-\alpha}$. Therefore, to have utility of one, every U_i must have an observed S_j as a child and consequently $\bigcup_{i: S_i \in A} P_i = U$ and we have a solution for SET-COVER. \square

Although we proved NP-Hardness for the specific value $\theta = 1$, we can expand the proof to any θ .

Corollary 2. *For any $\theta \in (0, 1]$, choosing a set of questions A with expected utility w is NP-Hard even if inference can be performed efficiently through a belief update oracle.*

Proof. The proof is very similar to the proof of Theorem 1 and follows from a reduction from set cover. See the full version of this paper for more details. \square

Approximation

Given the NP-hardness of the problem, the natural next question is whether the optimal solution can be approximated within some constant. We first note that the technique

used in (Krause and Guestrin 2005a) does not apply, since it relies on the assumption of submodularity, which fails in our case:

Claim 3. *In a network with $\theta = 1$, maximizing the expected utility is not submodular.*

Proof. See the full version for a simple network where submodularity doesn't hold. \square

However, in our case we observe that we obtain “for free” a simple approximation. For an instantiation u of the \mathcal{U} nodes, let the expression $\#(u)$ denote $|\{x \in u : x = 1\}|$. If the prior on $\#(u) \geq \theta|\mathcal{U}|$ is $\sigma \in [0, 1]$ then observe that the simple rule that passes the student iff $\sigma \geq .5$, yields an expected utility that of $\text{Max}\{\sigma, 1 - \sigma\}$, giving us a 2-approximation in the worst case. This leaves open the question whether one can achieve an even better approximation, which we answer now.

Theorem 4. *For any $\theta \in (0, 1)$, choosing a set of observations A under the parallel mechanism is not approximable in polynomial time within a factor better than 2 unless $P=NP$.*

Proof. We take any instance of 3-SC (Garey and Johnson 1990) (SET-COVER where all sets are of size at most three) as shown in Figure 1 and convert it to a network G where, if we are able to select a set of observations that approximates the optimal expected utility closer than a factor of two in polynomial time, we will have solved the generating 3-SC instance. For each element $v_i \in V$, create a node U_i whose value is drawn uniformly at random from $\{0, 1\}$. We create nodes S_1, S_2, \dots, S_m that correspond to the sets P_1, P_2, \dots, P_m and whose parents are nodes U_j such that $v_j \in P_j$. In addition we create nodes $S_{m+1}, S_{m+2}, \dots, S_{m+l}$ that correspond to the non-empty proper subsets of each set in P_1, P_2, \dots, P_m and whose parents correspond to the nodes in that subset of P_j . Since in instances of 3-SC our sets are at most size three, then for each set P_i we will add at most seven S nodes. Let each S node be the exclusive or of its parents. In addition, we add a node labeled XOR that is the exclusive or of all the U nodes. The XOR node has $\frac{\text{max}(\theta, 1 - \theta)|V|}{1 - \text{max}(\theta, 1 - \theta)}$ $FILL$ nodes as children, whose values are deterministically equal to the value of the XOR node itself. The U , $FILL$, and XOR nodes are all the members of \mathcal{U} .

First, we note that a set-cover exists if and only if there is a set $A \subseteq S$, $|A| < b$, where the union of A 's parents equals U and the pairwise intersection between the parents of any two sets in A is empty. Also note that $P(XOR = 0|A) = .5$ unless we exactly know the exclusive or of the XOR node's parents in which case $P(XOR = 0|A) = 0$ or 1 . Because the value of the XOR node (and all of the $FILL$ nodes with the same value) uniquely determines if $\#(u) \geq \theta|\mathcal{U}|$, the expected utility given any possible valid set A , $\text{max}(P(\#(u) \geq \theta|\mathcal{U}| |A), P(\#(u) < \theta|\mathcal{U}| |A))$, can only take on the values 1 and .5. If the expected utility is 1, this must mean that we have observed the XOR node, and thus have observed a set of nodes that exactly reveals XOR , which must be a set that corresponds to a set cover in our original problem. Since any approximation algorithm with

a ratio better than two would have to guarantee an expected utility of 1, which also corresponds to exactly solving the 3-SC instance, there is no such algorithm unless $P=NP$. \square

A Greedy Algorithm

Even in the face of this approximation result, we find that by restricting the networks, we are able to find an optimal greedy algorithm.

