

Lecture 2

Lecturer: Aaron Roth

Scribe: Aaron Roth

Privacy as a Desideratum – Making the Exponential Mechanism Truthful with Payments.

1 Introduction

In the last lecture, we saw (one instance of) how differential privacy can be used as a tool in mechanism design. We used it to design a prior-free near-revenue optimal digital goods auction. As this course continues, we will see several other instances in which privacy is used purely as a tool to solve other mechanism design problems, but today I want to shift gears and think of differential privacy itself as a goal of the mechanism designer.

Why might a mechanism designer care about providing privacy for the participants in a mechanism? The simplest answer is that participants *might not participate* if their privacy is not protected. But why? – Shouldn't we just model all concerns of the agents in their utility function?

The issue is that this might be difficult to do. Ideally we can come up with a utility function that accurately models an agent's gains in some restricted setting: for example, in an auction, we imagine the agent has some utility v_i for receiving an item, which might accurately capture the agent's true utility while interacting with the auctioneer. However, an agent might unbeknownst to the auctioneer be participating in some other interaction tomorrow. For example, the agent might be involved in a negotiation for some other good, from some other buyer. If the agent's value for the good he is purchasing *today* is correlated with his value for the good he might be purchasing *tomorrow*, then he has a very real concern about how information leaked about his interaction today might affect his interaction tomorrow. We might worry that this unmodeled utility might substantially affect the agents incentives in the otherwise dominant strategy truthful mechanism that we are running today¹

We can view privacy as a way of making mechanisms *robust to unmodeled future utility* that might be affected by agents' participation in a mechanism today. Suppose that we are designing a mechanism $\mathcal{M} : \mathcal{T}^n \rightarrow \mathcal{O}$, and we know that agents have utility functions $u_i^1(o) \equiv u^1(t_i, o)$ over outcome space \mathcal{O} . However, suppose that (unbeknownst to us), these agents will also experience some outcome $o'_i \in \mathcal{O}'_i$ tomorrow, that will be chosen in a way that might *depend* on the outcome that our mechanism chooses today. We make no assumptions about how this event is chosen – we imagine only that there is some (randomized) function $f_i(o) : \mathcal{O} \rightarrow \mathcal{O}'_i$ that determines the distribution over outcomes $o' \in \mathcal{O}'_i$. In total, agents obtain utility $u_i(o, o') = u_i^1(o) + u_i^2(o')$. We observe that if the mechanism we design \mathcal{M} is dominant strategy truthful with respect to the known utility functions u_i^1 , then it remains ϵ -dominant strategy truthful even with respect to the *unknown* utility function u_i . This is important, since we must always expect that our models are at most good approximations of people's true utilities.

Claim 1 *Suppose $\mathcal{M} : \mathcal{T}^n \rightarrow \mathcal{O}$ is dominant strategy truthful with respect to utility functions $u_i^1 : \mathcal{O} \rightarrow \mathbb{R}$. Then if \mathcal{M} is also ϵ -differentially private, it is also ϵ -approximately truthful with respect to the utility function $u_i : \mathcal{O} \times \mathcal{O}'_i \rightarrow \mathbb{R}$ defined as $u_i(o, o') = u_i^1(o) + u_i^2(o')$, for any \mathcal{O}'_i and for any utility function $u_i^2 : \mathcal{O}'_i \rightarrow [0, 1]$ when o' is chosen according to a distribution $f(o)$, where f is any function $f : \mathcal{O} \rightarrow \Delta \mathcal{O}'_i$.*

Proof Consider agent i 's expected utility for reporting truthfully:

$$\mathbb{E}_{o \sim \mathcal{M}(t), o' \sim f(o)}[u_i(o, o')] = \mathbb{E}_{o \sim \mathcal{M}(t)}[u_i^1(o)] + \mathbb{E}_{o' \sim f(\mathcal{M}(t))}[u_i^2(o')]$$

¹One example of a similar real-world concern arises in the sale of sugar-beet contracts in Denmark. Danish farmers sell sugar beets to Danisco, a monopolist sugar producer. Because this is a repeated interaction, the farmers are concerned about their true valuations becoming known. As a result, the auction is currently implemented using secure multiparty computation (which is a technology which addresses similar, but partially orthogonal privacy concerns to differential privacy) [FKOS13]

$$\begin{aligned}
&\geq \mathbb{E}_{o \sim \mathcal{M}(t_{-i}, t'_i)}[u_i^1(o)] + \exp(-\epsilon) \mathbb{E}_{o' \sim f(\mathcal{M}(t_{-i}, t'_i))}[u_i^2(o')] \\
&\geq \mathbb{E}_{o \sim \mathcal{M}(t_{-i}, t'_i), o' \sim f(o)}[u_i(o, o')] - \epsilon
\end{aligned}$$

