Parham Gohari, Bo Wu, Calvin Hawkins, Matthew Hale\*, Ufuk Topcu

Abstract-As members of network systems share more information among agents and with network providers, sensitive data leakage raises privacy concerns. Motivated by such concerns, we introduce a novel mechanism that privatizes vectors belonging to the unit simplex. Such vectors can be found in many applications, such as privatizing a decision-making policy in a Markov decision process. We use differential privacy as the underlying mathematical framework for this work. The introduced mechanism is a probabilistic mapping that maps a vector within the unit simplex to the same domain using a Dirichlet distribution. We find the mechanism well-suited for inputs within the unit simplex because it always returns a privatized output that is also in the unit simplex. Therefore, no further projection back onto the unit simplex is required. We verify and quantify the privacy guarantees of the mechanism for three cases: identity queries, average queries, and general linear queries. We establish a trade-off between the level of privacy and the accuracy of the mechanism output, and we introduce a parameter to balance the trade-off between them. Numerical results illustrate the proposed mechanism.

#### I. Introduction

In many decision-making problems, a policy-maker forms a control policy based on data collected from individuals in a network. The gathered data often contains sensitive information, which raises privacy concerns, e.g., for smart appliances [1]. In some applications, privatizing sensitive data has been achieved by adding carefully calibrated noise to sensitive data and functions thereof [2], [3], [4]. These noise-additive approaches are well-suited to some classes of numerical data, though sensitive data may take a form ill-suited to them. For example, developments in [5] explored symbolic control systems in which additive noise cannot be meaningfully implemented.

In this work, we privatize data that belongs to the unit simplex, *i.e.*, the set of vectors with non-negative entries that sum to one. We are primarily motivated by two uses of simplex-valued data, both from Markov decision processes (MDPs): (i) privatizing the decision policy and (ii) privatizing the transition probabilities among states (technical details for each can be found in Section II-A). Decision policies are computed to take

Parham Gohari is with the Department of Electrical and Computer Engineering, University of Texas at Austin, Austin, TX. Bo Wu and Ufuk Topcu are with the Department of Aerospace Engineering and Engineering Mechanics, and the Oden Institute for Computational Engineering and Sciences, University of Texas at Austin, Austin, TX. Email: {pgohari, bwu3, utopcu}@utexas.edu.PG, BW, and UT were supported by grants no. AFRL FA9550-19-1-0169, DARPA D19AP00004, and ARL ACC-APG-RTP W911NF.

Calvin Hawkins and Matthew Hale are with the Department of Mechanical and Aerospace Engineering at the University of Florida, Gainesville, FL. Email: {calvin.hawkins,matthewhale}@ufl.edu. CH and MH were supported by the AFOSR Center of Excellence on Assured Autonomy in Contested Environments (grant no. FA9550-19-1-0169) and by NSF CAREER Grant #1943275.

\*Corresponding author.

actions that maximize a reward in the MDP [6], [7]. In many cases, the optimal policy is a randomized function that maps each of an MDP's states to a probability distribution on the set of actions available at that state, see, *e.g.*, [8], [9], [10]. Finite action sets give rise to discrete, finitely supported probability distributions, which can be represented as vectors with nonnegative entries summing to one, which are elements of the unit simplex. Policies of this kind arise in applications such as autonomous driving [11] and the smart power grid [12], and revealing them can therefore reveal individuals' behaviors. Thus, there is a need to privatize such policies, and this is one application of privacy in the unit simplex.

The other motivating application in this work is privatizing the transition probabilities among an MDP's states. In each state, transitions to other states are given by a probability distribution. Finite numbers of states again give rise to discrete, finitely supported probability distributions, which are typically represented as elements of the unit simplex. We are interested in providing privacy when data is used to inform these facets of an MDP model. This data can be from individuals' healthcare [13], travel patterns [11], and other sensitive information. It is also vulnerable to privacy attacks [14], and we are motivated to privatize it as well.

In this paper, we use differential privacy as the underlying mathematical framework for privacy. Differential privacy, first introduced in [15], is designed to protect the exact values of sensitive pieces of data, while preserving their usefulness in statistical analyses. Two desirable properties of differential privacy are (i) that it is immune to post-processing [16], in the sense that arbitrary post-hoc transformations of privatized data do not weaken its privacy guarantees, and (ii) that it is robust to side information, in that gaining additional information about data-producing entities does not weaken its privacy guarantees by much [17]. As a result, differential privacy has been frequently used as the mathematical formulation of privacy in both computer science and, more recently, in control theory [18], [19], [20], [21], [22]. Existing noiseadditive approaches will not, in general, produce a privatized vector in the unit simplex. Projecting these privatized vectors back onto the simplex leads to poor accuracy of privatized data (which we illustrate in Section II-D). We therefore propose a new approach to privatization for this context.

As the main contribution of this paper, we introduce a novel mechanism that privatizes a vector within the unit simplex. A mechanism is a probabilistic mapping from some predefined domain to a pre-defined range, and a mechanism is used to privatize sensitive data. This paper develops a novel mechanism using the Dirichlet distribution, and we therefore call it the Dirichlet mechanism. The Dirichlet distribution is a multivariate distribution supported on the unit simplex, which makes it a natural choice for this setting because its outputs

are always elements of the unit simplex.

In our developments, we use probabilistic differential privacy, which is known to imply that the conventional form of differential privacy also holds [23]. Then, we show that the Dirichlet mechanism satisfies probabilistic differential privacy for identity queries. By an identity query, we mean privatizing a single vector within the unit simplex. In the course of proving these privacy guarantees, based on the assumptions we provide, we prove the log-concavity of the cumulative distribution function of a Dirichlet distribution. The proof that we present may be of independent interest in ongoing research on convexity analysis of special functions such as [24].