There will be three restrictions on networks, which we will call a *simple networks*. First, the \mathcal{U} nodes will have no edges to other \mathcal{U} nodes. Second, every node in \mathcal{U} will have one and exactly one node in \mathcal{S} as a descendant and every \mathcal{S} node will have exactly one parent. The value of the nodes in \mathcal{S} will always be deterministically equal to the value of its parent. Finally, our utility function will have $\theta = 1$.

Our greedy algorithm is also easy to define. At every step, it will add the node x , where $x = \operatorname{argmax}_{u \in (\mathcal{U}/A)} P(u = 0)$ until the budget is used up and where x is not already in our selected set A . Let $E[u(A)]$ be the expected utility of choosing set A and then following our optimal decision rule based on our observed values over the possible instantiations of the nodes in the network.

Theorem 5. *In a simple network, the greedy algorithm is optimal.*

Proof. Although this claim might seem intuitively straightforward, the proof is surprisingly involved. See the full version. \square

5. From Written to Oral Exams

The other situation we will model is an interview where the professor hears the answer to each question before asking the next one. Formally our new sequential mechanism proceeds as follows. One node is observed. After updating the observation in the network, inference can be performed and the next node will be selected. Once the budget has been reached, inference can be performed again and the optimal decision will be selected. This is just a change in the design space \mathcal{A} from the parallel space used in the basic model to a sequential space.

Separation

The expected utility under the sequential mechanism is at least as good as the expected utility under the parallel mechanism. This is easily seen since in the worst case, the sequential mechanism can mirror the node selection of parallel mechanism. Given this, the natural question is how much better might the separation be? Recall that the utility is $\in [.5, 1]$ under our optimal decision rule, so the upper bound on the separation is .5. The following theorem shows how close we can approach the maximal separation.

Theorem 6. *For any $\theta \in (0, 1)$ and any budget b , the difference in expected utility between the sequential and parallel mechanisms can be at least $.5 - \frac{1}{2^b}$*

Proof. See the full version for a class of networks with this large separation. \square

Complexity

The same reduction used for the parallel mechanism applies here as the \mathcal{U} nodes in Theorem 1 are all independent of each other, so the sequentiality will make no difference.

6. A Smoother Utility Model

It might seem that the specific form of our threshold utility function underlies the difficulty of the problem, but now we will show that this phenomenon is more general by defining a new class of symmetric monotonic utility functions.

Definition 7. *Let $f_\theta(d, u)$ be a symmetric monotonic utility function when it satisfies the following restriction: $f_\theta(\text{Pass}, u) \geq f_\theta(\text{Fail}, u)$ for all $u \in \text{Val}(\mathcal{U})$ such that $|\{x \in u : x = 1\}| \geq \theta|\mathcal{U}|$ and otherwise $f_\theta(\text{Pass}, u) < f_\theta(\text{Fail}, u)$.*

It is easy to see that $u_\theta(\cdot)$ is a member of this class of utility functions, however there are many others. We now show that our hardness result still holds for any member of this class.

Theorem 8. *Choosing observations under the parallel and sequential mechanisms that has at least expected utility w as measured by any monotonic utility function $f_\theta(d, u)$, is NP-Hard even if inference can be performed efficiently through a belief update oracle.*

Proof. Due to space constraints, we present a proof sketch here. We reduce from SET-COVER by generating the same network we did in the proof of Theorem 1. The maximum utility we can receive if we made all correct decisions is $u_{max} = \sum_{u \in \text{Val}(\mathcal{U})} P(u) \max(f_\theta(0, u), f_\theta(1, u))$. After observing a set A with observations a , our optimal decision $d^*(a) = \operatorname{argmax}_{d \in D} (P(\#(u) \geq (\theta|\mathcal{U})|a), P(\#(u) < (\theta|\mathcal{U})|a))$. Then the utility of observing a set A is $\sum_{a \in \text{Val}(A), u \in \text{Val}(\mathcal{U})} P(a, u) f_\theta(d^*(a), u)$. If we had an algorithm to efficiently select A , we could set $w = u_{max}$. In our network, the only way we can achieve u_{max} is if we find a set cover. \square

This is a powerful result since most reasonable utility functions can be expected to be symmetric monotonic. In fact, there is also an analog of our hardness of approximation in the space of symmetric monotonic utility functions.