Here we bound the change in the expected value u_i^1 using the dominant strategy truthfulness condition, and the change in the expected value of u_i^2 using differential privacy. We have implicitly used the simple (but extremely useful) fact that states that for any ϵ differentially private mechanism $\mathcal{M} : \mathcal{T}^n \rightarrow \mathcal{O}$, and for any randomized mapping $f : \mathcal{O} \rightarrow \mathcal{O}'$, the mechanism $f(\mathcal{M}(\cdot)) : \mathcal{T}^n \rightarrow \mathcal{O}'$ is ϵ -differentially private. ■

2 Designing Truthful, Private Mechanisms

Now that we have established that differential privacy imparts a desirable robustness property to our mechanisms, we might ask how to design mechanisms that incorporate it. Specifically, given some social choice problem, we might ask for a mechanism \mathcal{M} that satisfies (at least) the following 3 desiderata:

1. \mathcal{M} should be (exactly) dominant strategy truthful.
2. $\mathcal{M}(t)$ should choose an outcome $o \in \mathcal{O}$ that (approximately) maximizes social welfare, defined as $SW(t) = \sum_{i=1}^n u(t_i, o)$
3. \mathcal{M} should be ϵ -differentially private.

(We can also ask for some other nice properties, like individual rationality and a no-deficit property. You can check that the mechanism we give will satisfy these too (at least in expectation), but I won't worry about them explicitly in these notes).

So how can we do all this? For the first two, lets recall our old-favorite, the (slightly generalized) VCG mechanism. Here, we imagine bidders have *quasilinear* utility functions. They have some valuation function $v_i(o) \equiv v(t_i, o) : \mathcal{T} \times \mathcal{O} \rightarrow \mathbb{R}$. If outcome o occurs and they are asked to pay p_i , then they experience *utility* $u_i(o, p_i) = v_i(o) - p_i$.

Algorithm 1 The Generalized VCG Mechanism, parameterized by a range $\mathcal{O}' \subset \mathcal{O}$ and a function $\beta : \mathcal{O} \rightarrow \mathbb{R}$.

VCG($t; \mathcal{O}', \beta$)

OUTPUT $o^* \in \arg \max_o (\sum_{i=1}^n v(t_i, o) + \beta(o))$

Let $o_{-i}^* \in \arg \max_o (\sum_{j \neq i}^n v(t_j, o) + \beta(o))$

Charge Agent i :

$$p_i = \left(\sum_{j \neq i} v(t_j, o_{-i}^*) + \beta(o_{-i}^*) \right) - \left(\sum_{j \neq i} v(t_j, o^*) + \beta(o^*) \right)$$

Note that this mechanism is ever so slightly “generalized” beyond the standard statement of the VCG mechanism in two ways. First, it is not maximizing social welfare, but rather the sum of social welfare, and some input-independent function $\beta : \mathcal{O} \rightarrow \mathbb{R}$. Second, its range is not necessarily the whole outcome space \mathcal{O} , but possibly some restricted outcome space $\mathcal{O}' \subseteq \mathcal{O}$. We show the simple fact that this mechanism is dominant strategy truthful.

Theorem 2 For any $\mathcal{O}' \subseteq \mathcal{O}$ and for any $\beta : \mathcal{O} \rightarrow \mathbb{R}$, **VCG**($t; \mathcal{O}', \beta$) is dominant strategy truthful.

Proof Note that the payment rule used by this mechanism has the form $p_i = h(t_{-i}) - \left(\sum_{j \neq i} v(t_j, o^*) + \beta(o^*) \right)$, where h is some function that is independent of player i 's report. We will prove that any mechanism (i.e. generalized Groves mechanisms) that have this form are truthful, no matter what $h()$ is².