Beyond identity queries, we further show that the Dirichlet mechanism is differentially private for both average queries and general linear queries, in which we privatize operations over collections of vectors, each of which is contained in the unit simplex. We derive analytic expressions for privacy levels of both cases.

We also analyze the accuracy of the output of the Dirichlet mechanism. In particular, we evaluate the accuracy of the Dirichlet mechanism in terms of the expected value and the variance of its outputs. Similar to additive noise methods, the Dirichlet mechanism output has the same expected value as its input, which implies that its privatized outputs obey a distribution centered on the underlying sensitive data. We show that there exists a trade-off between privacy levels and the extent to which privatized data is concentrated around the underlying sensitive data.

We emphasize that additive noise privacy mechanisms are ill-suited to privacy on the unit simplex. The standard Laplace and Gaussian mechanisms add noise of infinite support [16], and these mechanisms will output vectors that do not belong to the unit simplex. Projecting back onto the simplex leads to poor accuracy, as we show in Section II. Recent work has rigorously established that finite-support Laplacian noise can be used for scalar-valued queries [25]. However that distribution does not have a closed form in general, which makes accuracy guarantees difficult to provide. An extension to vector-valued queries with dependent coordinates (such as summing to one for the simplex) appears quite difficult. It is for these reasons that we develop the Dirichlet mechanism.

Although its form appears quite different from existing mechanisms, they are related through membership in a broad class of probability distributions. In particular, the Laplacian, Gaussian, and exponential mechanisms all use distributions belonging to a parameterized family of exponential distributions. The outputs of the Dirichlet distribution can be generated using exponential distributions, which means the Dirichlet mechanism also belongs to the same family. This connection reveals why we should expect the Dirichlet mechanism to be well-suited to differential privacy, and this paper formalizes and confirms this intuition.

We also point out here that the exponential mechanism is another widely used differential privacy mechanism which can be used for sensitive data ill-suited to additive approaches [16]. However, the exponential mechanism can be computationally demanding to implement for privacy applications with many possible outputs. The output space here is the unit simplex,

which contains uncountably many elements. The resulting complexity of such an implementation therefore makes it infeasible [26], especially in large dimensions, and we avoid it here

A preliminary version of this work appeared in [27]. The current paper develops additional privatization techniques for general linear queries, provides concentration bounds to assess the accuracy of the Dirichlet mechanism, and provides full proofs of all results.

The rest of the paper is organized as follows. Section II establishes the privacy preliminaries needed in the rest of the paper. Then, Section III establishes privacy guarantees for identity queries, Section IV establishes privacy guarantees for averaging queries, and Section V establishes privacy guarantees for general linear queries. Section VI provides accuracy bounds on the Dirichlet mechanism's outputs, and Section VII provides simulations to illustrate our results. Finally, Section VIII concludes the paper.

## II. MOTIVATION AND PRELIMINARIES

We begin by briefly providing technical details associated with privacy concerns for simplex data. Then we establish the mathematical preliminaries needed for our developments. Below, we represent the real numbers by  $\mathbb{R}$  and the positive reals by  $\mathbb{R}_+$ . As described in the introduction, we consider privacy over the unit simplex. We denote the unit simplex in  $\mathbb{R}^n$  by  $\Delta_n$  where

$$\Delta_n := \left\{ x \in \mathbb{R}^n \mid \sum_{i=1}^n x_i = 1, x_i \ge 0 \text{ for all } i \in [n] \right\}.$$

## A. Technical Motivation: Sequential Decision-Making

Privacy concerns for simplex data arise, for example, in sequential decision-making problems. In particular, we consider decision-making problems that model their environment as a Markov decision process (MDP). An MDP  $\mathcal{M}=(\mathcal{S},\mathcal{A},\mathcal{P},r)$  models an environment with state space  $\mathcal{S}$ , action space  $\mathcal{A}$ , transition probabilities  $\mathcal{P}$ , and reward function  $r:\mathcal{S}\times\mathcal{A}\to\mathbb{R}$ . The goal with an MDP is to find a reward-maximizing policy  $\pi:S\to A$  that specifies the probability of taking each action in a particular state.

In a given state s, the probability of transitioning to a new state s' when taking action a is P(s,a,s'). Given the pair (s,a), we use the vector P(s,a) to denote the vector of all transition probabilities to other states when taking action a in state s. This vector is in the simplex: by virtue of being a finitely supported, discrete probability distribution, its entries are non-negative and they sum to one. Then  $P(s,a) \in \Delta_n$  and  $P \subseteq \Delta_n$ . As for the decision policy  $\pi$ , in a state s, the probability of taking action a is given by  $\pi(s,a)$ . Using  $\pi(s)$  to denote the vector of all such probabilities, it too is in the unit simplex: it is a discrete probability distribution on a finite set, and thus its entries are non-negative and sum to one. Then for all MDPs we find  $\pi(s) \in \Delta_n$ .

Both  $\mathcal{P}$  and  $\pi$  are sensitive. The transition probabilities in  $\mathcal{P}$  can reveal the internal dynamics of an MDP or the knowledge that drove modeling decisions for its environment. In fact, [14]

defines an attack for inferring  $\mathcal{P}$  in an MDP appearing in a reinforcement learning context. In addition, a decision policy can reveal the intentions of an agent modeled as an MDP, and [28] outlines this concern. These vulnerabilities and the wide use of MDPs motivate our developments for privacy over the unit simplex.