For the theorem, we define the following constants: let $\beta = \min(\theta, 1 - \theta)|\mathcal{U}|$,
 $FailLow = \max_{u \in \text{Val}(\mathcal{U}): (|\{x \in u: x=1\}| \leq \beta)} (f_\theta(0, u))$,
 $FailHigh = \max_{u \in \text{Val}(\mathcal{U}): (|\{x \in u: x=1\}| \geq |\mathcal{U}| - \beta)} (f_\theta(0, u))$,
 $PassLow = \max_{u \in \text{Val}(\mathcal{U}): (|\{x \in u: x=1\}| \leq \beta)} (f_\theta(1, u))$,
 $PassHigh = \max_{u \in \text{Val}(\mathcal{U}): (|\{x \in u: x=1\}| \geq |\mathcal{U}| - \beta)} (f_\theta(1, u))$,
and $z = \max(\frac{1}{2}FailLow + \frac{1}{2}FailHigh, \frac{1}{2}PassLow + \frac{1}{2}PassHigh)$.

Theorem 9. *For any symmetric monotonic utility function $f_\theta(\cdot)$ with $\theta \in (0, 1)$, choosing a set of observations A under the parallel mechanism is not approximable in polynomial time within a factor better than $\frac{u_{max}}{z}$ unless $P=NP$*

We can see that by applying this bound to our threshold utility function we get that it is NP-hard to approximate with a ratio better than 2.

Proof. We take any instance of 3-SC and convert it to a network \mathcal{B} where, if we are able select a set of observations that approximates the optimal expected utility closer than a factor of $\frac{u_{max}}{z}$ in polynomial time, we will have solved the generating 3-SC instance. We follow the reduction in the proof of Theorem 4 to generate \mathcal{B} .

First, we note that a set-cover exists if and only if there is a set $A \subseteq S$, $|A| < b$, where the union of A 's parents equals U and the pairwise intersection between the parents of any two sets in A is empty. Also note that $P(XOR = 0|A) = .5$ unless we exactly know the exclusive or of the XOR node's parents in which case $P(XOR = 0|A) = 0$ or 1 .

Now let us partition the possible solutions our algorithm can return into two groups. First, where we know XOR (i.e., $P(XOR = 0|A) = 0$ or 1), and second, where we don't know XOR (i.e., $P(XOR = 0|A) = .5$). If we have a solution where we know XOR, then we must have also returned a set cover.

If we don't know XOR, then our utility will be bounded by a constant dependent on our utility function. Let the maximum utility we can receive if we made all correct decisions be, $u_{max} = \sum_{u \in Val(\mathcal{U})} P(u) \max(f_\theta(0, u), f_\theta(1, u))$. Because our set of $FILL$ nodes take up a $\max(\theta, 1 - \theta)$ fraction of our total \mathcal{U} nodes, then for any $u \in Val(\mathcal{U})$ such that $P(u) \neq 0$, $|\{x \in u : x = 1\}| \leq \beta$ when $XOR = 0$ and $|\{x \in u : x = 1\}| \geq |\mathcal{U}| - \beta$ when $XOR = 1$. This, along with the fact that we have $P(XOR = 1) = .5$ implies,

$$\begin{aligned} & \sum_{u \in Val(\mathcal{U}): (|\{x \in u: x=1\}| \leq \beta)} P(u) \\ = & \sum_{u \in Val(\mathcal{U}): (|\{x \in u: x=1\}| \geq |\mathcal{U}| - \beta)} P(u) \\ = & P(XOR = 1) = P(XOR = 0) = .5 \end{aligned} \quad (1)$$

Thus, if we don't know XOR, then our utility will be bounded above by $z = \max(\frac{1}{2}FailLow + \frac{1}{2}FailHigh, \frac{1}{2}PassLow + \frac{1}{2}PassHigh)$. Since any approximation algorithm with a ratio better than $\frac{u_{max}}{z}$ would have to produce an expected utility of u_{max} , which also corresponds to exactly solving the 3-SC instance, there is no such algorithm unless $P=NP$. \square

7. The Multiagent Case

The multiagent setting is varied, and here we study it under the threshold utility function $u_\theta(\cdot)$. We can have multiple professors, multiple students, or both. We discuss a sample problem that arises in each of these settings: committees in the multiple professor case and fairness in the multiple student case. The case of having multiples of both does not appear to add much to our model.

