Incentivizing Truthfulness Through Audits in Strategic Classification

Andrew Estornell

Computer Science & Engineering Washington University in St. Louis aestornell@wustl.edu

Sanmay Das

Computer Science George Mason University sanmay@gmu.edu

Yevgeniy Vorobeychik

Computer Science & Engineering Washington University in St. Louis yvorobeychik@wustl.edu

Abstract

In many societal resource allocation domains, machine learning methods are increasingly used to either score or rank agents in order to decide which ones should receive either resources (e.g., homeless services) or scrutiny (e.g., child welfare investigations) from social services agencies. An agency's scoring function typically operates on a feature vector that contains a combination of self-reported features and information available to the agency about individuals or households. This can create incentives for agents to misrepresent their self-reported features in order to receive resources or avoid scrutiny, but agencies may be able to selectively audit agents to verify the veracity of their reports.

We study the problem of optimal auditing of agents in such settings. When decisions are made using a threshold on an agent's score, the optimal audit policy has a surprisingly simple structure, uniformly auditing all agents who could benefit from lying. While this policy can, in general be hard to compute because of the difficulty of identifying the set of agents who could benefit from lying given a complete set of reported types, we also present necessary and sufficient conditions under which it is tractable. We show that the scarce resource setting is more difficult, and exhibit an approximately optimal audit policy in this case. In addition, we show that in either setting verifying whether it is possible to incentivize exact truthfulness is hard even to approximate. However, we also exhibit sufficient conditions for solving this problem optimally, and for obtaining good approximations.

1 Introduction

Algorithmic decision-making systems are increasingly used to make high-stakes resource allocation decisions by social services agencies. This includes both scarce resource settings, where the demand for a limited pool of resources exceeds supply (for example, housing for the homeless (Kube, Das, and Fowler, 2019)), as well as risk-scoring settings, where only those who fall above or below a certain threshold are either given a resource (for example, a loan (Agarwal, Skiba, and Tobacman, 2009)) or targeted for further scrutiny (for example, parents suspected of child maltreatment or neglect (Chouldechova et al., 2018)). As is standard in classification and ranking settings, each individual or household (henceforth *agent*) is associated with a feature

Copyright © 2021, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

vector. In many such settings, the feature vector will combine information submitted by the agents themselves with information about them available from other sources. For example, in prioritizing households for homeless services, agencies make decisions based on self-reported items (e.g., history of alcohol or drug use) as well as on information available to them in government records (e.g., child-support or welfare payments received) (Brown et al., 2018). Naturally, this creates incentives for agents to try and game the system by strategically choosing their self-reported features in order to maximize their chances of receiving the resource or avoiding scrutiny.

Prior work on strategic or adversarial classification has considered a closely related problem where agents subject to classification can modify feature values at some cost or subject to a constraint on the total magnitude of such modification, with the goal of inducing an incorrect prediction (Athalye, Carlini, and Wagner, 2018; Carlini and Wagner, 2017; Hardt et al., 2016; Milli et al., 2019; Papernot et al., 2018; Tong et al., 2019; Vorobeychik and Kantarcioglu, 2018). This research has typically focused on either assessing how vulnerable particular families of classifiers are to such attacks (often termed adversarial examples) (Athalye, Carlini, and Wagner, 2018; Carlini and Wagner, 2017; Lowd and Meek, 2005; Xu, Qi, and Evans, 2016), or on designing classifiers that are robust in the sense that the prediction remains unchanged even after budget-constrained feature modifications (Brückner and Scheffer, 2011, 2012; Hardt et al., 2016; Li and Vorobeychik, 2018; Madry et al., 2018; Tong et al., 2019; Wong and Kolter, 2018). In this literature, the interests of the agents are commonly viewed as opposed to those of the decision-maker (e.g., learner), often motivated by security considerations (Šrndic and Laskov, 2014; Xu, Qi, and Evans, 2016). Moreover, the typical models representing costs to agents of modifying features are at times not adequate at capturing realistic limits on what agents can do (Tong et al., 2019; Wu, Tong, and Vorobeychik, 2020). In contrast, in the kinds of social services settings we describe, and potentially numerous others (e.g., tax filing), the costs of misrepresenting one's self-reported features are better captured by the risk associated with being audited than, say, a hard constraint on how much the features are modified. Moreover, the agents' interests are not fundamentally opposed to the principal's; rather, this is a case of misaligned incentives more akin to that studied in the incentive design literature (Haeringer, 2018; Nisan et al., 2007).

We consider a principal who has a limited budget of audits and can use these to determine whether an agent is telling the truth, with the cost of failing an audit the primary deleterious consequence to dishonest agents. For example, caseworkers can interview associates of the agent and ask about behavioral issues, alcohol or drug use, and the like, and impose restrictions or fines on the agent if the results reveal dishonesty. We suppose that the principal uses a score function f(for example, learned risk scores) that takes agent features as input in order to decide whether an agent is subject to further scrutiny whenever their score exceeds a predefined threshold (we term this the *threshold* setting), or to allocate resources to the agents with the top k values of f (we call this the top-k setting). We specifically focus on two problems: 1) designing an audit policy for the principal that minimizes incentives to lie, defined in terms of approximate Bayes-Nash incentive compatibility (ε -BNIC), and 2) verifying whether it is possible to ensure truthful reporting of features.

We show that in the threshold allocation setting an optimal policy audits uniformly at random all agents who are above the threshold, with special consideration for those who are either obviously lying or telling the truth. Although this policy is in general hard to compute, we present sufficient conditions under which it is tractable. In the top-k setting, we prove that auditing all agents who receive the scarce resource uniformly at random (again, modulo special treatment of agents who are either certainly truthful or dishonest) yields an additive approximation bound, although the problem is hard in general. Furthermore, we show that this audit policy is optimal if we consider dominant strategy incentive compatibility as a solution concept instead of ε -BNIC.

Surprisingly, the verification problem is even harder: determining if any audit policy can incentivize truthful reporting is #P-hard even for a uniform prior over features and only two agents. However, we give sufficient conditions under which verification becomes tractable in the threshold setting for both piecewise linear and logistic scoring functions. Our corresponding results are weaker for the top-k setting, where we require the distribution over features to be uniform to obtain a tractable algorithm for checking incentives to lie *assuming* that a uniform audit policy is used. Finally, we show that for distributions for which we can efficiently approximate integrals over intervals, we can also approximately verify incentive compatibility.

Our results are important for understanding the potential for audits to be useful in various social services settings. Of perhaps the most practical importance is the clear distinction we find between the threshold (modeling unlimited, but costly, deployment of resources) and top-k (modeling scarce resource allocation) settings in terms of the difficulty of finding a good audit policy, and the simplicity of the optimal audit policy in the threshold setting.

2 Preliminaries

We consider a setting with a collection of n agents in which either a scarce resource is distributed among k of them using a score function, or each agent is scored to determine

whether they are selected to receive a resource. Each agent is associated with a vector of attributes (features) which are grouped into two categories: "known", denoted by x, and "self-reported", denoted by z. Throughout, we refer to (x, z)as an agent's true type, to contrast it with (x, z') in which z' is self-reported and may be different from the true corresponding characteristics of the agent. For example, the agent may have a history of substance abuse, corresponding to "true" $z_j = 1$, but reports that they do not, with "reported" $z'_i = 0$. Let d be the number of known and s the number of self-reported features. We assume that each feature in either category either belongs to a continuous or discrete interval, i.e., each $x_j, z_k \in I = [a, b] \cap S$, where $S = \mathbb{R}$ (continuous interval) or $S = \mathbb{Z}$ (discrete interval). We further assume that the true types of each of the n agents are i.i.d. according to a (common knowledge) prior distribution D with PDF (or PMF, in the discrete case) denoted by $h: I^d \times I^s \to [0,1]$. We will use $\mathbb{P}(\cdot)$ to denote the associated probability measure.

Let $\mathcal{A} = \{\mathbf{a}_1,...,\mathbf{a}_n\}$ denote the collection of n agents, where $\mathbf{a}_i = (\mathbf{x}_i,\mathbf{z}_i) \in I^d \times I^s$ represents the agent's *true* type, and let $\mathcal{A}' = \{\mathbf{a}'_1,...,\mathbf{a}'_n\}$ be the collection of reported types, $\mathbf{a}'_i = (\mathbf{x}_i,\mathbf{z}'_i)$. We assume that each agent knows their own type, but only knows the common prior h about the types of other agents.

The principal publishes a score function $f:I^d\times I^s\to\mathbb{R}$ that takes each agent's *reported* type \mathbf{a}_i' as input, and returns a real-valued score. For example, f may represent the probability (learned from historical data) that a homeless person will be safely and stably housed in 1 year if allocated a housing resource. There are two common ways that f is used in resource allocation: (1) **Threshold allocation:** all agents scoring above a threshold θ are allocated a resource (e.g., not chosen for further scrutiny in a child neglect case), and (2) **Top-**k **allocation:** agents with the highest k scores based on reported types are allocated a resource (e.g., housing).

The principal can *audit* up to B agents and thereby verify whether their reported type matches their true type. Let ϕ denote the audit policy, which is a function of the full collection of n reported types \mathcal{A}' . We consider stochastic audit policies, where $\phi_i(\mathcal{A}') \in [0,1]$ is the probability that agent i is audited. If an audit of agent i determines that the agent has lied, i.e., $\mathbf{z}'_i \neq \mathbf{z}_i$, there are two consequences: 1) the agent does not receive the resource, and 2) the agent pays a penalty (fine) $c \ge 0$. Let α denote the allocation policy with $\alpha_i(f, \mathcal{A}', \phi) = 1$ if agent i receives the resource, and 0 otherwise. Further, let $\mathcal{L}_i = 1$ if agent i is audited and $\mathbf{z}_i' \neq \mathbf{z}_i$ (the agent is caught lying) and 0 otherwise; note that since the audit policy is stochastic, \mathcal{L}_i is a random variable. We assume that an agent obtains a value of 1 for receiving the resource and 0 otherwise. Consequently, the agent's utility is $u_i(\mathcal{A}') = \alpha_i(f, \mathcal{A}', \phi)(1 - \mathcal{L}_i) - c\mathcal{L}_i$.

This game between a principal and agents can be expressed as the following sequence of events:

- 1. The principal knows D, n, $I^d \times I^s$, α , c, and f, and announces an audit policy ϕ .
- 2. Realizations of n agents are drawn i.i.d. from D. Each agent knows its own type $(\mathbf{x}_i, \mathbf{z}_i)$, D, n, $I^d \times I^s$, α , c, ϕ ,

and f, but does not know the types of other agents.

- 3. Agents simultaneously submit their reported type $(\mathbf{x}_i, \mathbf{z}'_i)$, where \mathbf{z}'_i need not equal \mathbf{z}_i .
- 4. The principal audits up to B agents, according to ϕ . Any agent i found to have reported $\mathbf{z}_i' \neq \mathbf{z}_i$ is removed from consideration (not allocated the resource), and pays a fine of c.
- The remaining agents are distributed a resource according to α.

Note that if an agent i is found to be dishonest through an audit in the top-k allocation setting, another agent would receive the resource in place of i.

The goal of the principal is to achieve truthful reporting of types by the agents in an (approximate) Bayes-Nash equilibrium, or (approximate) *Bayes-Nash incentive compatibility* (BNIC). Formally:

Definition 1. $(\varepsilon$ -BNIC) An audit policy ϕ is ε -Bayes-Nash incentive compatible $(\varepsilon$ -BNIC) if for all i and \mathbf{a}_i ,

$$\mathbb{E}_{\mathcal{A}_{-i} \sim D}[u_i(\mathbf{a}_i, \mathcal{A}_{-i})|f, \phi, \alpha]$$

$$\geq \mathbb{E}_{\mathcal{A}_{-i} \sim D}[u_i(\mathbf{a}_i', \mathcal{A}_{-i})|f, \phi, \alpha] - \varepsilon \quad \forall \mathbf{a}_i' : \mathbf{x}_j' = \mathbf{x}_j.$$

$$\phi \text{ is BNIC if it is 0-BNIC.}$$

We consider two problems in this setting. First, since it is in general impossible to induce BNIC, as we show below, we aim to identify an *optimal* audit policy, defined as follows.

Definition 2. (Optimal) An audit policy ϕ is optimal if ϕ induces an ε^* -BNIC, and there does not exist another policy ϕ' for which truthful reporting is an ε -BNIC with $\varepsilon < \varepsilon^*$.

In other words, the optimal ϕ induces the least incentive to lie among all policies.¹ As a consequence, if we find that an optimal policy is not BNIC, then no policy can be. Our second problem is to determine the smallest ϵ that can be induced by an audit policy. We show that in general, these problems have differing complexity.

Before proceeding with a general analysis, we make three observations about our model: 1) if B=n, any score function f can be made BNIC; 2) if $k \in \{0,n\}$, the top-k case is trivially BNIC; and 3) if $(1+c)(B/k) \ge 1$, the top-k case is again trivially BNIC.