## B. Notation

For a positive integer n, let  $[n]:=\{1,\ldots,n\}$ . As above, we use  $\Delta_n$  to denote the unit simplex in  $\mathbb{R}^n$ , and we use  $\Delta_n^{\circ}$  to represent the interior of  $\Delta_n$ . Letting  $W\subseteq [n-1]$  with  $|W|\geq 2$  we then define the set

$$\Delta_{n,W}^{(\eta,\bar{\eta})} := \bigg\{ p \in \Delta_n^{\circ} \mid \sum_{i \in W} p_i \le 1 - \bar{\eta}, p_i \ge \eta \text{ for all } i \in W \bigg\}.$$

We impose the following assumption on  $\eta$  and  $\bar{\eta}$  that will be used below to ensure that ratios of Dirichlet distributions remain bounded when showing that they provide differential privacy.

**Assumption 1.** In 
$$\Delta_{n,W}^{(\eta,\bar{\eta})}$$
,  $\eta > 0$ ,  $\bar{\eta} > 0$ , and  $\eta + \bar{\eta} < \frac{1}{2}$ .

Letting p be a vector in  $\mathbb{R}^n$ , we use the notation  $p_{(i,j)}$  to denote the vector  $(p_i,p_j)^T \in \mathbb{R}^2$ , where  $(\cdot)^T$  is the transpose of a vector, and  $p_{-(i,j)} \in \mathbb{R}^{n-2}$  to denote the vector p with  $i^{th}$  and  $j^{th}$  entries removed.  $\mathbb{P}[\cdot]$  denotes the probability of an event. For a random variable,  $\mathbb{E}[\cdot]$  denotes its expectation and  $\mathrm{Var}[\cdot]$  denotes its variance. We use the notation  $|\cdot|$  for the cardinality of a finite set.  $||\cdot||_1$  denotes the 1-norm of a vector. We also use the special functions

$$\Gamma(x) = \int_0^\infty z^{x-1} \exp(-z) dz, \qquad x \in \mathbb{R}_+$$

$$\mathrm{beta}(a,b) = \int_0^1 t^{a-1} (1-t)^{b-1} dt = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}, \quad \ a,b \in \mathbb{R}_+$$

$$\psi^{(0)}(x) = \frac{d}{dx} \log \left(\Gamma(x)\right), \ \psi^{(1)}(x) = \frac{d^2}{dx^2} \log \left(\Gamma(x)\right), \ x \in \mathbb{R}_+$$

which are the gamma, beta, digamma, and trigamma functions, respectively.

# C. Differential Privacy

Intuitively, differential privacy guarantees that two *nearby* pieces of sensitive data will have statistically similar privatized values. In differential privacy, the notion of "nearby" is formally defined by an adjacency relation, and we define adjacency over the unit simplex as follows.

**Definition 1.** For a constant  $b \in (0,1]$  and fixed set  $W \subseteq [n-1]$ , two vectors  $p,q \in \Delta_{n,W}^{(\eta,\bar{\eta})}$  are said to be b-adjacent if there exist indices  $i,j \in W$  such that

$$p_{-(i,j)} = q_{-(i,j)}$$
 and  $||p - q||_1 \le b$ .

We express this condition with the binary symmetric adjacency relation

$$\mathrm{Adj}_b(p,q) = \begin{cases} 1 & p \text{ and } q \text{ are adjacent} \\ 0 & \text{otherwise} \end{cases}.$$

In words, two vectors are adjacent if they differ in two entries by an amount not more than b. Conventional differential privacy considers sensitive data differing in a single entry, e.g., one entry in a database [16]. However, it is not possible to do so for an element of the unit simplex because changing only a single entry would violate the condition that vectors' entries sum to one. We therefore consider privacy with the above adjacency relation. Differential privacy itself is defined next.

**Definition 2.** (Probabilistic differential privacy; [29]) Let  $b \in (0,1]$  and  $W \subseteq [n-1]$  be given. Fix a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ . A mechanism  $\mathcal{M}: \Delta_{n,W}^{(\eta,\bar{\eta})} \times \Omega \to \Delta_n$  is said to be probabilistically  $(\epsilon, \delta)$ -differentially private if we can partition the output space  $\Delta_n$  into two disjoint sets  $\Omega_1, \Omega_2$ , such that, for all  $p \in \Delta_{n,W}^{(\eta,\bar{\eta})}$ ,

$$\mathbb{P}[\mathcal{M}(p) \in \Omega_2] \le \delta_2$$

and for all  $q \in \Delta_{n.W}^{(\eta,\bar{\eta})}$  b-adjacent to p and for all  $x \in \Omega_1$ ,

$$\log\left(\frac{\mathbb{P}[\mathcal{M}(p) = x]}{\mathbb{P}[\mathcal{M}(q) = x]}\right) \le \epsilon.$$

We note that  $(\epsilon, \delta)$ -probabilistic differential privacy is known to imply conventional  $(\epsilon, \delta)$ -differential privacy [29], which requires that

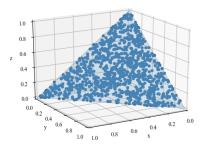
$$\mathbb{P}[\mathcal{M}(p) \in S] \le e^{\epsilon} \mathbb{P}[\mathcal{M}(q) \in S] + \delta$$

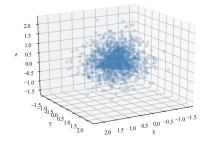
for all measurable subsets S of the range of  $\mathcal{M}$  and all adjacent p and q. Probabilistic differential privacy has a complicated interpretation under post-processing (and may not hold for the same  $\epsilon$  and  $\delta$  depending on the post-processing steps undertaken), while conventional differential privacy is immune to post-processing. We are interested in providing privacy guarantees that hold regardless of what analysis is undertaken on private data, and thus we are primarily interested in conventional differential privacy. Proving that probabilistic differential privacy holds is one way to show that conventional differential privacy holds, and we are able to use this fact in this work.

