

Game-Theoretic Problem Selection for Tests

Abstract

Conventionally, the problems on a test are assumed to be kept secret from test takers until the test. However, for tests that are taken on a large scale, particularly asynchronously, this is very hard to achieve. For example, example TOEFL iBT and driver's license test problems are easily found online. This also appears likely to become an issue for Massive Open Online Courses (MOOCs).

In this paper, we take the loss of confidentiality as a fact. Even so, not all hope is lost as the test taker can memorize only a limited set of problems, and the tester can randomize which problems appear on the test. We model this as a Stackelberg game, where the tester commits to a mixed strategy and the follower responds. We provide an exponential-size linear program formulation, prove several NP-hardness results, and give efficient algorithms for special cases.

1 Introduction

2 General Test Game

A test game $G = (A, \mathbb{B}, p, v, m, t)$ is a 2-player Bayes game between the tester and the test taker. The tester has a set of potential problems $A = \{a_1, a_2, \dots, a_n\}$ to test and the test taker has an uncertain type (thus Bayes) characterized by $\mathbb{B} = \{B_1, B_2, \dots, B_L\}$. Set $B_l \subseteq A$ is the unsolvable problem set of the l -th type test taker. Function $p : \mathbb{B} \rightarrow [0, 1]$ characterizes the probability that a particular type test taker occurs. The memorize capacity function $m : \mathbb{B} \rightarrow \mathbb{N}$ denotes how many problems a particular test taker can memorize and the value function $v : \mathbb{B} \rightarrow \mathbb{R}^+$ denotes how much utility the tester gets if that test taker failed the test (recall that the tester can only fail test takers who cannot solve all problems and he want to fail as many of them as possible). For simplicity, we write $p(B_l), m(B_l), v(B_l)$ as p_l, m_l, v_l when the context is clear. During the test, t out of n problems will be tested and a particular test taker will pass the test if all tested problems are either not in the unsolvable problem set or memorized.

The game only has two outcomes, pass or fail. The tester gets 0 utility if a test taker passes and he gets v_l utility if a type- l test taker fails. On the type- l test taker's side, he gets

0 utility if he passes and $-v_l$ if he fails. From the test taker's perspective, any utility function that has a higher pass utility than fail utility will give him the same incentive to pass the test at his best. We define them as 0, $-v_l$ specifically for the zero-sum property of the game.

Our goal is to find the optimal strategy for the tester to maximize his utility under Stackelberg settings. That is, the tester firstly reveals his strategy about how t problems are picked to test, then the test taker chooses the best memorizing strategy. This corresponds to our confidential-free assumption where test takers know what problems might be tested, their answers, and how likely they are tested. The only constraint for test takers is that they cannot memorize too many (greater than m_l) problem answers.

2.1 General Linear Program (LP) Formulation

General 2-player Bayes Stackelberg games are NP-hard in terms of game matrix size so previous work used mixed integer LP (e.g. DOBSS [cite]) to solve them. We model test games as zero-sum so we can bypass the NP-hardness and use the following maximin LP to solve our 2-player zero-sum Bayes Stackelberg games in polynomial time with respect to the game matrix size:

$$\begin{aligned} \max \quad & \sum_l p_l V_l \\ \text{s.t.} \quad & (\forall l, \forall c_l \in C_l), \\ & \sum_r x_r u^l(r, c_l) \geq V_l \\ & \sum_r x_r = 1 \\ & (\forall r \in R) x_r \geq 0 \end{aligned} \tag{1}$$

where

- V_l : utility that the test taker can get from type l test takers
- C_l : action space of type l test takers, i.e. all combinations of problems that they can memorize.
- R : action space of the tester, i.e. all combinations of problems that he can test.
- x_r : probability to test problem set r

The dual of the above LP is:

$$\begin{aligned}
\min U & \quad (2) \\
s.t. (\forall r) U & \geq \sum_{l, c_l} u^l(r, c_l) y_{l, c_l} \\
(\forall l) \sum_{c_l} y_{l, c_l} & = p_l \\
(\forall l, c_l) y_{l, c_l} & \geq 0
\end{aligned}$$

which is equivalent to

$$\begin{aligned}
\min \max_r (U_r) & \quad (3) \\
s.t. (\forall r) U_r & = \sum_{l, c_l} u^l(r, c_l) y_{l, c_l} \\
(\forall l) \sum_{c_l} y_{l, c_l} & = p_l \\
(\forall l, c_l) y_{l, c_l} & \geq 0
\end{aligned}$$

The dual LP is as if that all types of test takers share the same goal to lower the tester's best testing utility ($\max_r(U_r)$) where U_r is the utility of testing r) as much as possible by memorizing problems and revealing their memorizing strategies y_{l, c_l} to the tester. This also shows that a Stackelberg strategy is also a Nash equilibrium if test takers play the dual minimax strategy.