We begin by showing that without auditing the self-reported features (equivalently, when the audit budget B=0), ensuring BNIC amounts to ignoring ${\bf z}$ altogether whenever we use a deterministic scoring function f. Since self-reported features may be important in determining priority of individuals for resources, this impossibility motivates a careful treatment of optimal auditing, which follows.

Proposition 1. Suppose B=0. Then, both the top-k and threshold mechanism are incentive compatible iff $f(\mathbf{x}, \mathbf{z}) = f(\mathbf{x})$. Moreover, in the threshold setting, BNIC can be achieved only if c > 0.

Due to space constraints, this and other full proofs are deferred to the supplement.

3 Design of Optimal Audit Policies

The problem of incentivizing truthfulness via auditing can be broken into two primary components: design and verification. The first component, design, is the construction of *optimal* or approximately *optimal* audit policies. The second, verification, focuses on computing the maximum incentive to lie under an *optimal* audit policy, denoted as ε^* . Although both problems are in general hard, we show that verification is intrinsically "harder" in the sense that in a wide range of settings optimally auditing agents is tractable, but computing ε^* remains hard. The focus of this section is on design. In particular, we exhibit a simple audit policy which is guaranteed to be optimal under the threshold allocation setting, and approximately optimal under the top-k allocation setting.

We begin with some remarks and notation that will be subsequently used in characterizing the optimal audit policies. When selecting which agents to audit, the principal is unaware of each agent's true type $\mathbf{a}_i = (\mathbf{x}_i, \mathbf{z}_i)$, and sees only the reported type $\mathbf{a}_i' = (\mathbf{x}_i, \mathbf{z}_i')$. Since the principal is interested in minimizing the marginal gain that *any* agent can achieve from lying, agents' true types must be considered through the lens of worst-case analysis. Note that the type with the largest incentive to report $(\mathbf{x}_i, \mathbf{z}_i')$ is the type with the lowest scoring \mathbf{z} , given *known* type \mathbf{x}_i (denoted as $\mathbf{a}_i^* = (\mathbf{x}, \mathbf{z}_i^*)$. From the principal's perspective, any agent reporting $(\mathbf{x}_i, \mathbf{z}_i')$ must be assumed to have true type $(\mathbf{x}_i, \mathbf{z}_i^*)$.

With this in mind, agent reports can be classified as one of the following: a sure-truth, a sure-lie, or suspicious. Sure-truths are reports which are guaranteed to be honest (e.g. $\mathbf{z}_i' = \mathbf{z}_i^*$). Sure-lies are reports which are guaranteed to be false (these are only of the form $h(\mathbf{x}_i, \mathbf{z}_i') = 0$). Suspicious reports are those with an unknown truth value. The following two definitions formalize these observations.

Definition 3. (Minimum Type) For any known partial type \mathbf{x}_i , we say the minimum type of \mathbf{x}_i is $\mathbf{a}_i^* = (\mathbf{x}_i, \mathbf{z}_i^*) = \arg\min_{\mathbf{z} \in I^s: h(\mathbf{x}_i, \mathbf{z}) > 0} f(\mathbf{x}_i, \mathbf{z}).$

Definition 4. (Suspicious) We say a type \mathbf{a}_i' is suspicious if the minimum type \mathbf{a}_i^* has a strictly lower chance of being allocated a resource barring auditing, i.e., $\mathbb{E}_{\mathcal{A}_{-i}}\left[\alpha_i(f, \mathcal{A}_{-i} \cup \{\mathbf{a}_i'\})\right] > \mathbb{E}_{\mathcal{A}_{-i}}\left[\alpha_i(f, \mathcal{A}_{-i} \cup \{\mathbf{a}_i^*\})\right]$

The key point here is that the principle should never waste an audit on a sure-truth, and when looking at incentive compatibility (i.e. single deviations from collective truth-telling), there is at most one sure-lie in any set of reports, which should be audited with probability 1. The more interesting question regarding audit polices is; what to do with *suspicious* reports.

3.1 Threshold Allocation

Recall that in the threshold allocation setting, an agent receives a resource if $f(\mathbf{x}, \mathbf{z}') \geq \theta$, where $(\mathbf{x}, \mathbf{z}')$ is the agent's reported type. We first show that, in general, optimal auditing under threshold allocation is NP-hard in general, but is tractable if and only if identifying sure-truths is tractable. The hardness of auditing stems from the possibly arbitrary relationship between the distribution D and the score function f.

¹To avoid confusion, note that the principal could have other objectives, and our definition of optimality is specific to inducing the "best" approximation of BNIC.

Theorem 1. For a given set of n reports A' and a budget B, computing an optimal audit policy is NP-hard.

Proof Sketch. This result stems from the observation that the principal would never want to "waste" an audit on an agent whose report is guaranteed to be truthful. For example, suppose agent $\mathbf{a}_1 = (\mathbf{x}_1, \mathbf{z}_1)$ reports type $\mathbf{a}_1' = (\mathbf{x}_1, \mathbf{z}_1')$ with $f(\mathbf{x}_1, \mathbf{z}_1') \geq \theta$. Suppose further that for all \mathbf{z} with $f(\mathbf{x}_1, \mathbf{z}) < \theta$, we have $h(\mathbf{x}_1, \mathbf{z}) = 0$. Then the principal is certain that agent 1 is truthful since this agent's true type could not have scored below the threshold. Due to this dependency on the underlying distribution, one can encode a SAT formula into the distribution such that determining if there exists a \mathbf{z} such that $h(\mathbf{x}_1, \mathbf{z}) > 0$ and $f(\mathbf{x}_1, \mathbf{z}) < \theta$ is equivalent to determining the satisfiability of the SAT instance.

To better understand the nature of the problem of characterizing an optimal audit policy, consider the following simple example.

Example 1. Suppose there are two agents with one *known* and one *self-reported* binary feature, and suppose that z=1 if x=1, and can be either 0 or 1 according to some prior distribution if x=0. Further, suppose that f(x,z)=z and $\theta=1/2$, which means that an agent receives the resource iff z=1. Now, suppose that B=1 and the principal observes two types: (1,1) and (0,1). Clearly, the principal would not audit the former, since x=1 already implies that the agent is honest, but would audit the latter. This simple example suggests that one could expect an optimal audit policy to depend in rather complex ways on the observed types \mathcal{A}' .

However, we show that a simple policy of uniformly auditing all *suspicious* agents (Definition 4), is optimal. We call this policy UNIFORM, and define it formally next.

Definition 5. (UNIFORM) For a given set of reports \mathcal{A}' , let $G(\mathcal{A}')$ be the set of all agent's whose reports are suspicious. Given budget B, the UNIFORM audit policy audits each $\mathbf{a}' \in \mathcal{A}'$ with probability

$$\phi_{i}(\mathcal{A}') = \begin{cases} 1 & \text{if } h(\mathbf{a}'_{i}) = 0\\ \min\left(\frac{B}{|G(\mathcal{A}')|}, 1\right) & \text{if } \mathbf{a}'_{i} \in G(\mathcal{A}'),\\ 0 & \text{otherwise} \end{cases}$$

Next, we show that in the threshold allocation setting, this UNIFORM audit policy is optimal.

The intuition for the optimality of UNIFORM comes from the fact that any type ${\bf z}$ can report any other type ${\bf z}'$ at no cost. This means that any lie that gets an agent above the threshold is equivalent, modulo auditing. Thus, if an audit is non-uniform (as long as the reported type is above the threshold), some lies become more valuable than others, and we should shift auditing to those lies (more precisely, to agents who feature such lies). The discontinuity arises by observing a sure-lie (i.e., h(x,z)=0); only in this case do we know which agent was dishonest, and can thus place higher audit weight on this agent without increasing the value of lying for any other agent.

Note that this implies optimal auditing is equivalent to identifying sure-truths.

Theorem 2. In the threshold allocation setting, for any score function f, UNIFORM is an optimal audit policy.

Proof Sketch. For the sake of illustration, we demonstrate how this result holds in the cases of a discrete distribution over agent types. An identical idea holds for continuous features, although the technical details differ.

When analyzing ε -BNIC, we are considering the value that any agent gains when deviating from a truthful reporting, while all other agents remain truthful, i.e. we consider this case when at most one report is dishonest. In any set of reports \mathcal{A}' , if the principal sees a sure-lie, they are immediately aware of the dishonest agent's identity and should exclusively audit that agent, since all other agents are guaranteed to be truthful.

The principal's objective is to minimize the expected gain of any type \mathbf{a}_i misreporting their type as \mathbf{a}_i' , when all other agents are truthful. Note that when all agents, aside from agent i are honest, the set of reported types $\mathcal{A}' = \mathcal{A}_{-i} \cup \{\mathbf{a}_i'\}$ (where \mathcal{A}_{-i} is the set true types for all other agents). As such, we can express the minimum expected gain of misreporting, achievable by any audit policy ϕ , as

$$\varepsilon = \min_{\phi} \max_{\mathbf{a}_i', \mathbf{a}_i} \left(\mathbb{E}_{\mathcal{A}_{-i}} \left[\alpha_i(f, \mathcal{A}') - \alpha_i(f, \mathcal{A}) \right] \right)$$
 (1)

$$- \underset{\mathcal{A}_{-i}}{\mathbb{E}} \left[\left(\alpha_i(f, \mathcal{A}') + c \right) \phi_i(\mathcal{A}') \right] \right)$$
 (2)

Where term (1) the expected difference in the allocation decision between agent i falsely reporting \mathbf{a}_i' or truthfully reporting \mathbf{a}_i , and term (2) represents the expected cost of being caught lying when reporting \mathbf{a}_i' . Making use of two simple observations, we can simplify this equation. First, in the threshold setting, agents know both their own type and the threshold θ , thus agent i knows the allocation decision on both the true type \mathbf{a}_i , and reported type \mathbf{a}_i' , meaning that the expectations on α can be dropped. Second, we need only consider this term for *suspicious* agents, so we may assume that $\alpha_i(f, \mathcal{A}') = 1$ and $\alpha_i(f, \mathcal{A}) = 0$. With this, the equation can be simplified to

$$\varepsilon = \min_{\boldsymbol{\phi}} \max_{\substack{\mathbf{a}_i', \mathbf{a}_i \\ f(\mathbf{a}_i') \geq \boldsymbol{\theta} > f(\mathbf{a}_i)}} 1 - (1+c) \mathbb{E}_{\mathcal{A}_{-i}} \big[\phi_i(\mathcal{A}') \big]$$

Thus ε is solely determined by the value of $\mathbb{E}_{\mathcal{A}_{-i}}[\phi_i(\mathcal{A}')]$ for any *suspicious* type \mathbf{a}'_i . Let

$$G(\mathcal{A}') = \{ (\mathbf{x}, \mathbf{z}') \in \mathcal{A}' : f(\mathbf{x}, \mathbf{z}') \ge \theta \text{ and } \exists \mathbf{z}^* \text{s.t. } f(\mathbf{x}, \mathbf{z}^*) < \theta \text{ and } h\big((\mathbf{x}, \mathbf{z}^*)\big) > 0 \text{ and } h\big((\mathbf{x}, \mathbf{z}')\big) > 0 \}$$

be the set of *suspicious* types in \mathcal{A}' . In the case when agent features are distributed according to a discrete distribution, this expectation can be expressed as

$$\mathbb{E}_{\mathcal{A}_{-i}} \left[\phi_i(\mathcal{A}') \right] = \sum_{\mathcal{A}_{-i}} \phi_i(\mathcal{A}') Q(\mathcal{A}_{-i})$$
$$= \sum_{\mathcal{A}_{-i}} \min \left(1, \frac{B}{G(\mathcal{A}')} \right) Q(\mathcal{A}_{-i}),$$

where $Q(\mathcal{A}_{-i})$ is the probability of any realization of the specified types of agents other than i induced by D. The probability which sure-lies are audited has no effect on the value of other lies, and thus sure-liescan be audited with probability 1. Moreover, the sum is equal for any two suspicious agents with $h(\mathbf{a}') > 0$. In each set of reports \mathcal{A}' , the principle fully spends their budget (or audits all suspicious types with probability 1) and $Q(\mathcal{A}_{-i})$ is independent of the type agent i reports. Thus, for any policy different from UNIFORM, at least one audit weight must be changed, i.e., $\phi_i(\mathcal{A}') \neq \min\left(1, \frac{B}{G(A')}\right)$, for some i and some \mathcal{A}' . As a result of the tightness and independence of $Q(\mathcal{A}_{-i})$, this change of audit weight could only result in a (not necessarily strict) increase in the expected gain of misreporting for any agent type.

Note that while optimal, UNIFORM is in general intractable because of the combinatorial structure of such policies that may be induced by $h(\cdot)$. However, we now show that for sufficiently well-behaved h and f we can compute UNIFORM efficiently.

Theorem 3. The audit policy UNIFORM can be computed in polynomial time if for any report $(\mathbf{x}, \mathbf{z}')$, it can be efficiently determined if $(\mathbf{x}, \mathbf{z}')$ is a sure-truth, (i.e. there exists a self reported type \mathbf{z}^* , such that $h(\mathbf{x}, \mathbf{z}^*) > 0$ and $f(\mathbf{x}, \mathbf{z}^*) < \theta$).