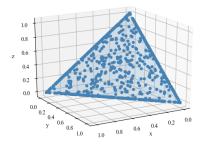
Specifically, in Sections III-V we prove that the Dirichlet mechanism satisfies Definition 2 as a means of showing that it provides conventional differential privacy. Regarding the relative strength between the two forms of privacy, it has been noted in the literature [30, Section 4] that providing  $(\epsilon, \delta)$ -probabilistic differential privacy implies  $(\epsilon, \delta')$ -differential privacy is also provided for some  $\delta' < \delta$ . Therefore, the differential privacy guarantees provided by the Dirichlet mechanism are strictly stronger than what is implied by the probabilistic differential privacy parameters we derive. The parameters we derive lie in conventionally desirable ranges for privacy, and hence we provide ordinary differential privacy of conventionally desirable strength as well.

## D. On Additive Noise Approaches

It was noted in the Introduction that differential privacy (in both its probabilistic and ordinary forms) can be enforced with additive Gaussian or Laplacian noise. For simplex-valued data, adding noise with infinite support can perturb data outside







- (a) Vectors in the unit simplex
- (b) Vectors after adding Gaussian noise

(c) Projections of noisy vectors

Fig. 1: Left-hand figure: a random sampling of 1,000 vectors from the unit simplex in  $\mathbb{R}^3$ . Center figure: the privatized forms of all data points after adding Gaussian noise from the distribution  $\mathcal{N}(0,0.25I)$ . Right-hand figure: the results of projecting all privatized vectors back onto the unit simplex. Adding Gaussian noise and projecting back onto the simplex results in many points accumulating at the boundary of the simplex, which harms accuracy.

the simplex, and some form of projection would be required to ensure membership in the unit simplex after privatization. However, we show in Figure 1 that adding Gaussian noise to elements of the simplex and projecting back onto the simplex results in biases in private data. Specifically, Figure 1 shows elements of the simplex to which Gaussian noise has been added, after which the projection algorithm in [31] is used to project them onto the simplex. The variance of noise added is  $\sigma^2=0.25$ . Using [32, Theorem 3], this level of noise provides (1,0.01)-differential privacy (for an adjacency relation that classifies vectors as adjacent if their 2-norm distance is bounded above by 0.1).

It can be seen that the privatized-then-projected forms of many elements of the simplex are simply mapped to its boundary. The reason is that adding infinite-support noise to a point in the simplex will almost always move it out of the simplex, and the projection step maps such points to the boundary of the simplex. We note here that the projection back onto the simplex is necessary for the data to be valid. For example, for Markov decision processes (discussed in Section II-A), the data of interest is vectors of transition probabilities and decision policies, and both must have their privatized forms contained in the simplex, which requires projecting onto it. This both impairs the accuracy of the mechanism itself, because the privatized forms of data points can be far from their original, sensitive values, and harms any downstream uses of this data. Rather than taking this approach, accuracy could be improved by effectively utilizing the interior of the simplex to enforce the approximate indistinguishability required by differential privacy, and that is the subject of the next subsection.

# E. Dirichlet Mechanism

One contribution of this paper is to present a differentially private mechanism that, without any need of projection, maps elements of  $\Delta_n$  to  $\Delta_n$ . In order to do so, we first introduce the Dirichlet mechanism. A Dirichlet mechanism with parameter  $k \in \mathbb{R}_+$ , denoted by  $\mathcal{M}_D^{(k)}$ , takes as input a vector  $p \in \Delta_n^o$ 

and outputs  $x \in \Delta_n$  according to the Dirichlet probability distribution function (PDF) centered on p, *i.e.*,

$$\mathbb{P}[\mathcal{M}_D^{(k)}(p) = x] = \frac{1}{B(kp)} \prod_{i=1}^{n-1} x_i^{kp_i - 1} \left( 1 - \sum_{i=1}^{n-1} x_i \right)^{kp_n - 1},$$
(1)

where

$$B(kp) := \frac{\prod_{i=1}^{n} \Gamma(kp_i)}{\Gamma\left(k\sum_{i=1}^{n} p_i\right)}$$
 (2)

is the multi-variate beta function. For brevity, we will use the notation  $Dir_k$  to denote the PDF on the right-hand side of (1). We impose the following assumption on the parameter k.

**Assumption 2.** For the Dirichlet mechanism  $\mathcal{M}_D^{(k)}$ , the parameter k satisfies

$$k \ge \max\left\{\frac{1}{\eta}, \frac{1}{1-\eta-\bar{\eta}}\right\}.$$

We later use the parameter k to adjust the trade-off that we establish between the accuracy and the privacy level of the Dirichlet mechanism. Next, we establish the privacy guarantees that the Dirichlet mechanism provides.

# III. DIRICHLET MECHANISM FOR DIFFERENTIAL PRIVACY OF IDENTITY QUERIES

We begin by analyzing identity queries under the Dirichlet mechanism. Here, a sensitive vector p is directly input to the Dirichlet mechanism to make it approximately indistinguishable from other adjacent sensitive vectors. The space of sensitive data of interest is  $\Delta_{n,W}^{(\eta,\bar{\eta})}$ , and it is over this space that we provide privacy. To show the level of privacy that holds, we first bound  $\delta$ , then bound  $\epsilon$ .

# A. Computing $\delta$

Fix  $W \subseteq [n-1]$ . In accordance with Definition 2, we partition the output space of the Dirichlet mechanism into two sets  $\Omega_1, \Omega_2$  defined by

$$\Omega_1 := \{ x \in \Delta_n \mid x_i \ge \gamma \text{ for all } i \in W \}$$
 (3)

and  $\Omega_2 := \Delta_n \backslash \Omega_1$ , where  $\gamma \in (0,1)$  is a parameter that defines these sets, upon which we impose the following.