Consider player i 's utility for truth telling. Let (o^*, p_i) be the outcome payment pair $(o^*, p_i) = \text{VCG}(t; \mathcal{O}', \beta)$ that results when player i truthfully reports:

$$\begin{aligned} u_i(o^*, p_i) &= v(t_i, o) - p_i \\ &= v(t_i, o) + \left(\sum_{j \neq i} v(t_j, o^*) + \beta(o^*) \right) - h(t_{-i}) \\ &= \sum_{j=1}^n v(t_j, o^*) + \beta(o^*) - h(t_{-i}) \end{aligned}$$

But o^* is defined to be the outcome that maximizes $\sum_{j=1}^n v(t_j, o^*) + \beta(o^*)$, and $h(t_{-i})$ is independent of player i 's reported type, and so truthtelling is the strategy that maximizes $u_i(o^*, p_i)$. ■

Remark Note that everything above goes through if we enlarge the outcome space to include *lotteries* over outcomes: $\mathcal{O}' \subseteq \Delta \mathcal{O}$, and we extend players valuation functions to lotteries by defining them to be: $v_i(D) = E_{o \sim D}[v_i(o)]$ for any distribution $D \in \Delta \mathcal{O}$. Here we assume bidders are risk neutral and evaluate lotteries according to their expected value.

So we know how to get dominant strategy truthful mechanisms. Is there a choice of \mathcal{O}' and β that makes the VCG mechanism private, while continuing to output a high welfare outcome?

Lets think about it in the reverse direction. We already know a mechanism that is private, and that outputs a high welfare outcome: the exponential mechanism, instantiated with the social-welfare quality score:

Algorithm 2 The Exponential Mechanism with Social Welfare Quality Score

$\mathcal{M}_E(t; \mathcal{O}, \epsilon)$:

Output $o \in \mathcal{O}$ with probability:

$$\frac{\exp\left(\frac{\epsilon \sum_{i=1}^n v(t_i, o)}{2}\right)}{Z(x)}$$

where $Z(x) = \sum_{o \in \mathcal{O}} \exp\left(\frac{\epsilon \sum_{i=1}^n v(t_i, o)}{2}\right)$

Here we have simply taken the exponential mechanism and plugged in quality score $q(t, o) = \sum_{i=1}^n v_i(t_i, o)$, noting that this quality score has sensitivity 1 if $v_i : \mathcal{T} \times \mathcal{O} \rightarrow [0, 1]$. Recall that this mechanism also outputs an element $o \in \mathcal{O}$ with high probability that achieves high welfare:

Theorem 3 ([MT07]) Define $OPT(t) = \max_{o^* \in \mathcal{O}} \sum_{i=1}^n v(t_i, o^*)$ denote the optimal social welfare when player types are t . Let $o = \mathcal{M}_E(t; \mathcal{O}, \epsilon)$. Then:

$$\Pr\left[\sum_{i=1}^n v(t_i, o) \leq OPT(t) - \frac{2}{\epsilon} (\log |\mathcal{O}| + r)\right] \leq \exp(-r)$$

²The fact that we have chosen $h(t_{-i}) = \left(\sum_{j \neq i} v(t_j, o_{-i}^*) - \beta(o_{-i}^*) \right)$ is what makes this a generalized VCG mechanism. This choice ensures that the algorithm also always has non-negative payments, and is individually rational (i.e. never causes any player to have negative utility). However, we won't talk about those properties here.

This is very good if we expect $\text{OPT}(t)$ to be growing with n , since the loss from the optimal welfare is independent of n . But can this mechanism be made truthful?

We follow [HK12] and argue that in fact, the exponential mechanism instantiated with the social welfare quality score is an instantiation of the generalized VCG mechanism.

Theorem 4 Let $\mathcal{D}_E \in \Delta\mathcal{O}$ denote the distribution over outcomes induced by the exponential mechanism $\mathcal{M}_E(t, \mathcal{O}, \epsilon)$.

$$\mathcal{D}_E \in \arg \max_{\mathcal{D} \in \Delta\mathcal{O}} \sum_{i=1}^n \mathbb{E}_{o \sim \mathcal{D}}[v(t_i, o)] + \frac{2}{\epsilon} H(\mathcal{D})$$

where $H(\mathcal{D})$ is the Shannon Entropy of distribution \mathcal{D} : $H(\mathcal{D}) = \sum_{o \in \mathcal{O}} \Pr_{\mathcal{D}}[o] \cdot \log \frac{1}{\Pr_{\mathcal{D}}[o]}$.

Remark We remark that this theorem implies that the exponential mechanism (instantiated with the social-welfare quality score) is an instantiation of the VCG mechanism, with range equal to lotteries over \mathcal{O} , $\Delta\mathcal{O}$, and $\beta(\mathcal{D}) = \frac{2}{\epsilon} H(\mathcal{D})$. In other words, given a set of reports $t \in \mathcal{T}^n$, the exponential mechanism is exactly the distribution that maximizes (expected) social welfare, plus the *entropy* of the distribution. For risk-neutral agents (i.e. agents who care only for their expected utility), sampling an outcome from this distribution does not affect the incentive properties.