3.2 Top-k Allocation

We now turn our attention to selecting the *optimal* audit policy when resources are given to the k highest scoring agents. In this case, the optimal policy no longer admits a clean characterization. The main challenge is that now there are far more complex interdependencies among agents' benefits from lying, other agents' reports, and the audit policy. For example, if an agent in the top-k is caught lying, another agent would now receive the resource. Instead, we study a natural adaptation of UNIFORM to this setting, and exhibit an additive approximation bound for its optimality. We then show that if we use dominant strategy incentive compatibility (defined formally below) as a solution concept in place of BNIC, uniform auditing is optimal even in this setting.

We begin by showing that optimal auditing in the top-k setting is NP-hard even when sure-truths are identifiable in constant time.

Theorem 4. In the top-k allocation setting determining which agents should be audited is NP-hard, even for n=4 agents, monotone f, uniform D, and even if sure-truth can be identified in constant time.

Proof Sketch. We can encode an instance of Vertex Cover into f such that agents with a self-reported type, which constitutes a vertex cover, ranks in the top-k with extremely low probability, while all other types have score proportional to number of vertices that their self-reported type "covers". For a sufficiently small budget and penalty for lying, there will be agents whose expected value of lying (even if never audited) is smaller than agents who receive the highest probability weight. As such, the principal must determine which agents should receive zero audit weight, which is NP hard due to the encoding of VC.

Now, consider a variant of the UNIFORM policy in the topk setting where we uniformly at random audit agents who have scores in the top k. We first define this policy formally.

Definition 6. (UNIFORM-K) For any set of reported types \mathcal{A}' , let UNIFORM-K denote the policy of auditing each of the top-k agents (refereed to as the set $T_k \subset \mathcal{A}'$) with probability $\min(1, B/k)$.

Next, we show that UNIFORM-K admits an additive approximation of an optimal audit policy in the top-k setting. Recall that multiplicative approximations are in general NP-hard to achieve.

Theorem 5. Let ϕ denote the audit policy UNIFORM-K. Then the maximum utility gained by lying under ϕ is no more than $\max\left(0,1-\frac{(1+c)B}{k}\right)$ greater than that of the optimal audit policy, and this bound is tight.

Proof. This is the result of simple worst case analysis on the expected value of lying, which can be expressed as

$$\mathbb{E}_{\mathcal{A}_{-i}} \left[\alpha_i(\mathcal{A}') \left(1 - \phi_i(\mathcal{A}') \right) - c \phi_i(\mathcal{A}') - \alpha_i(\mathcal{A}) \right]$$

=\mathbb{P}(\mathbf{a}'_i \in T_k) \mathbb{E}[\phi_i(\mathcal{A}')|\mathbf{a}'_i \in T_k] - c \mathbb{E}[\phi_i(\mathcal{A}')] + \mathbb{P}(\mathbf{a}_i \in T_k)

In the worst case, the expected value of lying could be 0 for all agents. However, the uniform audit policy will have expected value of lying equal to $\mathbb{P}(\mathbf{a}_i' \in T_k)(1-(1+c)\frac{B}{k}) - \mathbb{P}(\mathbf{a}_i \in T_k)$. Which again in the worst case is equal to $((1-(1+c)\frac{B}{k})$

This bound is tight to within any small $\beta>0$. To see this, construct an instance with 3 agents of types $x\in\{0,1\},z\in\{0,1\}$. Where $f(x,z)=x\wedge z$. Let $\mathbb{P}(x=1,z=1)=\beta$ and the rest have probability $\frac{1-\beta}{3}$. Let B=1 and k=2. Then an optimal audit policy yields $\varepsilon^*=0$, but uniformly auditing the top-k yields $\varepsilon=((1-\beta^2)(1-(1+c)\frac{B}{k}).$

A major part of what makes auditing difficult is the dependence on the distribution. We now consider an alternative solution concept which eliminates this dependence: ε -Dominant Strategy Incentive Compatibility (ε -DSIC). Specifically, under ε -DSIC the principal aims to design a policy under which truthful reporting is (approximately) optimal for agents regardless of other agents' types.

Definition 7. An audit policy ϕ is ε -DSIC if for all i and \mathbf{a}_i ,

$$\mathbb{E}[u_i(\mathbf{a}_i, \mathcal{A}_{-i})|f, \phi, \alpha, \mathcal{A}'_{-i}]$$

$$\geq \mathbb{E}[u_i(\mathbf{a}'_i, \mathcal{A}_{-i})|f, \phi, \alpha, \mathcal{A}'_{-i}] - \varepsilon \ \forall \mathbf{a}'_i : \mathbf{x}'_i = \mathbf{x}_i \ and \ \forall \mathcal{A}'_{-i}.$$

Theorem 6. In the top-k setting, UNIFORM-K yields ε^* -DSIC with an optimal ε^* .

Proof Sketch. In the top-k setting the key difference from ε -BNIC is that for any realization \mathcal{A}_{-i} and any set of corresponding reports \mathcal{A}'_{-i} , agent i knows the allocation decision on both their true type \mathbf{a}_i and any reported type \mathbf{a}'_i . This certainty of outcomes it precisely what made all *suspicious* reports equivalent in the threshold case. Using a similar argument for the optimality of UNIFORM in the threshold case, we can see that UNIFORM-K is optimal in the top-k case.

4 Verification of Policy Effectiveness

In the previous section we showed that in many circumstances we can fully characterize the optimal audit policy, and it can be efficiently computed for a broad range of settings. We now consider the problem of *verification*, that is, computing the smallest ϵ^* that we can achieve for an optimal audit policy. We show that this problem is hard even when auditing is easy. Subsequently, we first show that we can often effectively approximate this problem, and then exhibit special cases in which we can even compute this ϵ^* efficiently.

4.1 Complexity of Verification

In the threshold setting, we will show that computing the minimum ε^* inducible by any policy is #P-hard, even in cases when optimal auditing is tractable. This complexity stems from *both* the score function f and distribution D. Intuitively, these uniquely define both the set of agent types which are considered *suspicious* and the probability that a *suspicious* type will occur. As *suspicious* types are more likely to occur, the probability that any particular agent is audited decreases. Thus, we can encode "hard" problems into f or D where agent types (binary vectors) correspond to satisfying assignments of the encoded problem. We can also observe that if the number of possible agent types is polynomial, then the problem is trivially tractable through brute force search.

We show here hardness in terms of f; a similar construction works to show the hardness in terms of D. In this construction, optimal auditing is easy even in the top-k case.

Theorem 7. In both the threshold and top-k setting, computing the minimum ε inducible by any audit policy is #P-hard, for both continuous and discrete features, even when the feature distribution is uniform, there are only 2 agents, and f is both monotone and binary.

Proof Sketch. For this proof sketch we will work in the setting of threshold allocation and discrete features, similar logic holds in the other cases. We reduce from #VC. For a graph G=(V,E), let D be uniform and agents be $\mathbf{a}=\langle x_1,...,x_{|V|},z_1\rangle$, for $x,z\in\{0,1\}$. Let $\theta=\frac{1}{2}$ and set

$$f(\mathbf{x}, z) = \left(\bigwedge_{(v_r, v_t) \in E} (x_r \vee x_t) \right) \wedge z_1.$$

Under this construction of f we see that an agent scores $f(\mathbf{a}) = 1$ if and only if \mathbf{x} constitutes a vertex cover and z = 1. Thus when B = 1 and n = 2, if agent 1 scores below $\frac{1}{2}$ and is considering misreporting their type, they are audited with lower probability if $f(\mathbf{a}_2) = 1$. Since D is uniform, the probability of this occurring is equivalent to the number of vertex covers of G.

In addition to hardness of checking BNIC, we can show that it is even hard to multiplicatively approximate an ε -BNIC in the threshold and top-k settings.

Theorem 8. Multiplicatively approximating to any constant factor the smallest ε such that there is an ε -BNIC audit policy, in both threshold and top-k allocation is NP-hard even for $\Theta(1)$ agents.

Proof. This result is a straightforward consequence of the construction in the proof of Theorem 7. In that proof we encode an NP-hard problem into an instance of our problem, and show that determining if truthful reporting is BNIC is equivalent to counting the number of satisfying assignments of vertex covers. If we reduce instead from an Unambiguous-SAT instance (f is no longer monotone), then the mechanism is BNIC if and only if the formula has exactly one satisfying assignment. This would imply that $\varepsilon=0$ if and only if the U-SAT instance has no satisfying assignment, and any multiplicative factor ε would likewise be zero, immediately indicating the satisfiability of the U-SAT instance.

Note that UNIFORM-K is the optimal audit policy in these cases, implying that not only is verification of an optimal policy hard, but also verification of UNIFORM-K is also in general hard.

In summary, the problem of *checking* whether a particular setting is ε -BNIC is hard, even in instances when auditing is tractable. To further outline the relation of the complexity of both problems we make the following observation.

Theorem 9. In the threshold allocation setting, verification being in P implies optimal auditing is also in P.

Next, we turn to positive results. To begin, we now show that when agents' *minimum* type can be efficiently computed we can achieve a probabilistic bound on the value of lying in polynomial time, via Monte Carlo simulations.

Theorem 10. Suppose that ε^* is the minimum value for which UNIFORM is ε^* -BNIC. Then, for any $\gamma \in \Theta(1)$, performing n^{γ} rounds of Monte-Carlo sampling will yield a value ε' , such that $\varepsilon' = \varepsilon^* \pm \Theta(1/\sqrt{n^{\gamma}-3})$ with probability at least $1 - 1/n^2$. This can be done in time $\Theta(n^{\gamma+1})$.

Observe from Theorem 10 that as n increases, the error of approximation tends towards 0 with probability tending towards 1. Next, we consider special cases in which verification is tractable.

4.2 Tractable Special Cases

Thus far, our results are negative when it comes to checking incentive compatibility, and mixed in terms of devising an optimal audit policy. We now proceed to identify a number of special cases in which we can check incentive compatibility in polynomial time. In the threshold setting, we focus on checking ε -BNIC for a UNIFORM audit policy, which we showed earlier is optimal, while in the top-k setting we focus on the UNIFORM-K audit policy. We consider, in particular, three common machine learning models for f: linear, piecewise linear, and logistic (sigmoid) functions. Throughout, we assume that distributions over types are sufficiently well behaved, in that it is tractable to compute probabilities of intervals.

We begin by showing that verification is tractable in any instance in which the CDF (CMF) of h can be computed over the set of *suspicious* agent types. As can be surmised from the complexity results regarding verification, the

"hardness" of the problem stems from determining the probability that an agent's true type is suspicious. However, when this can be computed efficiently, so can ε^* .

Theorem 11. Let $U = \{(\mathbf{x}, \mathbf{z}') \in I^d \times I^s : f(\mathbf{x}, \mathbf{z}') \geq \theta, h(\mathbf{x}, \mathbf{z}) \neq 0, \text{ and } \exists (\mathbf{x}, \mathbf{z}^*) \text{ with } f(\mathbf{x}, \mathbf{z}^*) < \theta \text{ and } h(\mathbf{x}, \mathbf{z}^*) \neq 0\}.$ If $\mathbb{P}_{\mathbf{a} \sim D}(\mathbf{a} \in U)$ can be efficiently computed, then so can ε^* .

Proof Sketch. Let $p_U = \mathbb{P}_{\mathbf{a} \sim D}(\mathbf{a} \in U)$. Suppose an agent initially scores below the threshold, then this agent's only means for allocation is to report a type in U. Moreover, UNIFORM only audits agents in U and does so uniformly. Thus, for a given realization, the more agents with true types in U, the lower the probability that the dishonest agent is audited. More specifically, suppose that some agent $\mathbf{a}_i = (\mathbf{x}_i, \mathbf{z}_i)$, with $f(\mathbf{a}_i) < \theta$, is able to falsely submit $\mathbf{a}_i' = (\mathbf{x}_i, \mathbf{z}_i')$ with $f(\mathbf{a}_i') \geq \theta$. Then, this agent's expected marginal gain is, $\mathbb{E}_{\mathcal{A}_{-i}}[u_i(\mathbf{a}_i, \mathbf{a}_i')|f, \alpha, \phi] = \mathbb{E}_{\mathcal{A}_{-i}}[1 - (1+c)\phi_i(\mathcal{A}')]$

$$= 1 - (1+c) \sum_{\ell=0}^{n-1} \binom{n-1}{\ell} p_U^{\ell} (1-p_U)^{n-\ell-1} \min \left(1, \frac{B}{\ell+1}\right)$$

Since, under UNIFORM all dishonest reports have either value 0 or value $\mathbb{E}_{\mathcal{A}_{-i}}[u_i(\mathbf{a}_i,\mathbf{a}_i')|f,\alpha,\phi]$, we need only compute this single sum, for any agent type, and have found ε . Moreover, UNIFORMis optimal and thus $\varepsilon = \varepsilon^*$.