**Assumption 3.** Fix 
$$W \subseteq [n-1]$$
. Then  $\gamma \leq \frac{1}{|W|}$ .

Next, our goal is to show that the Dirichlet mechanism output belongs to  $\Omega_1$  with high probability. Let p be a vector in  $\Delta_{n,W}^{(\eta,\bar{\eta})}$ . In the next lemma we show how to calculate  $\mathbb{P}[\mathcal{M}_D^{(k)}(p) \in \Omega_1]$ . As in Definition 2, we will bound  $\mathbb{P}[\mathcal{M}_D^{(k)}(p) \in \Omega_2]$  by  $\delta$  and use  $\epsilon$  to bound the ratios of distributions of outputs in  $\Omega_1$ . Changing  $\gamma$  changes  $\Omega_1$  and  $\Omega_2$ , and its role in determining  $\epsilon$  and  $\delta$  will be elaborated upon below.

**Lemma 1.** Let Assumptions 1 and 3 hold. Let  $W \subseteq [n-1]$ , let  $p \in \Delta_{n,W}^{(\eta,\bar{\eta})}$ , and let

$$\mathcal{A}_r := \left\{ x \in \mathbb{R}^{r-1} \mid \sum_{i \in [r-1]} x_i \le 1, x_i \ge \gamma \text{ for all } i \in W \right\},\,$$

for all  $r \geq |W| + 1$ . Then, for a Dirichlet mechanism with parameter  $k \in \mathbb{R}_+$ , we have that  $\mathbb{P}[\mathcal{M}_D^{(k)}(p) \in \Omega_1]$  equals

$$\frac{\int_{\mathcal{A}_{|W|+1}} \prod_{i \in W} x_i^{kp_i - 1} \left(1 - \sum_{i \in W} x_i\right)^{k(1 - \sum_{i \in W} p_i) - 1} \prod_{i \in W} dx_i}{\mathbf{B}(k\tilde{p}_W)},$$

where  $\tilde{p}_W \in \Delta_{|W|+1}$  is equal to p after removing entries with indices outside W and with an additional entry equal to  $1 - \sum_{i \in W} p_i$  appended as its new final entry.

*Proof.* For concreteness we set W=[n-1], though the proof is identical for other cases. In order to find  $\mathbb{P}[\mathcal{M}_D^{(k)}(p) \in \Omega_1]$ , we need to integrate the Dirichlet PDF over the region  $\mathcal{A}_n$ . Therefore, we need to evaluate the (n-1)-fold integral

$$\frac{\int_{\mathcal{A}_n} \left( \prod_{i=1}^{n-1} x_i^{kp_i - 1} \right) \left( 1 - \sum_{i=1}^{n-1} x_i \right)^{kp_n - 1} dx_{n-1} \dots dx_1}{\mathbf{B}(kp)}.$$
 (4)

Using a method similar to the one adopted in [33], let  $y := \sum_{i=1}^{n-2} x_i$ . Then we can rewrite (4) as

$$\frac{1}{\mathbf{B}(kp)} \int_{\mathcal{A}_{n-1}} \int_0^{1-y} \left( \prod_{i=1}^{n-1} x_i^{kp_i - 1} \right) (1 - y - x_{n-1})^{kp_n - 1} dx_{n-1} \dots dx_1.$$
 (5)

Now let  $u:=\frac{x_{n-1}}{1-y}$  and take the inner integral with respect to u. Then (5) becomes

$$\frac{1}{B(kp)} \int_{\mathcal{A}_{n-1}} \prod_{i=1}^{n-2} x_i^{kp_i-1} (1-y)^{k(p_{n-1}+p_n)-1} \int_0^1 u^{kp_{n-1}-1} (1-u)^{kp_n-1} du \ dx_{n-2} \dots dx_1.$$

From the definition of the beta function, we have

$$\int_{0}^{1} u^{kp_{n-1}-1} (1-u)^{kp_{n}-1} du = \text{beta}(kp_{n-1}, kp_{n}).$$

Using the gamma function representation of beta functions, *i.e.*.

$$beta(a,b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}, \ a,b \in \mathbb{R}_+, \tag{6}$$

and (2), we find that  $\mathbb{P}[\mathcal{M}_D^{(k)}(p) \in \Omega_1]$  is equal to

$$\frac{1}{B(kp)} \frac{\Gamma(kp_{n-1})\Gamma(kp_n)}{\Gamma(k(p_{n-1}+p_n))} \int_{\mathcal{A}_{n-1}} \prod_{i=1}^{n-2} x_i^{kp_i-1} \left(1 - \sum_{i=1}^{n-2} x_i\right)^{k(p_{n-1}+p_n)-1} dx_{n-2} \dots dx_1.$$

Using the same idea, for the next step, let  $y:=\sum_{i=1}^{n-3}x_i$  and  $u:=\frac{x_{n-2}}{1-y}$ . Then  $\mathbb{P}[\mathcal{M}_D^{(k)}(p)\in\Omega_1]$  is equal to

$$\frac{1}{\mathbf{B}(kp)} \frac{\Gamma(kp_{n-2})\Gamma(kp_{n-1})\Gamma(kp_n)}{\Gamma(k(p_{n-2}+p_{n-1}+p_n))} \int_{\mathcal{A}_{n-2}} \prod_{i=1}^{n-3} x_i^{kp_i-1} \left(1 - \sum_{i=1}^{n-3} x_i\right)^{k(p_{n-2}+p_{n-1}+p_n)-1} dx_{n-3} \dots dx_1.$$

We continue to adopt the same change of variable strategy until we are left with an integral over the region  $A_{|W|+1}$ , which concludes the proof.