However, even if the LPs above are in polynomial size of the game matrix size (or the action space size), they are exponential to the input of our test games: the tester's action space is exponential to t (the number of problems to test) and the type- l test taker's action space is exponential to m_l (the number of problems he can memorize). So those LPs only give us polynomial algorithms when m_l, t are constant. Our next question is whether the test game in general can be solved in polynomial time in terms of the input size.

2.2 Hardness on Test Size

In order to formally characterize *test game's* hardness, define

Definition 1 (OPTIMAL TEST STRATEGY Problem). *Given a test game G and a value u , the OPTIMAL TEST STRATEGY problem is to decide whether the tester has a Stackelberg strategy with at least u utility.*

It's easy to show

Theorem 1. *Even if test takers cannot memorize any problems ($m_l = 0$), OPTIMAL TEST STRATEGY is NP-hard when test size t is non-constant.*

Proof. Reduce a VERTEX COVER instance to a OPTIMAL TEST STRATEGY instance as follows. Given a graph $G = (V, E)$, construct a test game G' with problem set $A = V$. For each edge $e = \{i, j\} \in E$, add one type of test takers l_e whose unsolvable problem set $B_{l_e} = e$. Let $p_{l_e} = 1/|E|$, $v_{l_e} = 1$ and $m_{l_e} = 0$ for all l_e . Graph G has a vertex cover of size k if and only if the tester has a strategy to test $t = k$ problems and gets $u = 1$ utility (i.e. all test takers will fail for sure) in test game G' . \square

2.3 Hardness on Memorize Capacity

Theorem 2. *Even if test size is 2, OPTIMAL TEST STRATEGY is coNP-hard when memorize capacity is non-constant.*

Proof. We reduce INDEPENDENT SET instances to test games with test size $t = 2$ and ask the complementary question: is u the highest utility that the tester can get. That question has a yes answer if and only if the dual minimax LP [reference] has a feasible solution with objective at most u . For a graph $G = (V, E)$, let V be our test problem set A . Add $L = |E| + |V| + 2$ types of test takers as follows:

- For each edge $e = (i, j) \in E$, add one type of test takers l_e with $v_{l_e} = 1$, $B_{l_e} = \{i, j\}$ and $m_{l_e} = 0$.
- For each vertex i , add one type of test takers l_i with $v_{l_i} = d_G - d(i)$, $B_{l_i} = \{i\}$ and $m_{l_i} = 0$. Here d_G is the maximum vertex degree in G and $d(i)$ is the degree of vertex i .
- Add one auxiliary type of test takers l_α with $v_{l_\alpha} = \alpha\varepsilon$, $m_{l_\alpha} = 2$ and $B_{l_\alpha} = V$ where $\alpha = \binom{|V|}{2} - |E| - \binom{k}{2}$ and $\varepsilon \ll 1$.
- Finally, add one type of test takers l_k with $v_{l_k} = \varepsilon$, $B_{l_k} = V$ and $m_{l_k} = k$.