In both the discrete and continuous case, when $\mathbb{P}(\mathbf{a} \in U)$ can be computed exactly, verification is tractable. Next, we give a sufficient condition on this, and present several tractable special cases.

Definition 8. We say a PDF h is well-behaved if h is zero on a polynomial number of s+d-dimensional maximal intervals, and over any any interval $[a,b] \subset \mathbb{R}$, the value of $\int_a^b h(\mathbf{x},\mathbf{z})dz_r$ and $\int_a^b h(\mathbf{x},\mathbf{z})dx_t$ for observed features r and unobserved features s have closed-form solutions derivable in polynomial-time w.r.t. $(n,B,s,d,\log(c))$.

Remark 1. Many commonly used distributions, such as uniform and exponential, are well-behaved. In many other common cases, such as Gaussian, we can obtain a good numerical approximation, so that the approaches below can approximately apply in these also. We formalize this below.

As we show next, in the threshold case, checking ε -BNIC is easy for piecewise linear and logistic score functions as long as the distribution over types is *well-behaved*. For top-k, we need a much stronger assumption that types are distributed uniformly to obtain comparable positive results. Each proof proceeds as follows (see the supplement for details). The score function f partitions $I^d \times I^s$ into two disjoint regions, one of which is U (the set of suspicious types). We then show that one of the two partitions it is easy to compute the CDF, as long as h is *well-behaved*. Once this is done, we have computed either p_U or $1-p_U$ and from here Theorem 11 directly implies that ε^* is computable in polynomial time.

Definition 9. A function $f: \mathbb{R}^n \to \mathbb{R}$ is said to be Piecewise Linear if for some partition of \mathbb{R}^{d+s} into disjoint rectangular regions, given by $P = \{L_1, ..., L_m\}$ the function $f|_{L_s}: \mathbb{R}^{d+s} \to \mathbb{R}$ is linear for each $L_s \in P$.

Corollary 1. Suppose the distribution of agent types is well-behaved and f is piecewise linear, or logistic, and $\alpha(f, \mathcal{A}')$ is threshold allocation. Then determining the minimum $\varepsilon \geq 0$ such that UNIFORM is ε -BNIC, can be done in polynomial time.

Corollary 2. Suppose the distribution of agent types is uniform. Suppose further that f is piecewise-linear, or logistic, and $\alpha(f, \mathcal{A}')$ is top-k allocation. Then determining the minimum $\varepsilon \geq 0$ such that UNIFORM-K is ε -BNIC, can be done in polynomial time.

For many common continuous distributions, such as Gaussian, only a numerical approximation of $p_U = \mathbb{P}(\mathbf{a} \in U)$ can be computed. Our final result is to quantify the error in ε^* , in terms of the additive numerical error γ in p_U .

Theorem 12. Suppose with error γ we have a numerical approximation $p'_U = p_U \pm \gamma$. Then we can compute $\varepsilon' = \varepsilon^* \pm (n-B)\binom{n-1}{B-1} \int_{p_U}^{p_U+\gamma} x^{B-1} (1-x)^{n-B} dx$.

Although the error term looks messy, it is tight and in general small relative to γ , which itself is also in general a small value. As an illustration, when we have error $\gamma = 4.44 {\rm E}^{-4}$, a typical absolute error for a standard Gaussian, and n = 1000, B = 250, and $p_U = 0.6$, then $\epsilon' = \epsilon^* \pm 6 {\rm E}^{-60}$.

5 Conclusion

We study the problem of auditing self-reported attributes in resource allocation settings from two perspectives: 1) the complexity of checking whether a particular audit policy is incentive compatible, and 2) characterizing and computing an audit policy that minimizes incentives to lie. We find that checking incentive compatibility is, in general, hard. However, in settings where resources are assigned by thresholding the individual's computed score, a uniform audit policy, particularly appealing for its simplicity, is optimal. In addition, we show that in two important classes of score functions, piecewise linear and logistic, we can check incentive compatibility in polynomial time under some assumptions on the distribution of agent types.

A number of open questions remain. While we show that computing an optimal audit policy in the setting where resources are allocated to the top-k scoring agents is hard, it may be possible to achieve a better approximation of optimal than what we exhibit for the uniform policy. Moreover, our model presumes that agents incur no direct costs of misreported preferences besides the endogenous costs of being audited. In practice, there may be both cognitive and tangible costs involved, and these can be considered as an extension to our model. Finally, we assume that the distribution over agent types is known a priori, whereas it likely needs to be learned from data.

Acknowledgments

This research was partially supported by the National Science Foundation (IIS-1910392, IIS-1939677, IIS-1905558, IIS-1903207, IIS-1927422, and ECCS-2020289), Army Research Office (W911NF-19-1-0241), and Amazon.

References

- Agarwal, S.; Skiba, P. M.; and Tobacman, J. 2009. Payday loans and credit cards: New liquidity and credit scoring puzzles? *American Economic Review Papers & Proceedings* 99(2): 412–17.
- Athalye, A.; Carlini, N.; and Wagner, D. 2018. Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples. In *International Conference on Machine Learning*, 274–283.
- Brown, M.; Cummings, C.; Lyons, J.; Carrión, A.; and Watson, D. P. 2018. Reliability and validity of the Vulnerability Index-Service Prioritization Decision Assistance Tool (VI-SPDAT) in real-world implementation. *Journal of Social Distress and the Homeless* 27(2): 110–117.
- Brückner, M.; and Scheffer, T. 2011. Stackelberg games for adversarial prediction problems. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 547–555.
- Brückner, M.; and Scheffer, T. 2012. Static Prediction Games for Adversarial Learning Problems. *Journal of Machine Learning Research* 13: 2617–2654.
- Carlini, N.; and Wagner, D. 2017. Towards Evaluating the Robustness of Neural Networks. In *IEEE Symposium on Security and Privacy*, 39–57.
- Chouldechova, A.; Benavides-Prado, D.; Fialko, O.; and Vaithianathan, R. 2018. A case study of algorithmassisted decision making in child maltreatment hotline screening decisions. In Conference on Fairness, Accountability and Transparency, 134–148.
- Haeringer, G. 2018. *Market Design: Auctions and Matching*. The MIT Press.
- Hardt, M.; Megiddo, N.; Papadimitriou, C.; and Wootters, M. 2016. Strategic classification. In ACM Conference on Innovations in Theoretical Computer Science, 111–122.
- Kube, A.; Das, S.; and Fowler, P. J. 2019. Allocating Interventions Based on Predicted Outcomes: A Case Study on Homelessness Services. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 622–629.
- Li, B.; and Vorobeychik, Y. 2018. Evasion-robust classification on binary domains. *ACM Transactions on Knowledge Discovery from Data*.
- Lowd, D.; and Meek, C. 2005. Adversarial Learning. In ACM SIGKDD International Conference on Knowledge Discovery in Data Mining, 641–647.
- Madry, A.; Makelov, A.; Schmidt, L.; Tsipras, D.; and Vladu, A. 2018. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*.
- Milli, S.; Miller, J.; Dragan, A.; and Hardt, M. 2019. The social costs of strategic classification. In *Conference on Fairness, Accountability, and Transparency*.
- Nisan, N.; Roughgarden, T.; Tardos, E.; and Vazirani, V. V., eds. 2007. *Algorithmic Game Theory*. Cambridge University Press.

- Papernot, N.; McDaniel, P.; Sinha, A.; and Wellman, M. 2018. Towards the Science of Security and Privacy in Machine Learning. In *IEEE European Symposium on Se*curity and Privacy.
- Šrndic, N.; and Laskov, P. 2014. Practical Evasion of a Learning-Based Classifier: A Case Study. In *IEEE Symposium on Security and Privacy*, 197–211.
- Tong, L.; Li, B.; Hajaj, C.; Xiao, C.; Zhang, N.; and Vorobeychik, Y. 2019. Improving Robustness of ML Classifiers against Realizable Evasion Attacks Using Conserved Features. In USENIX Security Symposium.
- Vorobeychik, Y.; and Kantarcioglu, M. 2018. *Adversarial Machine Learning*. Morgan and Claypool.
- Wong, E.; and Kolter, J. Z. 2018. Provable defenses against adversarial examples via the convex outer adversarial polytope. In *International Conference on Machine Learn*ing.
- Wu, T.; Tong, L.; and Vorobeychik, Y. 2020. Defending Against Physically Realizable Attacks on Image Classification. In *International Conference on Learning Repre*sentations.
- Xu, W.; Qi, Y.; and Evans, D. 2016. Automatically Evading Classifiers: A Case Study on PDF Malware Classifiers. In *Network and Distributed System Security Symposium*.

Supplementary Materials

Proof of Proposition 1

When resources are allocated via threshold, this result is straightforward. In the case of top-k, the need for auditing can be seen by a simple two agent example for any arbitrary score function. Let n=2 and k=1. Given $f:I^d\times I^s\to\mathbb{R}$, assume by way of contradiction that f induces truthful reporting to be an equilibrium, but has $f(\mathbf{x},\mathbf{z})\neq f(\mathbf{x})$. Then there exists some values $\mathbf{x}\in I^d,\mathbf{z}_1,\mathbf{z}_2\in I^s$ with $f(\mathbf{x},\mathbf{z}_1)\neq f(\mathbf{x},\mathbf{z}_2)$. WLoG assume $f(\mathbf{x},\mathbf{z}_1)>f(\mathbf{x},\mathbf{z}_2)$. Let $\mathbf{a}_1=(\mathbf{x},\mathbf{z}_1)$ and $\mathbf{a}_2=(\mathbf{x},\mathbf{z}_2)$. Since k=1, agent 2 will never receive the resource unless they report $\mathbf{z}_2'=\mathbf{z}_1\neq\mathbf{z}_2$. Therefore reporting truthful \mathbf{z} can be an equilibrium iff $\forall (\mathbf{x},\mathbf{z})\in I^d\times I^s, f(\mathbf{x},\mathbf{z})=f(\mathbf{x})$.

Proof of Theorem 1

To show the hardness of optimally auditing we will reduce from SAT. Before showing this reduction, we make the following observation about optimal auditing. Suppose an agent reports type $\mathbf{a}_i' = (\mathbf{x}_i, \mathbf{z}_i')$, by definition the principal is sure that \mathbf{x}_i is reported truthfully, but may be sure about the veracity of \mathbf{z}_i' . However, in some cases the principal may know for certain if \mathbf{z}_i' was reported truthfully. One such case, which we will make use of, is when for all $\mathbf{z} \in I^s$ with $\mathbf{z} \neq \mathbf{z}_i'$, we have $h(\mathbf{x}_i, \mathbf{z}) = 0$.

For a given Boolean Φ , over variables $b_1, \ldots b_m$, we will encode Φ into the distribution over agent types, given by D, such that it is NP-hard to determine if for some reported $(\mathbf{x}_i, \mathbf{z}_i')$ there exists another $(\mathbf{x}_i, \mathbf{z})$ with $\mathbf{z} \neq \mathbf{z}_i$ and $h(\mathbf{x}_i, \mathbf{z}) > 0$. Suppose that $I = \{0, 1\}$, i.e. agent features are binary, set d = 2, and s = m, the agents are of the form $\langle x_1, x_2, z_1, \ldots z_m \rangle$. Define the PDF of D over $I^2 \times I^s$ as

$$h(\mathbf{x}, \mathbf{z}) = \begin{cases} \frac{1}{2^{m+1}} & \text{if } x_1 = 1, x_2 = 1, \text{ and } \mathbf{z} = 1\\ \frac{1}{2^{m+1}} & \text{if } x_1 = 1, \mathbf{z} \neq 1, \text{ and } x_2 = \Phi(\mathbf{z})\\ \frac{1}{2^{m+1}} & \text{if } x_1 = 0\\ 0 & \text{otherwise} \end{cases}$$

Lastly, suppose that c = 0 (the cost of lying), B = 1 (the number of audits), and $f(\mathbf{x}, \mathbf{z}) = x_2 z_1 ... z_m$. Note that only agents with $x_2 = 1$ and $\mathbf{z} = 1$ are allocated a resource, and for non-constant n the value of lying will always be positive. Under this definition of h, we see that given $x_1 = 0$, types are distributed uniformly. However, when $x_1 = 1$ the PMF is defined conditionally on the relationship between x_2 and z. The hardness of this problem arises when the principal must decide how to audit an agent reporting (1, 1, ..., 1). If Φ has no satisfying assignments, other than potentially $\Phi(\mathbf{z}) = 1$, then an agent reporting $\langle 1, 1, \dots, 1 \rangle$ is guaranteed to be truthful. However, if Φ does have such an assignment, then the principal will have to audit that agent with nonzero probability. This is due to the fact that another agent, will have $x_1, x_2 = 1$, but $\mathbf{z} \neq 1$ and will be incentivized to lie and report z = 1. With out a satisfying assignment, no such incentivized type will exist.