Lemma 1 shows that instead of an (n-1)-fold integral of the Dirichlet PDF, the computations can be reduced to a |W|-fold integral. However, the expression still depends on the input vector p, which is undesirable and generally incompatible with differential privacy. The reason is that  $(\epsilon, \delta)$ -differential privacy must be a guarantee for all adjacent input data and not for a specific data point. In the next lemma, we show that  $\mathbb{P}[\mathcal{M}_D^{(k)}(p) \in \Omega_1]$  is a log-concave function of p over  $\Delta_{n,W}^{(\eta,\bar{\eta})}$ , which we will use to derive a bound for  $\delta$  that holds for all p of interest.

**Lemma 2.** Let Assumption 1 hold, fix  $W\subseteq [n-1]$ , and let  $\mathcal{M}_D^{(k)}$  be the Dirichlet mechanism with parameter k. Then  $\mathbb{P}[\mathcal{M}_D^{(k)}(p)\in\Omega_1]$  is a log-concave function of p over the domain  $\Delta_{n,W}^{(\eta,\bar{\eta})}$ .

*Proof:* See Appendix A. Revisiting the definitions of  $\Omega_1, \Omega_2$  above, we find that

$$\mathbb{P}[\mathcal{M}_D^{(k)}(p) \in \Omega_2] = 1 - \mathbb{P}[\mathcal{M}_D^{(k)}(p) \in \Omega_1]$$

$$\leq 1 - \min_{p} \mathbb{P}[\mathcal{M}_D^{(k)}(p) \in \Omega_1] = \delta$$
(7)

is the smallest possible choice of  $\delta$ , and we use this value for the remainder of the paper. From this, we see that bounding  $\delta$  can be done by minimizing  $\mathbb{P}[\mathcal{M}_D^{(k)}(p) \in \Omega_1]$ , an explicit form of which was given in Lemma 1. In Lemma 2, we established the log-concavity of the function that we seek to minimize. As a result, instead of minimizing  $\mathbb{P}[\mathcal{M}_D^{(k)}(p) \in \Omega_1]$  over the entirety of  $\Delta_{n,W}^{(\eta,\bar{\eta})}$ , we can only consider the extreme points. Note that the points within  $\Delta_{n,W}^{(\eta,\bar{\eta})}$  form a polyhedron with at most |W|(|W|+1)/2 vertices. As the minimum of an unsorted list of n entries can be found in linear time, the time

complexity of finding  $\min \mathbb{P}[\mathcal{M}_D^{(k)}(p) \in \Omega_1]$  is  $\mathcal{O}(|W|^2)$ . This analytical bound will be further explored through numerical results in Section VII. Next, we develop analogous bounds for  $\epsilon$ .

## B. Computing $\epsilon$

As above, fix  $\eta, \bar{\eta} \in (0,1)$  satisfying Assumption 1,  $b \in (0,1]$ , and  $W \subseteq [n-1]$ . Then, for a given  $k \in \mathbb{R}_+$ , bounding  $\epsilon$  requires evaluating the term

$$\log \left( \frac{\mathbb{P}[\mathcal{M}_D^{(k)}(p) = x]}{\mathbb{P}[\mathcal{M}_D^{(k)}(q) = x]} \right)$$

for all  $x \in \Omega_1$ , where p and q are any b-adjacent vectors in  $\Delta_{n,W}^{(\eta,\bar{\eta})}$ . Let  $i,j \in W$  be the indices in which p and q differ. Using the definition of the Dirichlet mechanism, we find

$$\log \left( \frac{\mathbb{P}[\mathcal{M}_{D}^{(k)}(p) = x]}{\mathbb{P}[\mathcal{M}_{D}^{(k)}(q) = x]} \right) = \log \left( \frac{\mathbf{B}(kq) \prod_{i=1}^{n} x_{i}^{kp_{i}-1}}{\mathbf{B}(kp) \prod_{i=1}^{n} x_{i}^{kq_{i}-1}} \right)$$

$$= \log \left( \frac{\Gamma(kq_{i})\Gamma(kq_{j})x_{i}^{kp_{i}-1}x_{j}^{kp_{j}-1}}{\Gamma(kp_{i})\Gamma(kp_{j})x_{i}^{kq_{i}-1}x_{j}^{kq_{j}-1}} \right)$$

$$= \log \left( \frac{\Gamma(kq_{i})\Gamma(kq_{j})}{\Gamma(kp_{i})\Gamma(kp_{j})}x_{i}^{k(p_{i}-q_{i})}x_{j}^{k(p_{j}-q_{j})} \right).$$

Since p and q are b-adjacent, we have that  $p_i + p_j = q_i + q_j$ . Therefore, we can compute  $\epsilon$  by evaluating the term

$$\log \left( \frac{\Gamma(kq_i)\Gamma(kq_j)}{\Gamma(kp_i)\Gamma(kp_j)} \left( \frac{x_i}{x_j} \right)^{k(p_i - q_i)} \right). \tag{8}$$

Note that if either  $x_i$  or  $x_j$  goes to 0, then the term in (8) would be unbounded. Recalling that the indices at which p and q can differ are restricted to the set W, we find that the values at these indices must be bounded below by  $\eta$ , and therefore the ratios of interest remain bounded as well.

Lemma 4 below will provide an explicit value of  $\epsilon$ , aided in part by the following lemma.

**Lemma 3.** Let Assumptions 1 and 2 hold. Let W be a given set of indices which is used to construct  $\Delta_{n,W}^{(\eta,\bar{\eta})}$  and let p,q be any b-adjacent vectors in  $\Delta_{n,W}^{(\eta,\bar{\eta})}$  with their  $i^{th}$  and  $j^{th}$  entries different. Then, for a constant  $k \in \mathbb{R}_+$ , we have that

$$\frac{\mathrm{beta}(kq_i,kq_j)}{\mathrm{beta}(kp_i,kp_j)} \leq \frac{\mathrm{beta}(kq_i,k(1-\bar{\eta}-q_i))}{\mathrm{beta}(kp_i,k(1-\bar{\eta}-p_i))}$$

Proof: See Appendix B.