Let all types of test takers occur with uniform probability $p = 1/L$. Then we show that deciding whether $u = (2d_G + a \cdot \varepsilon)/L$ is an upper-bound for the tester's utility, or equivalently whether we can find a dual solution with objective as low as $u = (2d_G + a \cdot \varepsilon)/L$ in our dual minimax LP, is equivalent to finding a size- k independent set in G . Recall that a dual solution is a strategy for test takers to memorize problems. If no one memorizes, the tester's utility U_r to test a problem set r is $U_r = (2d_G + a \cdot \varepsilon + \varepsilon)/L$ if $r \notin E$ and $U_r = (2d_G - 1 + a \cdot \varepsilon + \varepsilon)/L$ if $r \in E$. To lower the objective, test takers l_α, l_k should memorize problem sets that bring the maximum of U_r down. No matter what the strategy is (as long as they only memorize unsolvable problems), $\sum_r U_r$ is a constant. And because U_r for $r \notin E$ is much larger than U_r for $r \in E$ as $\varepsilon \ll 1$, it's better to only decrease U_r for $r \notin E$. The only way to do that is to let l_α test takers only memorize pairs of problems that are not edges and let l_k test takers only memorize size- k independent sets of problems. If there's a size- k independent set, let l_k test takers memorize one size- k independent set all the time and let l_α test takers memorize all pairs of problems that are not covered by that independent set with uniform probability. That will make $U_r = (2d_G + a \cdot \varepsilon)/L$ for $r \notin E$ and $U_r = (2d_G - 1 + a \cdot \varepsilon + \varepsilon)/L$ for $r \in E$. This is the best they can achieve and this can only be achieved when a size- k independent set exists in G . \square

2.4 Summary on General Test Games

Theorem 3. *The OPTIMAL TEST STRATEGY for a test game can be computed in polynomial time when test size and memorize capacity are constant. It's NP-complete when test size is non-constant and coNP-complete when memorize capacity is non-constant. Hence it's both NP-hard and coNP-hard when both test size and memorize capacity are non-constant.*

Proof. When test size and memorize capacity are both non-constant, OPTIMAL TEST STRATEGY is both NP-hard and coNP-hard by Theorem 1 and 2. LP (1) shows that OPTIMAL TEST STRATEGY can be computed in polynomial time when test size and memorize capacity are constant. When memorize capacity is constant or test size is constant, either the primal LP (1) or the dual LP (3) has constant number of constraint. Such LP has a polynomial size optimal solution even if the number of variables is exponential. That optimal solution's feasibility and objective can be checked in polynomial time. Therefore OPTIMAL TEST STRATEGY is in NP when memorize capacity is constant and it is in coNP when test size is constant. Hence they are NP-complete and coNP-complete respectively by Theorem 1 and 2. \square

[HARDNESS TABLE]

3 One-Problem Test Game

We showed the hardness of test games when either one of test size or memorize capacity is non-constant. When the test size is non-constant, there is little we can improve because it is NP-hard even if the memorize capacity is zero. When the memorize capacity is non-constant, our proof for Theorem 2 only showed that NP-hardness for testing 2 problems. Therefore an efficient algorithm to compute optimal one-problem test may exist. A quick evidence for the existence of such polynomial algorithm is constraint generation. Though there are exponential many constraints, we can find the most violated constraint in polynomial time by: enumerate over all types of test takers; for type- l test takers select top m_l problems M_l according to x_a that maximize $\sum_{a \in M_l} x_a$; finally see whether $V_l - \sum_{a \in B_l \setminus M_l} x_a$ updates the most violated constraint.

We also investigated such tests by conducting experiments using LP (1). Surprisingly, no matter how input varies, the strategy we get always has the following format: among a subset T of problems A , test any one of them with uniform probability. That is counter-intuitive as some problems in T seem to be obviously more superior than others so intuitively they should be tested with more probability. For example, when n problems can be sorted by difficulty, a test taker who can solve a harder problem can always solve an easier problem. In that case, it seems that the hardest problem should be tested strictly more likely than the second hardest problem. But in most cases, the optimal strategy will test one of them with equal probability.

[EXAMPLE TABLE]

Inspired by that uniform test property, we present an algorithm to compute optimal strategies and use it to prove uniform test property consequently. Let U_a^0 be the utility of testing a without memorizing: $U_a^0 = \sum_{B_l \ni a} v_l p_l$. The algo-

rithm's first part is based on the following LP:

$$\begin{aligned} \min \quad & U \\ \text{s.t.} \quad & (\forall a \in A) U \geq U_a^0 - \sum_{B_l \ni a} p_l v_l z_{l,a} \\ & (\forall l) \sum_{a \in B_l} z_{l,a} \leq m_l \\ & (\forall l, a) 0 \leq z_{l,a} \leq 1 \end{aligned} \tag{4}$$

It modifies LP (2) by:

1. Change joint probability y_{l,c_l} to marginal probability $z_{l,a}$ of memorizing one problem a . We did this modification because marginal probability is more meaningful for one-problem test and we can always restore a joint probability from a marginal probability[reference].
2. Relax equality $\sum_{a \in B_l} z_{l,a} = m_l$ to inequality $\sum_{a \in B_l} z_{l,a} \leq m_l$ (allow using less memorize capacity). We can do this because $v_l \geq 0$ (thus $U^l, V_l \geq 0$) and this modification is essentially adding $V_l \geq 0$ to the primal LP (1). Intuitively, memorizing less problems won't hurt tester's utility.