Proof of Theorem 2

First note that for any realization of agents, only agents whose true type scores below the threshold, but can report

a type scoring above the threshold, will have incentive to lie. Denote the set of these agent types as

$$U = \{ (\mathbf{x}, \mathbf{z}) \in I^d \times I^s : f(\mathbf{x}, \mathbf{z}) \ge \theta \text{ with } h(\mathbf{x}, \mathbf{z}) > 0 \text{ and}$$
$$\exists \mathbf{z}' \text{ s. t. } f(\mathbf{x}, \mathbf{z}') < \theta \text{ with } h(\mathbf{x}, \mathbf{z}') > 0 \}.$$

That is, U is the set of all *suspicious*-types. For any set of reports \mathcal{A}' , UNIFORM considers reports in \mathcal{A}' to be in one of three categories; *impossible*, *suspicious*, or neither. First note that by definition, any agent regardless of their true type receives utility at most 0 from reporting a type which falls into the "neither" category, even when that agent is audited with probability 0. Hence any optimal audit policy need only focus on how to audit agents whose reports fall into the other two categories, *impossible* or *suspicious*.

Since the notion of optimality is defined in terms of ε -BNIC, we are examining the case when all agents report their type truthfully, and one agent, say agent i is considering deviating while all other agent's strategies remain fixed. As such, from the principal's perspective, there is at most one dishonest agent in any set of reports. Thus, if agents i reports type \mathbf{a}_i' and $h(\mathbf{a}_i')=0$, the principal is immediately aware of the identity of the dishonest agent. Moreover, if some other agent \mathbf{a}_j , considers falsely reporting some \mathbf{a}_j' , agent j knows that type \mathbf{a}_i' will never in appear in the set \mathcal{A}_{-j}' and thus \mathbf{a}_j 's utility of falsely reporting is independent of the probability with which any *impossible* type is audited. Therefore auditing any *impossible* type with probability 1 is optimal.

Now we need only show that the way in which UNIFORM audits *suspicious* types constitutes an optimal audit policy. Let

 $L = \{(\mathbf{x}, \mathbf{z}) \in I^d \times I^s : h(\mathbf{x}, \mathbf{z}) > 0 \text{ and } \exists \mathbf{z}' \text{ s.t. } (\mathbf{x}, \mathbf{z}') \in U\}$ That is, L is the set of all true types that could possibly report a *suspicious* type. Under UNIFORM, for any realization of agents, all agents with true type in L have the same expected value of lying when misreporting a type which scores above the threshold. That is, for any realization \mathcal{A} and for any $\mathbf{a}_i, \mathbf{a}_j \in \mathcal{A}$,

$$\mathbb{E}[u_i(\mathbf{a}_i', \mathcal{A}_{-i}')|f, \phi] = \mathbb{E}[u_i(\mathbf{a}_j', \mathcal{A}_{-j}')|f, \phi]$$
$$\forall \mathbf{a}_i', \mathbf{a}_j' \text{ with } f(\mathbf{a}_i'), f(\mathbf{a}_j') \ge \theta$$

Let $G(A_{-i})$ be the set of agents whose true type scores above the threshold. Then the previous equivalence can be seen by the fact that under threshold, the expected value of lying for agent \mathbf{a}_i when all other agents are truthful is

$$\begin{split} & \mathbb{E}[u_{i}(\mathbf{a}_{i}^{\prime},\mathcal{A}_{-i}^{\prime})|f,\phi] = \mathbb{E}_{\mathcal{A}_{-i}^{\prime}}[1-\phi_{i}(\mathcal{A}^{\prime})-c\phi_{i}(\mathcal{A}^{\prime})] \\ = & \mathbb{E}_{\mathcal{A}_{-i}}[1-\phi_{i}(\mathcal{A}_{-i}\cup\{\mathbf{a}_{i}^{\prime}\})-c\phi_{i}(\mathcal{A}_{-i}\cup\{\mathbf{a}_{i}^{\prime}\})] \\ & = 1-(1+c)\mathbb{E}_{\mathcal{A}_{-i}}\bigg[\min\bigg(1,\frac{B}{|G(\mathcal{A}_{-i})|+1}\bigg)\bigg] \\ & = 1-(1+c) \\ & \bigg(\sum_{\ell=B}^{n-1}\frac{B}{\ell+1}\binom{n-1}{\ell}\bigg)\mathbb{P}_{\mathbf{a}}\big(f(\mathbf{a})\geq 0\big)^{\ell}\big(1-\mathbb{P}_{\mathbf{a}}\big(f(\mathbf{a})\geq 0\big)\big)^{n-\ell} \\ & + \sum_{\ell=0}^{B-1}\binom{n-1}{\ell}\mathbb{P}_{\mathbf{a}}\big(f(\mathbf{a})\geq 0\big)^{\ell}\big(1-\mathbb{P}_{\mathbf{a}}\big(f(\mathbf{a})\geq 0\big)\big)^{n-\ell}\bigg). \end{split}$$

Thus the expected value of lying has no dependence on the agents particular type, or the particular type they misreport, and depends only on the fact that the agent's true type is in L and their misreported type scores above the threshold. This implies that agents \mathbf{a}_i and \mathbf{a}_j have equivalent expected values of lying. The significance of this fact is that UNIFORM induces an ε^* BNIC where the value of ε^* is tight for all agents with true type in L.

We will now leverage the tightness of ε^* in order to show that no other policy can achieve $\varepsilon BNIC$ for $\varepsilon < \varepsilon^*$. Let $\mathcal A$ again be any possible realization of agent types, which has at least one type in L, denote the agent of this type as $(\mathbf x_i, \mathbf z_i)$. Then, given any audit policy ψ and assuming other agents are truthful, the expected value of agent i misreporting type $\mathbf a_i'$, with $f(\mathbf a_i') \geq 0$, is given as

$$\mathbb{E}[u_i(\mathbf{a}_i', \mathcal{A}_{-i})|f, \psi] = 1 - (1+c)\mathbb{E}_{\mathcal{A}_{-i}}[\psi_i(\mathcal{A} \cup \{\mathbf{a}_i'\}]]$$

$$= 1 - (1+c)$$

$$\left(\mathbb{E}_{\mathcal{A}_{-i}}[\psi_i(\mathcal{A}'||G(\mathcal{A}_{-i})|+1>B]\mathbb{P}(|G(\mathcal{A}_{-i})|+1>B)\right)$$

$$+ \mathbb{E}_{\mathcal{A}_{-i}}[\psi_i(\mathcal{A}'||G(\mathcal{A}_{-i})|+1\leq B]\mathbb{P}(|G(\mathcal{A}_{-i})|+1\leq B))$$

In the above term, the probability of being audited is broken into two terms conditioned on the number of agents in $G(\mathcal{A}_{-i})$. The two events, $|G(\mathcal{A}_{-i})|+1\leq B$ and $|G(\mathcal{A}_{-i})|+1>B$ represent a partition on the possible outcomes of \mathcal{A} , meaning that ψ can be independently defined for events in the first term and events in the second term. The ability to define ψ independently over these two events is of note due to the fact that when $|G(\mathcal{A}_{-i})|+1\leq B$ the principle has enough resources to audit each agent reporting above the threshold. Thus, it is feasible to audit each agent reporting above the threshold, with probability 1, maximizes the term $\mathbb{E}_{\mathcal{A}_{-i}}[\psi_i(\mathcal{A}\cup\{\mathbf{a}_i'\}||G(\mathcal{A}_{-i})|+1\leq B]\mathbb{P}(|G(\mathcal{A}_{-i})|+1\leq B)$. This is identical to UNIFORM when $|G(\mathcal{A}_{-i})|+1\leq B$.

It remains to be shown only that UNIFORM is the unique maximizer of the term $\mathbb{E}_{\mathcal{A}_{-i}}[\psi_i(\mathcal{A} \cup \{\mathbf{a}_i'\}||G(\mathcal{A}_{-i})|+1>B]$ Assuming D is discrete (an identical argument works for continuous by replacing the sum with an integral), the expected value can be further dissected as

$$\mathbb{E}_{\mathcal{A}_{-i}} \left[\psi_i(\mathcal{A} \cup \{\mathbf{a}_i'\} \big| |G(\mathcal{A}_{-i})| + 1 > B \right]$$

$$= \sum_{A \in \binom{I^d \times I^s}{n-1} : G(A)+1 > B} \mathbb{P}(\mathcal{A}_{-i} = A) \psi_i(A \cup \{\mathbf{a}_i'\}).$$

The expected value of lying is monotonically decreasing with respect to the above term. To show the optimally of UNIFORM consider any other policy ψ that differs from ϕ . Since ψ is not UNIFORM, there must exist a set of reported types and an agent in that set for which ψ and ϕ are different. Since \mathbf{a}_i was chosen arbitrarily, suppose the agent type is \mathbf{a}_i' and the realization is some $A_1 \cup \{\mathbf{a}_i'\}$. Then we can express $\psi_i(A_1 \cup \{\mathbf{a}_i'\}) = \phi_i(A_1 \cup \{\mathbf{a}_i'\}) + \gamma_{i,1}$, for $\gamma_{i,1} > 0$. Thus, there must exist some other reported type $\mathbf{a}_j' \in A_1$ for which $\psi_j(A_1 \cup \{\mathbf{a}_i'\}) = \phi_j(A_1 \cup \{\mathbf{a}_i'\}) - \gamma_{j,1}$ for $\gamma_{j,1} > 0$. As shown previous, the expected value of lying is tight for

all agents, meaning ε is strictly greater than ε^* if the term $\psi_j(A_1 \cup \{\mathbf{a}_i'\}) = \phi_j(A_1 \cup \{\mathbf{a}_i'\}) - \gamma_{j,1}$ is not offset in the above equation, for any agent a_i which is capable of reporting type \mathbf{a}_i' . Therefore, there must exist some other realization A_2 , such that $\psi_j(A_2 \cup \{\mathbf{a}_j'\}) = \phi_j(A_2 \cup \{\mathbf{a}_j'\}) + \gamma_{j,2}$ for $\gamma_{i,2} > 0$. Continuing this line of reasoning, there must be some other agent type in A_2 which has a lower audit weight under ψ than ϕ . This continues until we take weight from an agent whose audit probability has been given greater weight, i.e. an agent we have already seen before in this weight transferring process. This can be thought of as a weighted directed graph, where the nodes are agents the edges represent how much weight is shifted from one agent under a particular realization to another agent under that same realization. By the previous reasoning, this graph has no edge whose tip does not connect to the tail of another edge, i.e. all edges are part of a cycle. Assume that \mathbf{a}'_i and \mathbf{a}'_i are part of a two cycle, identical reasoning hold for any cycle length. We can write the terms of the expected probability of being audited, for types \mathbf{a}'_i and \mathbf{a}'_j which are affected by the weight shift as follows, first for \mathbf{a}'_i

$$\mathbb{P}(\mathcal{A}_{-i} = A_1)\psi_i(A_1 \cup \{\mathbf{a}_i'\}) + \mathbb{P}(\mathcal{A}_{-i} = A_2)\psi_i(A_2 \cup \{\mathbf{a}_i'\})$$

$$= \mathbb{P}(\mathcal{A}_{-i} = A_1)(\phi_i(A_1 \cup \{\mathbf{a}_i'\}) + \gamma_{i,1})$$

$$+ \mathbb{P}(\mathcal{A}_{-i} = A_2)(\phi_i(A_2 \cup \{\mathbf{a}_i'\} - \gamma_{i,2})$$
(3)

and for \mathbf{a}_i'

$$\mathbb{P}(\mathcal{A}_{-j} = A_3)\psi_i(A_1 \cup \{\mathbf{a}_i'\}) + \mathbb{P}(\mathcal{A}_{-j} = A_4)\psi_i(A_2 \cup \{\mathbf{a}_i'\}) \\
= \mathbb{P}(\mathcal{A}_{-j} = (A_2 \setminus \{\mathbf{a}_j'\}) \cup \{\mathbf{a}_j'\})\psi_j(A_1 \cup \{\mathbf{a}_i'\}) \\
+ \mathbb{P}(\mathcal{A}_{-j} = (A_2 \setminus \{\mathbf{a}_j'\}) \cup \{\mathbf{a}_j'\})\psi_j(A_2 \cup \{\mathbf{a}_i'\}) \\
= \mathbb{P}(\mathcal{A}_{-j} = (A_2 \setminus \{\mathbf{a}_j'\}) \cup \{\mathbf{a}_j'\})(\phi_j(A_1 \cup \{\mathbf{a}_i'\}) - \gamma_{j,1}) \\
+ \mathbb{P}(\mathcal{A}_{-j} = (A_2 \setminus \{\mathbf{a}_j'\}) \cup \{\mathbf{a}_j'\})(\phi_j(A_2 \cup \{\mathbf{a}_i'\}) + \gamma_{j,2}) \\
= \mathbb{P}(\mathcal{A}_{-i} = A_1)\frac{h(\mathbf{a}_i')}{h(\mathbf{a}_j')}(\phi_j(A_1 \cup \{\mathbf{a}_i'\}) - \gamma_{j,1}) \\
+ \mathbb{P}(\mathcal{A}_{-i} = A_1)\frac{\mathbb{P}(\mathbf{a}_i')}{\mathbb{P}(\mathbf{a}_j')}(\phi_j(A_2 \cup \{\mathbf{a}_i'\}) + \gamma_{j,2}) \\
= \frac{h(\mathbf{a}_i')}{h(\mathbf{a}_j')} \left(\mathbb{P}(\mathcal{A}_{-i} = A_1)(\phi_j(A_1 \cup \{\mathbf{a}_i'\}) - \gamma_{j,1}) \\
+ \mathbb{P}(\mathcal{A}_{-i} = A_1)(\phi_j(A_2 \cup \{\mathbf{a}_i'\}) + \gamma_{j,2})\right)$$
(4)