**Lemma 4.** Let Assumptions 1, 2, and 3 hold, let  $\Omega_1$  be as defined in (3), let  $W \subseteq [n-1]$ , and let  $\mathcal{M}_D^{(k)}$  be a Dirichlet mechanism with parameter k. Then, for all adjacent p and q and for all  $x \in \Omega_1$  we have that

$$\begin{split} \log \left( \frac{\mathbb{P}[\mathcal{M}_D^{(k)}(p) = x]}{\mathbb{P}[\mathcal{M}_D^{(k)}(q) = x]} \right) \leq \\ \log \left( \frac{\operatorname{beta}(k\eta, k(1 - \bar{\eta} - \eta))}{\operatorname{beta}(k(\eta + \frac{b}{2}), k(1 - \bar{\eta} - \eta - \frac{b}{2}))} \right) \\ + \frac{kb}{2} \log \left( \frac{1 - (|W| - 1))\gamma}{\gamma} \right), \end{split}$$

where the parameter  $\gamma \in (0,1)$  defines the set  $\Omega_1$  as in Section III-A.

*Proof.* Because p and q are adjacent, we suppose they differ in indices  $i, j \in W$ . Then from (8) we know that

$$\log\!\left(\frac{\mathbb{P}[\mathcal{M}_D^{(k)}(p)\!=\!x]}{\mathbb{P}[\mathcal{M}_D^{(k)}(q)\!=\!x]}\right)\!=\!\log\!\left(\!\frac{\Gamma(kq_i)\Gamma(kq_j)}{\Gamma(kp_i)\Gamma(kp_j)}\!\left(\!\frac{x_i}{x_j}\right)^{k(p_i-q_i)}\!\right)\!.$$

Let

$$v := \max_{p,q,x \in \mathbb{R}^n} \log \left( \frac{\Gamma(kq_i)\Gamma(kq_j)}{\Gamma(kp_i)\Gamma(kp_j)} \left( \frac{x_i}{x_j} \right)^{k(p_i - q_i)} \right)$$
subject to  $|p_i - q_i| \le \frac{b}{2}$ ,
$$p_i + p_j = q_i + q_j,$$

$$p_i + p_j \le 1 - \bar{\eta},$$

$$p_{(i,j)} \in [\eta, 1 - \bar{\eta} - \eta]^2,$$

$$q_{(i,j)} \in [\eta, 1 - \bar{\eta} - \eta]^2,$$

$$x_{(i,j)} \in [\gamma, 1 - (|W| - 1)\gamma]^2,$$

and let  $\mathcal C$  denote the set of feasible points of the optimization problem in (9); we note that the first two constraints enforce adjacency, while the others encode  $p,q\in\Delta_{n,W}^{(\eta,\bar\eta)}$  and  $x\in\Omega_1$ . Assumptions 1-3 ensure that all intervals above are non-empty.

By sub-additivity of the maximum, we have

$$v \leq \max_{p,q,x \in \mathcal{C}} \log \left( \frac{\Gamma(kq_i)\Gamma(kq_j)}{\Gamma(kp_i)\Gamma(kp_j)} \right) + \max_{p,q,x \in \mathcal{C}} \log \left( \frac{x_i}{x_j} \right)^{k(p_i - q_i)}. \quad (10)$$

Now, with

$$v_1 := \max_{p,q,x \in \mathcal{C}} \log \left(\frac{x_i}{x_i}\right)^{k(p_i - q_i)},$$

we find

$$v_1 \le \max_{p,q,x \in \mathcal{C}} |k(p_i - q_i)| \left| \log \left( \frac{x_i}{x_j} \right) \right|$$
  
$$\le \frac{kb}{2} \log \left( \frac{1 - (|W| - 1)\gamma}{\gamma} \right).$$

The fact that  $|p_i - q_i| \leq \frac{b}{2}$  follows from adjacency in Definition 1, and the definition of  $\Omega_1$  directly implies both that  $x_i \leq 1 - (|W| - 1)\gamma$  and that  $x_j \geq \gamma$ . Assumption 3 then ensures that the argument of the logarithm is positive.

Next, let  $c:=p_i+p_j=q_i+q_j$  and substitute  $q_j,p_j$  with  $c-q_i$  and  $c-p_i$  respectively. Let

$$\begin{split} v_2 := & \max_{p_i,q_i,c \in \mathbb{R}} & \log \left( \frac{\Gamma(kq_i)\Gamma(k(c-q_i))}{\Gamma(kp_i)\Gamma(k(c-p_i))} \right) \\ & \text{subject to} & |p_i - q_i| \leq \frac{b}{2}, \\ & c \in [2\eta, 1 - \bar{\eta}], \\ & p_i \in [\eta, 1 - \bar{\eta} - \eta], \\ & q_i \in [\eta, 1 - \bar{\eta} - \eta], \end{split}$$

where the constraints again encode adjacency of p and q and their containment in  $\Delta_{n.W}^{(\eta,\bar{\eta})}$ .