To prove this theorem, we will make use of Jensen's inequality:

Lemma 5 (Jensen's inequality) For any real concave function ϕ , numbers x_1, \dots, x_n in its domain, and positive real numbers a_1, \dots, a_n :

$$\frac{\sum_{i=1}^n a_i \phi(x_i)}{\sum_{i=1}^n a_i} \leq \phi \left(\frac{\sum_{i=1}^n a_i x_i}{\sum_{i=1}^n a_i} \right)$$

That, together with the fact that $\log(x)$ is a concave function, is all we need to know to prove Theorem 4!

Proof We will equivalently show that \mathcal{D}_E maximizes $\frac{\epsilon}{2} \mathbb{E}_{o \sim \mathcal{D}}[\sum_{i=1}^n v(t_i, o)] + H(\mathcal{D})$. First, note that for any distribution $\mathcal{D} \in \Delta\mathcal{O}$:

$$\begin{aligned} \frac{\epsilon}{2} \mathbb{E}_{o \sim \mathcal{D}}[\sum_{i=1}^n v(t_i, o)] + H(\mathcal{D}) &= \frac{\epsilon}{2} \sum_{o \in \mathcal{O}} \Pr_{\mathcal{D}}[o] \cdot \sum_{i=1}^n v(t_i, o) + \sum_{o \in \mathcal{O}} \Pr_{\mathcal{D}}[o] \cdot \log \frac{1}{\Pr_{\mathcal{D}}[o]} \\ &= \sum_{o \in \mathcal{O}} \Pr_{\mathcal{D}}[o] \cdot \left(\frac{\epsilon}{2} \sum_{i=1}^n v(t_i, o) + \log \frac{1}{\Pr_{\mathcal{D}}[o]} \right) \\ &= \sum_{o \in \mathcal{O}} \Pr_{\mathcal{D}}[o] \cdot \log \left(\frac{\exp \left(\frac{\epsilon}{2} \sum_{i=1}^n v(t_i, o) \right)}{\Pr_{\mathcal{D}}[o]} \right) \\ &\leq \log \left(\sum_{o \in \mathcal{O}} \Pr_{\mathcal{D}}[o] \cdot \frac{\exp \left(\frac{\epsilon}{2} \sum_{i=1}^n v(t_i, o) \right)}{\Pr_{\mathcal{D}}[o]} \right) \\ &= \log \left(\sum_{o \in \mathcal{O}} \exp \left(\frac{\epsilon}{2} \sum_{i=1}^n v(t_i, o) \right) \right) \end{aligned}$$

where the inequality follows from Jensen's inequality, and the fact that $\sum_{o \in \mathcal{O}} \Pr_{\mathcal{D}}[o] = 1$

If we can show that the distribution \mathcal{D}_E attains this bound, then we will have proven that $\mathcal{D}_E \in \arg \max_{\mathcal{D} \in \Delta\mathcal{O}} \frac{\epsilon}{2} \mathbb{E}_{o \sim \mathcal{D}}[\sum_{i=1}^n v(t_i, o)] + H(\mathcal{D})$, which is what we want. So lets plug in \mathcal{D}_E and see what we get. Recall that:

$$\Pr_{\mathcal{D}_E}[o] = \frac{\exp\left(\frac{\epsilon \cdot \sum_{i=1}^n v(t_i, o)}{2}\right)}{\sum_{o \in \mathcal{O}} \exp\left(\frac{\epsilon \cdot \sum_{i=1}^n v(t_i, o)}{2}\right)}$$

Therefore, we have:

$$\begin{aligned} \frac{\epsilon}{2} \mathbb{E}_{o \sim \mathcal{D}_E} \left[\sum_{i=1}^n v(t_i, o) \right] + H(\mathcal{D}_E) &= \sum_{o \in \mathcal{O}} \Pr_{\mathcal{D}_E}[o] \cdot \log \left(\frac{\exp\left(\frac{\epsilon}{2} \sum_{i=1}^n v(t_i, o)\right)}{\Pr_{\mathcal{D}_E}[o]} \right) \\ &= \sum_{o \in \mathcal{O}} \Pr_{\mathcal{D}_E}[o] \cdot \log \left(\sum_{o \in \mathcal{O}} \exp\left(\frac{\epsilon \cdot \sum_{i=1}^n v(t_i, o)}{2}\right) \right) \\ &= \log \left(\sum_{o \in \mathcal{O}} \exp\left(\frac{\epsilon}{2} \sum_{i=1}^n v(t_i, o)\right) \right) \end{aligned}$$

which completes the proof. ■

To sum up, we have the following mechanism:

PrivateVCG($t; \mathcal{O}, \epsilon$):

Let \mathcal{D}^* be the distribution such that:

$$\Pr_{\mathcal{D}^*}[o] = \frac{\exp\left(\frac{\epsilon \cdot \sum_{i=1}^n v(t_i, o)}{2}\right)}{Z(x)}$$

where $Z(x) = \sum_{o \in \mathcal{O}} \exp\left(\frac{\epsilon \cdot \sum_{i=1}^n v(t_i, o)}{2}\right)$.

Output $o \sim \mathcal{D}^*$.

Let \mathcal{D}_{-i}^* be the distribution such that:

$$\Pr_{\mathcal{D}_{-i}^*}[o] = \frac{\exp\left(\frac{\epsilon \cdot \sum_{j \neq i}^n v(t_j, o)}{2}\right)}{Z_{-i}(x)}$$

where $Z_{-i}(x) = \sum_{o \in \mathcal{O}} \exp\left(\frac{\epsilon \cdot \sum_{j \neq i}^n v(t_j, o)}{2}\right)$.

Charge Agent i

$$p_i = \left(\sum_{j \neq i} \mathbb{E}_{o \sim \mathcal{D}_{-i}^*} [v(t_j, o)] + H(\mathcal{D}_{-i}^*) \right) - \left(\sum_{j \neq i} \mathbb{E}_{o \sim \mathcal{D}^*} [v(t_j, o)] + H(\mathcal{D}^*) \right)$$

For any $v : \mathcal{T} \times \mathcal{O} \rightarrow [0, 1]$ we have the following Theorem:

Theorem 6 ([HK12]) *PrivateVCG has the following properties:*

1. It is dominant strategy truthful
2. It is ϵ -differentially private
3. For any r , with probability at least $1 - \exp(-r)$ it produces an outcome o such that:

$$\sum_{i=1}^n v(t_i, o) \geq OPT(t) - \frac{2}{\epsilon} (\log |\mathcal{O}| + r)$$

Remark Again, suppose we are in a “large market” setting in which $\text{OPT}(t)$ is growing with n , at a rate faster than $\log |\mathcal{O}|$. Then in exchange for privacy, we have given up only a diminishing fraction of the optimal social welfare. In return, we have obtained a remarkable robustness property: Even if we have failed to account for agents having utilities for future events (that we have not modelled) that affect their utility possibly as some function of the outcome of our mechanism, our mechanism remains ϵ -approximately dominant strategy truthful. Of course, if we like, we can also choose ϵ to be some function that tends to zero as n grows.

Bibliographic Information The fact that the exponential mechanism distribution is the maximizer of the expected welfare plus entropy has been rediscovered in a number of fields, including physics (in which the exponential mechanism distribution is known as the Gibbs Measure), and online learning, in which it is known as the distribution arising from the “multiplicative weights” learning algorithm [AHK12]. This connection was first made to the exponential mechanism by Huang and Kannan [HK12], who used it to show that the exponential mechanism could be made exactly truthful with the introduction of payments. A note of encouragement – this result began as a course project of Zhiyi Huang’s in a differential privacy class here at Penn in 2011. For those taking the class, your goal should be to produce course projects that are taught in the next iteration of the course!

References

- [AHK12] Sanjeev Arora, Elad Hazan, and Satyen Kale. The multiplicative weights update method: A meta-algorithm and applications. *Theory of Computing*, 8(1):121–164, 2012.
- [FKOS13] Mark Flood, Jonathan Katz, Stephen Ong, and Adam Smith. Cryptography and the economics of supervisory information: balancing transparency and confidentiality. Technical report, 2013.
- [HK12] Zhiyi Huang and Sampath Kannan. The exponential mechanism for social welfare: Private, truthful, and nearly optimal. In *IEEE 53rd Annual Symposium on the Foundations of Computer Science (FOCS)*, pages 140–149. IEEE, 2012.
- [MT07] Frank McSherry and Kunal Talwar. Mechanism design via differential privacy. In *FOCS*, pages 94–103, 2007.