Our algorithm is shown in Algorithm 1.

Algorithm 1 Compute the optimal one-problem test strategy

```

1: function OPTIMALONEPROBLEMTTEST( $A, \mathbb{B}, m, v, p$ )
2:   Solve LP (4) and let the solution be  $U, z_{l,a}$ 
3:    $T \leftarrow \{a \mid U_a^0 \geq U\}$ 
4:    $Q \leftarrow$  all  $B_l$  such that  $\sum_{a \in B_l \cap T} z_{l,a} < m_l$ 
5:   for all  $B_l \in Q$  do  $\triangleright$  Potential Memorizing Power
6:      $G(B_l) \leftarrow m_l - \sum_{a \in B_l \cap T} z_{l,a}$ 
7:   end for
8:   while  $Q$  has an unmarked element do
9:      $B_l \leftarrow$  get and mark an unmarked element from  $Q$ 
10:    for all  $a \in B_l \cap T$  and  $z_{l,a} < 1$  do
11:       $T \leftarrow T \setminus \{a\}$ 
12:      for all  $B_{l'} \ni a$  and  $z_{l',a} > 0$  do
13:         $G(B_{l'}) \leftarrow \min(z_{l',a}, \frac{G(B_l)}{|B_l|L})$ 
14:        Push  $B_{l'}$  into  $Q$ 
15:      end for
16:    end for
17:  end while
18:  return Strategy that test  $T$  uniformly
19: end function

```

We prove the correctness of Algorithm 1 by firstly introducing the following two lemmas.

Lemma 1. T returned by Algorithm 1 is non-empty.

Proof. T is non-empty before we delete problems from T on line 11 of Algorithm 1. Denote that initial T as T_0 . If all problems in T_0 are deleted, then for every a in T_0 there is a B_l that is in Q and $z_{l,a} < 1$. Note that for every B_l in Q , it has a positive potential memorizing power $G(B_l)$ either because $\sum_{a \in B_l} z_{l,a} < m_l$ (on line 6) or because there is an unnecessary memorizing $z_{l,a} > 0$ that can be freed (on line 13). This means that we can bring U down by using those

potential memorizing powers to memorize every problem in T_0 more (recall that $z_{l,a} < 1$), a contradiction to that U is minimized. \square

Lemma 2. *Let $z_{l,a}$ be the solution of LP (4) and T be the final subset that Algorithm 1 returns. Then for each type of test takers B_l , either $\sum_{a \in B_l \cap T} z_{l,a} = m_l$ (they memorize a subset of T that are unsolvable all the time) or $\forall a \in B_l \cap T, z_{l,a} = 1$ (they memorize unsolvable problems in T all the time so they will pass the test for sure).*

Proof. Let T_0 be the initial T before any deletion on line 11 of Algorithm 1. Initially, for any B_l such that $\sum_{a \in B_l \cap T_0} z_{l,a} < m_l$, we push it to Q . After that any $a \in B_l$ where $z_{l,a} < 1$ is deleted from T on line 11. Therefore $\forall a \in B_l \cap T, z_{l,a} = 1$ must hold for that B_l . Afterwards, more B_l may violate $\sum_{a \in B_l \cap T} z_{l,a} < m_l$ because some a are deleted from T . But for those B_l , we must have pushed them into Q on line 14. So they will satisfy $\forall a \in B_l \cap T, z_{l,a} = 1$ afterwards. In the end, every $B_l \in Q$ will satisfy $\forall a \in B_l \cap T, z_{l,a} = 1$ and every other $B_l \notin Q$ will satisfy $\sum_{a \in B_l \cap T} z_{l,a} = m_l$. \square

Then we show our strategy is optimal by complementary slackness theorem, which in our context is essentially: 1) the tester only tests problems (non-zero primal variables) that gives him best utility (tight dual constraints); 2) the test taker only memorizes problems (non-zero dual variables) that gives him best utility (tight primal constraints).¹ We will need two more lemmas as follows to use the complementary slackness theorem.