If $\varepsilon > \varepsilon^*$, then it must be the case that for \mathbf{a}_i' and \mathbf{a}_j' Equations 3 and 4 have greater value for some $\gamma_{i,1}, \gamma_{i,2} > 0$, than $\gamma_{i,1} = \gamma_{i,2} = 0$. Equivalently, it must be the case that

for
$$\mathbf{a}_i'$$
: $\mathbb{P}(\mathcal{A}_{-i} = A_1)\gamma_{i,1} - \mathbb{P}(\mathcal{A}_{-i} = A_2)\gamma_{i,2} > 0$ for \mathbf{a}_j' :

$$\frac{h(\mathbf{a}_i')}{h(\mathbf{a}_j')} \left(\mathbb{P}(\mathcal{A}_{-i} = A_1)(-\gamma_{i,1}) - \mathbb{P}(\mathcal{A}_{-i} = A_2)(-\gamma_{i,2}) \right) > 0$$

If the first condition holds true, then $\mathbb{P}(A_{-i} = A_1)\gamma_{i,1} > \mathbb{P}(A_{-i} = A_2)\gamma_{i,2}$. However, this would imply that the

second condition is false, meaning that no policy ψ an strictly decrease $\mathbb{E}_{\mathcal{A}_{-i}} \big[\psi_i(\mathcal{A} \cup \{\mathbf{a}_i'\} \big| |G(\mathcal{A}_{-i})| + 1 > B \big]$ when compared to ϕ . Therefore, when $|G(\mathcal{A}_{-i})| + 1 > B$, the maximum expected value of lying for any agent type is achieved by ϕ and as shown in the other case, when $|G(\mathcal{A}_{-i})| + 1 \leq B$ auditing each agent above the threshold with probability 1 achieves maximum expected value of lying of 0. Therefore, UNIFORM is the optimal audit policy in the sense that for no other policy ψ , the maximum expected value of lying under ψ is lower than that of ϕ .

Proof of Theorem 3

This result is straightforward. To impalement UNIFORM on any set of reports \mathcal{A}' , the principal need only determine if there exists an *impossible* type, and if no such type exists, compute the set of *suspicious* reports in \mathcal{A}' . Checking if there exits and *impossible* type corresponds to checking if for each $\mathbf{a}_i' \in \mathcal{A}'$, that $h(\mathbf{a}_i') > 0$. Suppose for any *known* type \mathbf{x} we can check if there exists a \mathbf{z} such that $h(\mathbf{x}, \mathbf{z}) > 0$ and $f(\mathbf{x}, \mathbf{z}) < \theta$ in polynomial time. Then to check if a report, say \mathbf{a}_i' , is *suspicious* the principal need only check the existence of such a \mathbf{z} for \mathbf{x}_i , and then check if $h(\mathbf{a}_i') > 0$. Each of these can be done in polynomial time and for any set of reports \mathcal{A}' there are n such checks that need to be done. Once the principal has checked each $\mathbf{a}_i' \in \mathcal{A}'$, the set of *suspicious* agents has been determined and thus UNIFORM can be implanted in polynomial time.

Proof of Theorem 4

This reduction will be from vertex cover. The crux of this proof comes from the fact that for sufficiently small B and c, the gain from lying will not be tight among agents. Suppose agent 1 draws a type \mathbf{a}_1 which has very low probability of being in the top-k, but can report a type \mathbf{a}_1' which is almost certainly in the top-k. In contrast, suppose agent 2 has slightly less than 1/2 chance of being in the top-k, but can report a type which has only slightly more than $\frac{1}{2}$ chance of being in the top-k. Then, prior to auditing, the expected payoff of agent 1 is far greater than that of agent 2. If the principal has insufficient audit strength, i.e. small B and small c, then the expected value of agent 1 may be greater than agent 2, even if agent 2 is never audited. We will show that identifying agents similar to agent 2, those who are *suspicious* but should never be audited, is hard

Given a graph G=(V,E), let agents be of the form $\langle x_1,\dots,x_{|V|,z_1,z_2}\rangle$ where attributes are binary. Let $c=1/|V|,\,n=4,\,k=2,\,B=1,$ and D be uniform. Let

$$g(\mathbf{x}) = \bigwedge_{(v_r, v_t) \in E} (x_r \vee x_t).$$

i.e. g is an indicator of \mathbf{x} representing a vertex cover. Let $f(\mathbf{x}, \mathbf{z}) = g(\mathbf{x})z_1 + z_2$

Under this construction, agents can score values 0, 1, 2, 3. Prior to auditing the agents with the most incentive to lie will be of the form $g(\mathbf{x}) = 1, z_1 = 0, z_2 = 0$, and these agents' highest utility report will be $z_1 = z_2 = 1$. Thus any report with $g(\mathbf{x}) = z_1 = z_2 = 1$ should have the highest audit weight in any set of reports. Thus in some set of report

 $\mathcal{A}' = \{\mathbf{a}_1', \mathbf{a}_2', \mathbf{a}_3', \mathbf{a}_4'\}, \text{ if } \mathbf{a}_1' \text{ and } \mathbf{a}_2' \text{ have } g(\mathbf{x}) = 1, \mathbf{z}' = \mathbf{1}, \text{ but } \mathbf{a}_3' \text{ and } \mathbf{a}_4' \text{ have } g(\mathbf{x}) = 0, \mathbf{z}' = \mathbf{1}, \text{ the principal should never audit } \mathbf{a}_3' \text{ and } \mathbf{a}_4' \text{ if } \mathbf{a}_1' \text{ and } \mathbf{a}_2' \text{ are higher utility reports.}$ For each of these type, the *minimum* type is $\mathbf{z}^* = \mathbf{0}$. Thus, the increase in marginal gain, prior to auditing for agents 1 and 2 is Agents \mathbf{a}_1' and \mathbf{a}_2' have expected payoff proportional to the probability that any other agent score a 3. This probability is itself directly proportional to the number of vertex covers of G since $\mathbb{P}(f(\mathbf{a}) = 3) = \frac{\beta}{2^{|V|+2}}$ where β is the number of vertex covers. Thus determining the relationship between the payoffs for agent's 1, 2 and 3, 4 is equivalent to determining if the given graph has more than $2^{|V|} \frac{1}{(1+c)}$ vertex covers, for any c.

Proof of Theorem 6

Suppose the principal's objective is to induce ε dominant strategy incentive compatibility for the minimum value of ε . Then for any type a_1 , the value of falsely reporting any other type \mathbf{a}'_1 , given any set of true types \mathcal{A}_{-1} and reported types \mathcal{A}'_{-1} of the other n-1 agents, must be at most ε . Agent a1 knows the scores of each of the other agents' reports. As such, \mathbf{a}_1 knows the top-k scoring agents in \mathcal{A}'_{-i} , and knows if their true type a_1 or any false type a'_1 will score in the top-k. Thus, agent 1 has binary utility prior to auditing, either they are in the top-k or they are not. In this sense, similar to the threshold setting, all lies, barring auditing, have the same payoff. Thus the expected value of putting forth any report a'_1 that is in the top-k must have the same payoff after auditing. Otherwise, using a similar argument to the proof of Theorem 2, the audit weights could be shifted from some less incentivized type, to some more incentivized type. This would always result in a strict decrease in ε . Under UNIFORM-K all reports in the top-k have the same expected payoff since UNIFORM-K treats all suspicious reports in the top-k as being equal, regardless of their actual score.

Proof of Theorem 7

Hardness is shown via a reduction from #VC, which is concerned with counting the number of vertex covers of a given graph G=(V,E). We first show this claim for discrete agent features, and then show that the reduction can be trivially extended to include continuous features.

Let D be uniform, B=1 and agents be ${\bf a}=\langle x_1,...,x_{|V|},z_1\rangle$, for $x,z\in\{0,1\}$. Set

$$f(\mathbf{x}, z) = \left(\bigwedge_{(v_r, v_t) \in E} (x_r \vee x_t) \right) \wedge z_1.$$

Thus an report \mathbf{a}' yields $f(\mathbf{a}')=1$ if and only if the *known* type \mathbf{x} constitutes a vertex cover and z=1. In both threshold and top-k, any agent reporting \mathbf{a}' with $f(\mathbf{a}')=1$ is *suspicious*.

First, suppose we are in the threshold setting with $\theta = \frac{1}{2}$. Let $\mathbf{a}_1 = \langle 1, \dots, 1, 0 \rangle$, then this agent can simply report z = 1 to score above the threshold, and when doing so, the expected marginal gain is

$$\varepsilon = 1 - (1+c)\mathbb{E}_{\mathbf{a}_2}[\phi_1(\{\mathbf{a}_1', \mathbf{a}_2\})]$$

Auditing over the set of reports $\{\mathbf{a}_1', \mathbf{a}_2\}$, can be broken into two cases. The first, $f(\mathbf{a}_1') = 1$ and $f(\mathbf{a}_2) = 0$. In this case it is optimal for the principal to audit agent 1 with probability 1. The second case is when $f(\mathbf{a}_1') = f(\mathbf{a}_2) = 1$. In this case, it is optimal to audit both agents with probability 1/2 since both reports then have the same value of being misreports. Thus, the utility of agent 1 reporting z=1 is given by

$$\varepsilon = 1 - (1+c)(1-1/2\mathbb{P}(f(\mathbf{a}_2)=1))$$

Let β be the number of vertex covers of G. Then falsely reporting z=1 is optimal for agent 1 when

$$0 < 1 - (1+c)\left(1 - \frac{1}{2}\mathbb{P}(f(\mathbf{a}_2) = 1)\right)$$

$$\Longrightarrow \mathbb{P}(f(\mathbf{a}_2) = 1) > \frac{1}{2} - \frac{1}{2+2c}$$

$$\Longrightarrow \frac{\beta}{2^{|V|+1}} > \frac{1}{2} - \frac{1}{2+2c}$$

$$\Longrightarrow \beta > 2^{|V|}\left(1 - \frac{1}{1+c}\right)$$

Thus, $\varepsilon>0$ when $\beta>2^{|V|}\big(1-\frac{1}{1+c}\big)$. Note that for c=0 the inequality always holds, and never holds for $c=2^{|V|}$. Thus, if there existed a polynomial time algorithm to determine if $\varepsilon>0$, then it determining if $\beta>2^{|V|}\big(1-\frac{1}{1+c}\big)$ could also be done in polynomial time. If such an algorithm existed, then using binary search over $c\in\{1,\ldots,2^{|V|}\}$, the value of β , i.e. the number of vertex covers, could be found in polynomial time.

Next, we show that in the top-k a similar argument holds. The key difference is that the expected value of lying is slightly altered. Let k=1. When B=1, n=2 and f is binary, auditing can be considered in 3 cases. If both agents report $f(\mathbf{a})1$ then both should be audited with probability 1/2. If only one agent reports $f(\mathbf{a})=1$, then that agent should be audited exclusively. Lastly, when both agents report $f(\mathbf{a})=0$ both are guaranteed to be truthful, and auditing is not necessary. Assuming that ties for the resource are broken uniformly at random, the expected marginal gain of agent 1 falsely reporting z=1 is given by

$$\varepsilon = \mathbb{P}(f(\mathbf{a}_2) = 1) \frac{3(1 - 2c)}{4} - \frac{1}{2} = \frac{\beta}{2^{|V|}} \frac{3(1 - 2c)}{4} - \frac{1}{2}$$

and we can again use a similar searching technique over c to find the value of β .

In the case of continuous agent features, we can modify f such that feature values are "binned". For example, suppose that D is uniform over $[0,1]^{d+s}$. Then we can define a truncation function, $g(\mathbf{a}) = \langle \lfloor x_1 + 0.5 \rfloor, ..., \lfloor z_s + 0.5 \rfloor \rangle$. Then defining a new score function f_1 to be $f_1(a) = f(g(a))$, the problem is identical to the discrete version.

Proof of Theorem 9

Suppose for some $I^d \times I^s$, B, f, D, c, θ , and n, there exists a polynomial time algorithm that can compute the minimum ε such that the problem instance is ε -BNIC. Then this algorithm can be used to construct an optimal audit policy, specifically UNIFORM, in polynomial time. To see this we can make use of Theorem 3, which states that UNIFORM is

tractable if and only if for any *known*-type \mathbf{x} , the corresponding *minimum*-type $(\mathbf{x}, \mathbf{z}^*)$ can be determined to have scored below the threshold. Thus, for any set of reports \mathcal{A}' , we can determine how to audit each $(\mathbf{x}_i, \mathbf{z}_i') \in \mathcal{A}'$ by determining if their *minimum*-type scores below the threshold.