Multiple Professors

Once we introduce multiple professors to the model, there are many factors that we can consider – sequentiality and committees among others. We have found that the sequential multiagent model with committees and two extreme forms of it are the most interesting cases, and that is what we will focus on in this section.

We capture these multiagent models by duplicating each node in S as many times as there are professors, attaching a separate distribution to each node, and assigning each copy to a distinct professor. All the nodes duplicated from the same original node will be said to form a question group. This represents a question that any of the professors could ask, but that they each draw different conclusions from. Each professor has a budget of $\frac{b}{|P|}$, where P is the set of professors.

A partition of the set of professors captures our notion of committees, with the professors in the same partition working together. All the committees operate independently and in parallel. Inside each committee, observations must be made by question group, not by individual node, i.e., all the nodes in the same question group that belong to professors inside the committee must be observed simultaneously, or none observed at all. A committee will observe question groups until its budget is exhausted. After all the nodes in a particular question group have been observed, the choice of the next question group may be conditioned on the results of previous observations by the same committee. The results of all observations are then used to make a pass/fail decision.

The trade-off between using large committees that select the next question better, and fewer committees that query more question groups is our focus. In order to quantify the power of committees, we examine two settings—one where every professor is a member of a singleton committee and the other where we have a single grand committee consisting of all the professors. It is surprising that no setting is uniformly better; in fact, the difference can be maximally large:

Theorem 10.

- Given any budget b and set of professors P and any u_θ utility function, There are networks for which the grand committee outperforms singleton committees : the grand committee yields an expected utility of 1, the highest possible, and the singleton committees yield the expected utility of $.5 - \frac{|P|}{2^{b-|P|+1}}$.
- Given any budget b and set of professors P and any u_θ utility function, There are networks for which singleton committees outperform the grand committee, and do so by the widest possible margin: the singleton committees yield an expected utility of 1, the highest possible, and the grand committee yields the expected utility of the trivial mechanism which decides based on the prior.

Proof. See the full version for networks that exhibit these separations. \square

Although we know that finding optimal solutions for sequential multiagent mechanism is NP-hard because it is a generalization of the single agent case, we would still like to know which one has higher optimal expected utility. This can be thought of as a meta-decision, a decision on which mechanism to use to make our decision on the student. We show this even this is hard.²

²When we write $EU(X)$, where X is a mechanism, we are referring to the maximal value of this, $\max_{A \subseteq S} [EU(X, A)]$

Theorem 11. *On a specific network with sequential questioning, deciding if $EU(\text{grand committee}) \geq EU(\text{singleton committees})$ is NP-hard for u_θ utility functions.*

Proof. The proof here is fairly complicated, but is also a reduction of SET-COVER. We construct a network so that the sequential grand committee has expected utility of 1 if a set cover exists and .5 otherwise. The sequential singleton committees always have a utility slightly above .5, and thus if we had an algorithm to decide which mechanism is better, we would solve set cover. See the full version for more details. \square

Multiple Students

It also makes sense to consider examining multiple students. We will assume that these students are i.i.d. according to our original network \mathcal{B} . Our utility function is just the summation of the utility for the individual students – now equal to the expected number of students judged correctly. A simple way of running this mechanism would be to have each professor spend the same amount of time on each student. In this case it is easy to see that each student would have the same expectation of being judged correctly. In fact, this is our definition of fairness with multiple students.

Definition 12. *If we have a network \mathcal{B} , budget b , and n students all drawn independently from this distribution, then the vector of observation sets $\mathbf{A} = \{A_1, A_2, \dots, A_n\}$ where $A_i \subseteq \mathcal{S}_i$ is fair if for all i, j , $H(D_i|A_i) = H(D_j|A_j)$, where $H(\cdot)$ computes the entropy.*

This definition captures the fact that we would like every student to know that the decision on their performance had the same expected amount of randomness. Given that fairness is well motivated, we would like to know the potential cost of fairness, i.e., what is the difference between the optimal set A with fairness and the optimal set without fairness.

Clearly, fairness can only hurt the optimal score, but by how much can it hurt?