Theorem 4. *Uniformly testing T returned by Algorithm 1 is optimal.*

Proof. First of all, T is non-empty by Lemma 1 so we can test T uniformly. Then since memorizing strategy $z_{l,a}$ makes U_a (the utility to test a problem a) decrease to U if $U_a^0 > U$ and we only pick T from problems that $U_a^0 > U$, every problem that we might test gives the tester best utility. That is the first condition for the complementary slackness theorem. Finally, Lemma 2 shows that one type of test takers either pass the test for sure or memorize a subset in T all the time which is optimal for uniform test, the second condition of complementary slackness theorem holds. \square

Finally, instead of using general LP algorithm like simplex [cite] to solve LP (4), we can alternatively use binary search plus max flow to solve it as shown in Algorithm 2. With a given U , we use max flow to check its feasibility. Intuitively, a U is feasible if and only if there is a strategy $z_{l,a}$ such that for every a where $U_a^0 > U$, that strategy decreases it down to U . This is achievable if and only if the max flow can saturate all edges from A to t . Given a max flow f , we can restore the strategy by setting $z_{l,a} = f_{(B_l,a)} / c_{(B_l,a)}$ which is the flow of an edge divided by the capacity of an edge. If we use Push-Relabel algorithm [cite] to solve the max flow in $O((n+L)^3)$

¹This is essentially saying that in a zero-sum game, a minimax, or maximin state is reached, if and only if the Nash equilibrium is reached.

time, the whole Algorithm 1 runs in $O((n+L)^3 + \sum_l |B_l|)$ time as the remaining part can run in linear time with respect to the input size (assuming we only binary search for constant times to achieve a constant accuracy).

Algorithm 2 Use binary search and max flow to solve LP (4)

```

1:  $U_{lower} \leftarrow 0, U_{upper} \leftarrow \max_{a \in A} U_a^0$ 
2: while  $U_{upper} - U_{lower} > \epsilon$  do ▷ Binary search
3:    $U \leftarrow (U_{upper} + U_{lower}) / 2$  ▷ Objective utility
4:    $V \leftarrow \{s\} \cup \mathbb{B} \cup A \cup \{t\}$  ▷ Vertex set
5:    $E \leftarrow \emptyset$  ▷ Initialize edge set
6:   for all  $a \in A$  do
7:     if  $U_a^0 > U$  then
8:        $E \leftarrow E \cup \{(a, t)\}$ 
9:        $c_{(a,t)} \leftarrow U_a^0 - U$  ▷ Edge capacity
10:    end if
11:  end for
12:  for all  $B_l \in \mathbb{B}$  do
13:     $E \leftarrow E \cup \{B_l\} \times B_l$ 
14:     $c_{\{B_l\} \times B_l} \leftarrow p_l v_l$ 
15:     $E \leftarrow E \cup \{(s, B_l)\}$ 
16:     $c_{(s,B_l)} \leftarrow p_l m_l v_l$ 
17:  end for
18:   $f \leftarrow \text{MAXFLOW}(G = (V, E), c)$ 
19:  if  $f$  saturates all edges in  $A \times \{t\}$  then
20:     $U_{upper} \leftarrow U$ 
21:  else
22:     $U_{lower} \leftarrow U$ 
23:  end if
24: end while
25:  $(\forall l, a), z_{l,a} \leftarrow f_{(B_l,a)} / c_{(B_l,a)}$ 
26: return  $U, z_{l,a}$ 

```

4 Experiment

In this section, we conduct experiments to see how different approaches scale on one-problem test games. The first approach is using CPLEX to solve our general LP (1) which is exponential to the memorize capacity. The second one is using CPLEX to directly solve the marginal version LP (4). The third and fourth use Algorithm 2 to solve LP (4). They are using Edmonds-Karp [cite] and Push-Relabel² [cite] max flow algorithm in C++ boost library respectively.