Computing this indicator can be accomplished by defining a new problem instance given by the same agent domain $I^d \times I^s$ with $\hat{D} = D$, $\hat{n} = n$, and $\hat{\theta} = \theta$. Further, set $\hat{B} = 1$, $\hat{c} = 0$, and $\hat{f}(\mathbf{x}, \mathbf{z}) = f(\mathbf{x}, \mathbf{z})\mathbb{I}[\mathbf{x} = \mathbf{x}_i]$. Under this score function, the only agents with possible incentive to lie will be those will known-type x_i . When auditing, the principal must consider the incentive to lie of the minimumtype. More precisely, the for a given known-type x, the agent with the most incentive to lie will be, by definition, be of the form $(\mathbf{x}, \mathbf{z}^*)$ where \mathbf{z}^* is the minimum-self-reported-type with respect to the *known*-type x. Thus for any given domain of agent types, the incentive to lie is ultimately determined by the *known*-type. In the constructed problem instance, the only agents with an incentive to lie are those with knowntype x_i . Thus if $\varepsilon > 0$ for the constructed instance, then the minimum type of x_i has incentive to lie and any report $(\mathbf{x}_i, \mathbf{z})$ with $f(\mathbf{x}_i, \mathbf{z}) \geq \theta$ should be audited.

As outlined in the proof of the optimality of UNIFORM, one need only determine the set of agents which are to be audited for any given set of reports \mathcal{A}' , rather than over the set of all possible sets of reports. Thus, when determining which agents to audit, the principal would only need to run the verification algorithm on at most n constructed instances (one for each of the unique \mathbf{x} 's in \mathcal{A}'). Thus if verification could be done in polynomial time for a given instance, optimal auditing could also be done in polynomial time for that instance.

Proof of Theorem 10

For a given agent \mathbf{a}_i , with $f(\mathbf{a}'_i)$, the expected value of falsely reporting some \mathbf{a}'_i with $f(\mathbf{a}'_i)$ is given by

$$u_i = 1 - \mathbb{E}[\phi_i(\mathcal{A}')] - c\mathbb{E}[\phi_i(\mathcal{A}')].$$

By definition, this reported type is *suspicious*. Under UNIFORM the audit probability on agent *i*'s *suspicious* report, in a given set of reports \mathcal{A}' , is $\phi_i(\mathcal{A}') = \min\left(1, \frac{B}{|G(\mathcal{A}')|}\right)$ where $G(\mathcal{A}')$ is the set of all *suspicious* reports in \mathcal{A}' . Thus,

$$\varepsilon^* = \left(1 - \mathbb{E}\left[\min(1 - \frac{B}{|G(\mathcal{A}')|}\right]\right) - c\mathbb{E}\left[\min(1 - \frac{B}{|G(\mathcal{A}')|}\right]$$

By theorem 2, we know that not only is UNIFORM optimal, but all agents with a nonzero value of lying, have the exact same value of lying. Thus, the identity of the suspicious agent (agent i) does not matter, and so we can write $|G(\mathcal{A}')| = |G(\mathcal{A}_{-i})| + 1$ Since \mathcal{A}_{-i} is drawn in accordance with D, we can sample this set as a random variable from D. Moreover, since the *minimum* type of any agent can be computed in polynomial time, we can compute $G(\mathcal{A}_{-i})$ in polynomial time as well. This is due to the fact that any suspicious report in \mathcal{A}_{-i} will have a minimum type which scores below the threshold.

Thus, the audit probability on agent i, i.e. $\min\left(1,\frac{B}{|G(\mathcal{A}_{-i}|+1)}\right)$, can be sampled by simulating the

truthful type of n-1 agents, in accordance with D, and counting the fraction of agents that spawn in $G(\mathcal{A}')$.

Let $\gamma \in \Theta(1)$ and let $\bar{\phi}$ be the empirical average of $\min\left(1,\frac{B}{|G(A'|}\right)$ after n^{γ} samples. Then, Hoeffding's inequality yields,

$$\mathbb{P}\left(\left|\bar{\phi} - \mathbb{E}\left[\min\left(1, \frac{B}{|G(\mathcal{A}')|}\right)\right]\right| \le \frac{1}{\sqrt{n^{\gamma - 1}}}\right)$$

$$\ge 1 - 2e^{-n^{\gamma} \frac{1}{n^{\gamma - 1}}} = 1 - 2e^{-2n} \ge 1 - \frac{1}{n^2}$$

Thus, by taking the sample average, $\bar{\phi}$ as an approximation of $\min\big(1,\frac{B}{|G(A'|}\big),$ we obtain

$$\begin{split} &|\varepsilon'-\varepsilon^*| \\ = & \big| \big(1 - \mathbb{E} \big[\min \big(1, \frac{B}{|G(\mathcal{A}')|} \big) \big] - c \mathbb{E} \big[\min \big(1, \frac{B}{|G(\mathcal{A}')|} \big) \big] \\ & - (1 - \bar{\phi}) - c \bar{\phi} \big| \\ = & (1 + c) \big| \big(\mathbb{E} \big[\min \big(1, \frac{B}{|G(\mathcal{A}')|} \big) \big] - \bar{\phi} \big) \big| \le (1 + c) \frac{1}{\sqrt{n^{\gamma - 1}}} \end{split}$$

with probability at least $1 - \frac{1}{n^{\gamma}}$. As stated previously, if $c \ge n$, the mechanism is trivially 0-BNIC. For 0 < c < n the dependency of c in the run-time can be removed. Therefore, we have that

$$(1+c)\frac{1}{\sqrt{n^{\gamma-1}}} \le \frac{1+\sqrt{n^2}}{\sqrt{n^{\gamma-1}}} = \Theta\left(\frac{1}{\sqrt{n^{\gamma-3}}}\right)$$

Thus the additive difference in the approximation of ε^* is at most $\Theta\left(\frac{1}{\sqrt{n^{\gamma-3}}}\right)$, with probability at least $1-\frac{1}{n^2}$.

Proof of Corollary 1

We first show this result for linear functions, i.e. m = 1. For linear f, we can write $f(\mathbf{x}, \mathbf{z}) = \mathbf{w}_1 \mathbf{x} + \mathbf{w}_2 \mathbf{z}$ for some weight vectors $\mathbf{w}_1, \mathbf{w}_2$. We omit the bias term for simplicity but it does not effect the analysis. When f is of this from a report $\mathbf{a}' = (\mathbf{x}, \mathbf{z}')$ is suspicious if $f(\mathbf{x}, \mathbf{z}') =$ $\mathbf{w}_1\mathbf{x} + \mathbf{w}_2\mathbf{z}' \geq \theta$ and there exists some \mathbf{z} with $h(\mathbf{x}, \mathbf{z}) > 0$ and $f(\mathbf{x}, \mathbf{z}') = \mathbf{w}_1 \mathbf{x} + \mathbf{w}_2 \mathbf{z} < \theta$. Thus, independent of h, we know any true type x with $w_1x \geq \theta - w_2z^*$, where $\mathbf{z}^* = \operatorname{argmin}_{\mathbf{z}} f(\mathbf{x}, \mathbf{z})$ is non-suspicious. Since the value of $\theta - \mathbf{w}_2 \mathbf{z}^*$ is independent of \mathbf{x} , the set of reports which are guaranteed to be non-suspicious are given by a separating hyperplane. We will deal with the case of $h(\mathbf{x}, \mathbf{z}) = 0$ later. Let

$$G = \{(\mathbf{x}, \mathbf{z}') \in I^d \times I^s : f(\mathbf{x}, \mathbf{z}') \ge \theta \text{ and } \exists \mathbf{z} \text{ s.t. } f(\mathbf{x}, \mathbf{z}) < \theta\}$$

Then G is the set of all reports that would be *suspicious* when h > 0. Using Theorem 11, the expected value of any suspicious agent is given by

$$1 - (1+c) \sum_{\ell=0}^{n-1} {n-1 \choose \ell} p_U^{\ell} (1-p_U)^{n-\ell-1} \min\left(1, \frac{B}{\ell+1}\right)$$

where p_U is the probability that any agent drawn from D is U, the set of suspicious types. Thus, all that is required is to compute p_U . To do this, we will first compute $p_G = \mathbb{P}(\mathbf{a} \in$ G) and then show that $\mathbf{a} \in \mathbb{G} \setminus \mathbb{U}$ can be computed in an equivalent way. Thus giving $p_U = \mathbb{P}(\mathbf{a} \in G) - \mathbf{a} \in \mathbb{G} \setminus \mathbb{U}$.

The value of $p_G=\int_G h(\mathbf{x},\mathbf{z})d\mathbf{x}\mathbf{z}$ can be computed as follows. Let $(\mathbf{x}^*,\mathbf{z}^*)=\mathrm{argmax}\{f(\mathbf{x},\mathbf{z}):(\mathbf{x},\mathbf{z})\in$ $I^d \times I^s$. Note that $(\mathbf{x}^*, \mathbf{z}^*)$ must be in G. Moreover, the region we need to integrate over can be defined as $G = [a, b]^{d+s} \cap \{(\mathbf{x}, \mathbf{z}) \in \mathbb{R}^{d+s} : f(\mathbf{x}, \mathbf{z}) \geq 0\}$. The set $R = \{(\mathbf{x}, \mathbf{z}) \in \mathbb{R}^{s+d} : f(\mathbf{x}, \mathbf{z}) = \mathbf{w}^T(\mathbf{x}, \mathbf{z}) = 0\}$ defines a separating hyperplane. For any dimension of the vector (\mathbf{x}, \mathbf{z}) say t, and WLoG assume dimension t is associated with the x component of the vector. Then we know that the bounds of integration of integration over dimension twill have either upper bound b or lower bound a, since $x_t^* \in \{a, b\}$ by the monotonicity of f. Suppose we know the lower bound is $x_t^* = b$, then the upper bound will be either $\hat{x}_t = b$, or $\hat{x}_t \in R$. The particular value of \hat{x}_t depends on the the value of the other variables. As such, we can express the bounds for the $t^{\rm th}$ dimensions as $x_t^* = a$ and $\hat{x}_t = \min(b, \hat{b}_t)$ where $\hat{b}_t = \frac{1}{w_t} (\sum_{\ell \neq t} w_\ell x_\ell + \sum_{\ell=1}^s w_\ell z \ell)$. Symmetric definitions are given for \hat{a}_t are given when the upper bound of integration is guarantied to be b. Therefore the bounds of integration for each dimension are given by intervals of the form $[x_t^*, \hat{x}_t] = [a, \min(b, \frac{1}{w_t}(\sum_{\ell \neq t} w_\ell x_\ell + \sum_{\ell=1}^s w_\ell z_\ell)]$ or $[\hat{x}_t, x_t^*] = [\max(a, \frac{1}{w_t}(\sum_{\ell \neq t} w_\ell x_\ell + \sum_{\ell=1}^s w_\ell z_\ell), b]$. For each dimension the integral can be split on the min or max, yielding a linear function. This split will yield at most d + srectangular regions and one region which is defined entirely by bounds of the form $[a+\gamma_t,\frac{1}{w_t}(\sum_{\ell\neq t}w_\ell x_\ell+\sum_{\ell=1}^sw_\ell z_\ell)]$ or $[\frac{1}{w_t}(\sum_{\ell\neq t}w_\ell x_\ell+\sum_{\ell=1}^sw_\ell z_\ell),b-\gamma_t]$ for some constant γ_t which is given by the boundary of the rectangular regions. Therefore, the integral of h over any of the t dimensions is computable and G is able to be broken down into s+d+1 regions which each contain s+d simple integrals and therefore the integral over the region G is computable in polynomial time. Hence the maximum expected value of lying for any agent can be computed as

$$\varepsilon = 1 - (1+c) \sum_{\ell=0}^{n-1} \binom{n}{\ell+1} \left(\int_G h(\mathbf{x}, \mathbf{z}) d(\mathbf{x}, \mathbf{z}) \right)^{\ell+1}$$
$$\left(1 - \int_G h(\mathbf{x}, \mathbf{z}) d(\mathbf{x}, \mathbf{z}) \right)^{n-\ell-1} \min\left(1, \frac{B}{\ell+1} \right)$$

when h > 0. However, there may be some polynomial number of intervals $[a_1, a_2]^{d+s}$ over which h = 0. In this case, if $[a_1, a_2]^{d+s} \cap \{(\mathbf{x}, \mathbf{z}) : f(\mathbf{x}, \mathbf{z}) < \theta\}$, then the set G may have non-suspicious types and we have over counted the value of p_U . However, since these areas where h=0 are given in intervals. The sets over which we have miscounted p_U are also in intervals. As such we can use the exact technique used to find p_G to find $p_{G \setminus U}$ by integrating h over each such interval. Thus we can find $p_U = p_G - p_{G \setminus U}$ for linear f. For