Theorem 13. *The utility of an optimal observation set A can be twice that of a set A that is fair.*

Proof. Let \mathcal{B} be a network where each student is represented by a single internal node with a .5 chance of being 1 with utility function $u_1(\cdot)$. For $b = n - 1$, if A is fair, then A must be the empty set with utility $\frac{n}{2}$. The optimal A without fairness has utility $n - 1/2$, which is asymptotically close to twice as good as a fair A . \square

8. Summary

Motivated by a real-world problem of optimal testing of structured knowledge, we have developed a precise model in which to investigate this and related questions. Novel elements of the model include the fact that the information is geared towards making a specific decision, that utility is expressed by a wide range of prescriptive utility models, and that observations can be made in parallel or in sequence by individuals or groups. Among other things, we have shown the NP-hardness of deciding on the optimal question set even for the individual agent, the hardness of approximating this problem, and a greedy algorithm for solving it in

a restricted network. In addition, we have shown the non-comparability between multi-agent mechanisms with and without committees, and that it is hard to determine the optimal committee structure for a given network. This problem has many theoretical aspects which have yet to be explored, and being closely tied to practical problems in computer adaptive testing and sensor networks, we hope to see more work in this field in the future.

References

- Almond, R.; Mislevy, R.; Steinberg, L.; Breyer, F. J.; and Johnson, L. 2002. Making sense of data from complex assessments. *Applied Measurement in Education*.
- Bian, F.; Kempe, D.; and Govindan, R. 2006. Utility based sensor selection. In *Proc. of IPSN '06*, 11–18. New York, NY: ACM.
- D. Kempe, J. K., and Tardos, E. 2003. Maximizing the spread of influence through a social network. *SIGKDD*.
- Degroot, M. H. 2004. *Optimal Statistical Decisions (Wiley Classics Library)*. Wiley-Interscience.
- Dittmer, S. L., and Jensen, F. V. 1997. Myopic value of information in influence diagrams. In *In Proc. of UAI-13*, 142–149.
- Gaag, L. v. d. 1993. Selective evidence gathering for diagnostic belief networks. *RUU-CS (Ext. rep. 93-31)*.
- Garey, M. R., and Johnson, D. 1990. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. New York, NY: W. H. Freeman and Co.
- Heckerman, D.; Horvitz, E.; and Middleton, B. 1993. An approximate nonmyopic computation for value of information. *IEEE Trans. Pattern Anal. Mach. Intell.* 15(3):292–298.
- Krause, A., and Guestrin, C. 2005a. Near-optimal nonmyopic value of information in graphical models. In *Proc. of UAI-21*.
- Krause, A., and Guestrin, C. 2005b. Optimal nonmyopic value of information in graphical models - efficient algorithms and theoretical limits. In *Proc. of IJCAI-19*.
- Madigan, D., and Almond, R. 1993. *Test Selection Strategies for Belief Networks*. Learning from Data: AI and Statistics V. Springer-Verlag. chapter 9, 88–97.
- Mislevy, R. J. 2003. On the structure of educational assessments. *Measurement: Interdisciplinary Research and Perspectives*.
- Pearl, J. 1988. *Probabilistic Reasoning in Intelligent Systems*. San Mateo, CA: Morgan Kaufmann.
- Quinlan, J. R. 1986. Induction of decision trees. *Mach. Learn.* 1(1):81–106.
- Rish, I.; Brodie, M.; Odintsova, N.; Ma, S.; and Grabarnik, G. 2004. Real-time problem determination in distributed systems using active probing. *NOMS'04* 1:133–146.
- Rish, I.; Brodie, M.; Ma, S.; Odintsova, N.; Beygelzimer, A.; Grabarnik, G.; and Hernandez, K. 2005. Adaptive diagnosis in distributed systems. *Neural Networks, IEEE Transactions on* 16(5):1088–1109.
- Robert J. Mislevy, Russell G. Almond, D. Y., and Steinberg, L. 1999. Bayes nets in educational assessment: Where the numbers come from. In *Proc. of UAI-15*.
- Zheng, A. X.; Rish, I.; and Beygelzimer, A. 2005. Efficient test selection in active diagnosis via entropy approximation. *Proceedings of UAI-05*.
- Zubek, V. B., and Dietterich, T. G. 2005. Integrating learning from examples into the search for diagnostic policies. *JAIR* 24:263–303.