For each test case, we specify four parameters n, L, m_{max}, b_{max} , the number of problems, the number of test taker types, the maximum memorize capacity and the maximum size of B_l . Given those parameters, a test game instance is generated as follows: for each $1 \leq l \leq L$, generate m_l uniformly from 1 to m_{max} , generate $|B_l|$ uniformly from m_l to b_{max} , generate B_l by selecting $|B_l|$ elements from A uniformly and generate $w_l = v_l p_l$ (weight of type- l) from $[0, 1]$ uniformly. Here we only generate w_l because v_l and p_l always appear together as $v_l p_l$. For each parameter setting, we generate 5 test game instances and compute the average running time for each algorithm. During each run, we set timeout to 5 seconds.

²We modified its `is_flow` function for floating number issues.

Figure 1: Average running time (timeout after 5 seconds) for different one-problem test algorithms over the number of problems n . For each n , we set the other three parameters $L = n$, $m_{max} = \lfloor \sqrt{n} \rfloor$ and $b_{max} = 2m_{max}$ (see Table 1).

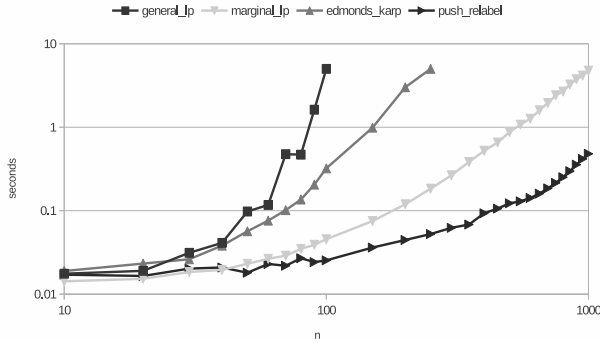


Table 1: Parameter settings for the benchmark

n	L	m_{max}	b_{max}
10	10	3	6
50	50	7	14
100	100	10	20
300	300	17	34
500	500	22	44
700	700	26	52
1000	1000	31	62

In our benchmark, we increase n and set the remaining three parameters L, m_{max}, b_{max} as following: $L = n$, $m_{max} = \lfloor \sqrt{n} \rfloor$ and $b_{max} = 2m_{max}$. Note that our L, b_{max} are quite small (they can be as large as $b_{max} = n/2$, $L = \binom{n}{b_{max}}$). That's because for large L, b_{max} , our uniform random input will almost always make the trivial strategy to uniformly test among all n problems optimal. So we make them small to avoid such trivial solutions. The realized parameter settings are shown in Table 1. The benchmark is shown in Figure 1. As expected, the general LP is exponential so it times out before n reaches 100. Unexpectedly, using CPLEX to directly solve LP (4) is much faster than Edmonds-Karp implementation. Finally, Push-Relabel version outperforms all others significantly. It solves $n = 1000$ case in about 0.5 second while the second best CPLEX version uses almost 5 seconds.

[Machine, GLPK and GPLEX spec]

5 Conclusion

Acknowledgments

The preparation of these instructions and the \LaTeX and Bib \TeX files that implement them was supported by Schlumberger Palo Alto Research, AT&T Bell Laboratories, and Morgan Kaufmann Publishers. Preparation of the Microsoft Word file was supported by IJCAI. An early version of this document was created by Shirley Jowell and Peter F. Patel-

Schneider. It was subsequently modified by Jennifer Balentine and Thomas Dean, Bernhard Nebel, and Daniel Pagenstecher. These instructions are the same as the ones for IJCAI-05, prepared by Kurt Steinkraus, Massachusetts Institute of Technology, Computer Science and Artificial Intelligence Lab.

A \LaTeX and Word Style Files

The \LaTeX and Word style files are available on the IJCAI-13 website, <http://www.ijcai-13.org/>. These style files implement the formatting instructions in this document.

The \LaTeX files are `ijcai13.sty` and `ijcai13.tex`, and the Bib \TeX files are named `.bst` and `ijcai13.bib`. The \LaTeX style file is for version 2e of \LaTeX , and the Bib \TeX style file is for version 0.99c of Bib \TeX (*not* version 0.98i). The `ijcai13.sty` file is the same as the `ijcai07.sty` file used for IJCAI-07.

The Microsoft Word style file consists of a single file, `ijcai13.doc`. This template is the same as the one used for IJCAI-07.

These Microsoft Word and \LaTeX files contain the source of the present document and may serve as a formatting sample.

Further information on using these styles for the preparation of papers for IJCAI-13 can be obtained by contacting pcchair13@ijcai.org.

References