To generalize this result to piece-wise linear functions we show that the above process for linear functions can be performed m number of times, where m is the number of piecewise regions over which f is defined. For each $L_t \in P$, define $f\big|_{L_s}$ as f_s , with the understanding that f_s is only applied to elements in L_t . Determining the minimum value of ε such that UNIFORM is ε -BNIC is equivalent again to determining the measure of $G=\{(\mathbf{x},\mathbf{z})\subset I^d\times I^s: f(\mathbf{x},\mathbf{z})\geq 0\}$ with respect to h. The expected value of lying for any agent whose true type scores below the threshold is again The value of $\int_G h(\mathbf{x},\mathbf{z})dG$, can again be computed by expressing the boundary conditions of G as simple limits of integration. The key difference in the setting of piecewise linear functions is that G is no longer contiguous and thus must first be broken down into contagious regions before integrating. For each of the partitions $L_s \in P$ we can again have a separating hyperplane $R_s = \{(\mathbf{x},\mathbf{z}) \in \mathbb{R}^{d+s}: f(\mathbf{x},\mathbf{z}) = \mathbf{w}^T(\mathbf{x},\mathbf{z}) = 0\}$. Over each of these regions we again know that the set value $(\mathbf{x}^*,\mathbf{z}^*)_s = \operatorname{argmax}_{(\mathbf{x},\mathbf{z})\in L_s} f(\mathbf{x},\mathbf{z})$ is in $L_s \cap R_s$. From here, since each partition is a rectangular region. Which will gives an efficient method to compute each integral of the form $\int_{G \cap L_s} h(\mathbf{x},\mathbf{z}) d(\mathbf{x},\mathbf{z})$.

$$\int_G h(\mathbf{x}, \mathbf{z}) d(\mathbf{x}, \mathbf{z}) = \sum_{s=1}^m \int_{G \cap L_s} h(\mathbf{x}, \mathbf{z}) d(\mathbf{x}, \mathbf{z})$$

As per the linear case, we have again over counted the set of *suspicious* types as G. But using an identical argument to the linear case, we can again construct some polynomial number of intervals which constitute the set $G \setminus U$. Integrating h over each of these regions gives us the difference in the measure of G and U. Thus we can compute $p_U = p_G - p_{G \setminus U}$, and by Theorem 11, we have that ε can be computed efficiently.

Proof of Corollary 2

This proof follows a similar line of reasoning to the proof of Theorem 11. As discussed previously, when $\frac{(1+c)B}{k} \geq 1$ UNIFORMis BNIC for all agents types independent of f and D. So assume that $\frac{(1+c)B}{k} < 1$ Let T_k be the set of the highest k scoring agents, then the expected value of an agent with true type \mathbf{a}_i , misreporting type \mathbf{a}_i' under UNIFORM can be expressed as

$$\mathbb{E}_{\mathcal{A}_{-i}}[\alpha_i(f, \mathcal{A}')(1 - \phi_i(\mathcal{A}')) - c\phi_i(\mathcal{A})' - \alpha_i(f, \mathcal{A})]$$

=\mathbb{P}(\mathbf{a}'_i \in T_k)(1 - (1 + c)\frac{B}{k}) - \mathbb{P}(\mathbf{a}_i \in T_k)

For any reported type \mathbf{a}_i' the value of $\mathbb{P}(\mathbf{a}_i' \in T_k)$ can treated as the CDF of a binomial random variable. For simplicity, suppose agents assume worst case tie-breaking, identical analysis holds for other simple tie breaking schemes such as random, and best case, tie breaking. If ties are broken in the worst case for agents, then an \mathbf{a}_i' receives the resource if there are at most k-1 agents with scores at least $f(\mathbf{a}_i')$. Thus $\mathbb{P}(\mathbf{a}_i' \in T_k)$ is associated with a binomial random variable with n-1 trials and probability of success equal to $\mathbb{P}(f(\mathbf{x},\mathbf{z}) \geq f(\mathbf{a}_i'))$. This mirrors the technique used in threshold allocation and we use an identical technique to compute those probabilities.

Once the value of $\mathbb{P}(\mathbf{a}_i' \in T_k)$ and $\mathbb{P}(\mathbf{a}_i \in T_k)$ are known we need only determine which agents has the highest incentive to lie. Since f is linear we can write $f(\mathbf{x}, \mathbf{z}) = \mathbf{w}_1^T \mathbf{x} + \mathbf{w}_2^T \mathbf{z}$. The domain of agent types is bounded and thus there exists $\mathbf{x}_{\max} = \max_{\mathbf{x}} \mathbf{w}_1^T$ and $\mathbf{x}_{\min} = \min_{\mathbf{x}} \mathbf{w}_1^T \mathbf{x}$. More over there also exists $\mathbf{z}_{\max} = \max_{\mathbf{z}} \mathbf{w}_2^T \mathbf{z}$ and $\mathbf{z}_{\min} = \min_{\mathbf{z}} \mathbf{w}_2^T \mathbf{z}$. Since D is uniform the $\mathbb{P}((\mathbf{x}_1, \mathbf{z}_{\max}) \in$

 $T_k) - \mathbb{P}((\mathbf{x}_1, \mathbf{z}_{\min}) \in T_k) = \mathbb{P}((\mathbf{x}_2, \mathbf{z}_{\max}) \in T_k) - \mathbb{P}((\mathbf{x}_2, \mathbf{z}_{\min}) \in T_k)$ for any $\mathbf{x}_1, \mathbf{x}_2$. More over, since f is continuous, $\mathbb{P}((\mathbf{x}, \mathbf{z}_{\max}) \in T_k) - \mathbb{P}((\mathbf{x}, \mathbf{z}_{\min}) \in T_k)$ will take on all values in the interval

$$\begin{split} \big[\mathbb{P}((\mathbf{x}_{\text{max}}, \mathbf{z}_{\text{max}}) \in T_k) - \mathbb{P}((\mathbf{x}_{\text{max}}, \mathbf{z}_{\text{min}}) \in T_k) \\ \mathbb{P}((\mathbf{x}_{\text{min}}, \mathbf{z}_{\text{max}}) \in T_k) - \mathbb{P}((\mathbf{x}_{\text{min}}, \mathbf{z}_{\text{min}}) \in T_k) \big]. \end{split}$$

We are interested in finding the agent type $(\mathbf{x}, \mathbf{z}_{\min})$ which has the the most incentive to lie. i.e. the largest gain for submitting $(\mathbf{x}, \mathbf{z}_{\max})$. We can express the value of this type as $v_1(1-(1+c)\frac{B}{k})-v_2$ Since this value is linear, and f (which determines v_1 , and v_2 is linear) this value is maximized at one of the extremes. Thus $v_1 = \mathbb{P}((\mathbf{x}_{\text{max}}, \mathbf{z}_{\text{max}}) \in T_k)$ and $v_2 = \mathbb{P}((\mathbf{x}_{\text{max}}, \mathbf{z}_{\text{min}}) \in T_k), \text{ or } v_1 = \mathbb{P}((\mathbf{x}_{\text{min}}, \mathbf{z}_{\text{max}}) \in T_k)$ and $v_2 = \mathbb{P}((\mathbf{x}_{\min}, \mathbf{z}_{\min}) \in T_k)$. Each of these values are computable in polynomial time, by directly taking the ideas in the proof of Corollary 1 and integrating the measure of h over the set $G = \{(\mathbf{x}, \mathbf{z}') :$ $f(\mathbf{x}, \mathbf{z}') \geq \theta$ and $\exists \mathbf{z}$ s.t. $f(\mathbf{x}, \mathbf{z}) < \theta$, where θ $\{f(\mathbf{x}_{\max}, \mathbf{z}_{\max}), f(\mathbf{x}_{\max}, \mathbf{z}_{\min}), f(\mathbf{x}_{\min}, \mathbf{z}_{\max}), f(\mathbf{x}_{\min}, \mathbf{z}_{\min})\}$ Once p_G is computed for each value. We know the probability that that each of these types scores in the top-k, via a binomial CDF with n-1 trials, and success rate p_G . Thus we have the values of v_1, v_2 , which give us the expected value of lying for the most incentivzed agent, i.e. $v_1(1-(1+c)\frac{B}{k})-v_2.$

Proof of Corollary ??

This result follows directly from the proof of Corollary 1 and Theorem 11. The key difference being that agents are now scored via a sigmoid function, rather than a linear function. However, since allocation decisions are made via a threshold, i.e. $f(\mathbf{x}, \mathbf{z}) \geq \theta$, sigmoid functions are equivalent to linear functions in the following sense

$$\theta = f(\mathbf{x}, \mathbf{z}) = \frac{1}{e^{\mathbf{w}^T(\mathbf{x}, \mathbf{z})} + 1} \iff \mathbf{w}^T(\mathbf{x}, \mathbf{z}) = \log(1/\theta - 1)$$

Thus, we can map any problem instance with a sigmoidal scoring function and threshold θ , to a problem instance with a linear scoring function and threshold $\log(1/\theta-1)$ Which again give the same hyper plane as in Corollary 1 and the proof follows identically from there.

Proof of Corollary ??

This is a direct result of the proofs from Corollary 2 and ??.

Proof of Theorem 12

As shown in several of the other proofs, the expected value of any agent is either 0 or is given by a the term

$$\varepsilon^* = 1 - (1+c) \sum_{\ell=0}^{n-1} \binom{n-1}{\ell} p_U^{\ell} (1-p_U)^{n-\ell-1} \min(1, B/\ell+1)$$

$$= 1 - (1+c) \left(\sum_{\ell=0}^{B} \binom{n-1}{\ell} p_U^{\ell} (1-p_U)^{n-\ell-1} + \sum_{\ell=B+1}^{n-1} \binom{n-1}{\ell} p_U^{\ell} (1-p_U)^{n-\ell-1} \frac{B}{\ell+1} \right)$$

If not from the multiplicative term $\frac{B}{\ell+1}$, the summation would constitute s a binomial sum, which would sum to 1. Using this fact, the complementary term $\sum_{\ell=B+1}^{n-1} \binom{n-1}{\ell} p_U^\ell (1-p_U)^{n-\ell-1} \frac{\ell+1-B}{\ell+1}$ can be used to simplify ε^* to

$$1 - (1+c)\left(1 - \sum_{\ell=B+1}^{n-1} \binom{n-1}{\ell} p_U^{\ell} (1-p_U)^{n-\ell-1} \frac{\ell+1-B}{\ell+1}\right)$$

$$\geq 1 - (1+c)\left(1 + \frac{n-B}{n}\sum_{\ell=B+1}^{n-1} \binom{n-1}{\ell}p_U^{\ell}(1-p_U)^{n-\ell-1}\right)$$

Using the standard technique of mapping a binomial CMF to a beta CDF, by taking the derivative of a binomial CDF w.r.t. p_U , and then reintegrating (via integration by parts), we can rewrite this term again as

$$1 - (1+c)\left(1 + (n-B)\binom{n-2}{B-1}\int_0^{p_U} x^{B-1}(1-x)^{n-B}\right)$$

Once in this form, suppose we have a numerical error in our calculation of p_U , say $\pm \gamma$. In the case when γ is positive, the negative case follows symmetrically, we can write ε' , the approximation of ε^* , as

$$1 - (1+c)\left(1 + (n-B)\binom{n-2}{B-1}\left(\int_0^{p_U} x^{B-1} (1-x)^{n-B} + \int_{p_U}^{p_U+\gamma} x^{B-1} (1-x)^{n-B}\right)\right)$$

Thus we can write

$$\begin{split} |\varepsilon' - \varepsilon^*| \\ &\leq (1+c)(n-B) \binom{n-2}{B-1} \left(\int_0^{p_U} x^{B-1} (1-x)^{n-B} \right. \\ &+ \int_{p_U}^{p_U + \gamma} x^{B-1} (1-x)^{n-B} \right) \\ &- (1+c)(n-B) \binom{n-2}{B-1} \int_0^{p_U} x^{B-1} (1-x)^{n-B} \\ &= (1+c)(n-B) \binom{n-2}{B-1} \int_{p_U}^{p_U + \gamma} x^{B-1} (1-x)^{n-B} \\ &\leq (n-B) \binom{n-1}{B-1} \int_{p_U}^{p_U + \gamma} x^{B-1} (1-x)^{n-B} \end{split}$$

Where the finally inequality comes from the assumption that B(c+1) < n, since otherwise incentive compatibility is trivially achieved by auditing *all* agents uniformly.

Thus, the additive error in ε' , when p_U has additive error γ , is no more than $(n-B)\binom{n-1}{B-1}\int_{p_U}^{p_U+\gamma}x^{B-1}(1-x)^{n-B}$. Although this error is not given in the most compact manner, it does offer some intuition as to the relative size of the error with respect to γ . The error term is roughly the the probability that a random beta variable is between p_U and $p_U+\gamma$.