Next, either  $q_i < p_i$  or  $q_j < p_j$ , and we assume without loss of generality that  $q_i < p_i$ . Then, from Lemma 3 and (6), we have that

$$v_{2} \leq \max_{p_{i},q_{i}\in\mathbb{R}} \log \left( \frac{\operatorname{beta}(kq_{i},k(1-\bar{\eta}-q_{i}))}{\operatorname{beta}(kp_{i},k(1-\bar{\eta}-p_{i}))} \right)$$
subject to  $|p_{i}-q_{i}| \leq \frac{b}{2}$ , (11)
$$p_{i} \in [\eta,1-\bar{\eta}-\eta],$$

$$q_{i} \in [\eta,1-\bar{\eta}-\eta].$$

Evaluating the gradient of the objective function in the optimization problem in (11), it can be shown that the Karush-Kuhn-Tucker (KKT) conditions of optimality are not satisfied in the interior of the set of feasible points except for points that lie on the line  $p_i = q_i$ , which are minima. Thus, since the KKT conditions are only necessary conditions (see Chapter 11 of [34]), satisfying them does not imply optimality, and we exclude points where  $p_i = q_i$  from the set of possible maximizers.

Evaluating points on the boundary of the feasible region shows that KKT conditions are also not satisfied. Thus, we need only to consider the extreme  $(p_i, q_i)$ 's in the set

$$\left\{ \left( \eta + \frac{b}{2}, \eta \right), \left( 1 - \bar{\eta} - \eta - \frac{b}{2}, 1 - \bar{\eta} - \eta \right), \left( \eta, \eta + \frac{b}{2} \right), \left( 1 - \bar{\eta} - \eta, 1 - \bar{\eta} - \eta - \frac{b}{2} \right) \right\}, (12)$$

which are the vertices of the feasible region. Note that since beta(a,b) = beta(b,a), the points in the first row give equal positive objectives and the points in the second row have equal negative objectives. Hence, we can choose the first point in (12) to find

$$v_2 = \log \left( \frac{\operatorname{beta}(k\eta, k(1 - \bar{\eta} - \eta))}{\operatorname{beta}(k(\eta + \frac{b}{2}), k(1 - \bar{\eta} - \eta - \frac{b}{2}))} \right).$$

Substituting  $v_1$  and  $v_2$  in (10) concludes the proof.

We now state the main theorem of this section, which formally establishes the  $(\epsilon, \delta)$ -differential privacy of the Dirichlet mechanism for identity queries.

**Theorem 1.** Fix  $\eta, \bar{\eta} \in (0, 1)$ ,  $b \in (0, 1]$ , and  $W \subseteq [n - 1]$ , and let Assumptions 1-3 hold. Let the adjacency relation in Definition 1 hold. Then the Dirichlet mechanism with parameter  $k \in \mathbb{R}_+$ , defined as  $\mathcal{M}_D^{(k)}(p) = \mathrm{Dir}_k(p)$ , is  $(\epsilon, \delta)$ -differentially private, where

$$\begin{split} \epsilon &= \log \left( \frac{\mathrm{beta}(k\eta, k(1-\bar{\eta}-\eta))}{\mathrm{beta}(k(\eta+\frac{b}{2}), k(1-\bar{\eta}-\eta-\frac{b}{2}))} \right) + \\ &\qquad \frac{kb}{2} \log \left( \frac{1-(|W|-1)\gamma}{\gamma} \right), \end{split}$$

and

$$\delta = 1 - \min_{p \in \Delta_{n,W}^{(\eta,\bar{\eta})}} \mathbb{P}[\mathcal{M}_D^{(k)}(p) \in \Omega_1].$$

*Proof.* The expression for  $\epsilon$  results immediately from Lemma 4 and the expression for  $\delta$  is a direct result of (7).

The expression given for  $\epsilon$  in Theorem 1 contains a logarithm of a ratio of beta functions, which can be difficult to reason about intuitively. In the following lemma we present upper and lower bounds for beta functions in terms of simpler functions, which we will use to provide a simplified upper bound for  $\epsilon$ .

**Lemma 5.** Let  $a, b \in \mathbb{R}$ . Then

$$\exp(2-a-b) \le \text{beta}(a,b) \le \frac{a+b-1}{(2a-1)(2b-1)}$$

*Proof:* See Appendix C.

Using Lemma 5, we can provide a simplified bound on  $\epsilon$  in exchange for that bound being somewhat looser.

**Corollary 1.** Let all conditions of Theorem 1 hold. Then, for identity queries over  $\Delta_{n,W}^{(\eta,\bar{\eta})}$ , the Dirichlet mechanism is  $(\epsilon,\delta)$ -differentially private, with

$$\epsilon = 2k(1 - \bar{\eta}) - 3 + \frac{kb}{2}\log\left(\frac{1 - (|W| - 1)\gamma}{\gamma}\right)$$

and

$$\delta = 1 - \min_{p \in \Delta_{n,W}^{(\eta,\bar{\eta})}} \mathbb{P}[\mathcal{M}_D^{(k)}(p) \in \Omega_1].$$

*Proof:* The value of  $\delta$  is the same as that in Theorem 1. For  $\epsilon$  from Theorem 1, we need to upper bound the term  $\mathrm{beta}\big(k\eta,k(1-\bar{\eta}-\eta)\big)$  and lower bound the term  $\mathrm{beta}\big(k(\eta+\frac{b}{2}),k(1-\bar{\eta}-\eta-\frac{b}{2})\big)$ . We thus apply Lemma 5 to find

$$\begin{split} &\frac{\mathrm{beta}(k\eta,k(1-\bar{\eta}-\eta))}{\mathrm{beta}(k(\eta+\frac{b}{2}),k(1-\bar{\eta}-\eta-\frac{b}{2}))} \\ &\leq \frac{k-k\bar{\eta}-1}{(2k\eta-1)(2k-2k\bar{\eta}-2k\eta-1)} \frac{1}{\exp(2-k+k\bar{\eta})}. \end{split}$$

Then, taking the logarithm of both sides, we find that

$$\begin{split} &\log\left(\frac{\mathrm{beta}(k\eta,k(1-\bar{\eta}-\eta))}{\mathrm{beta}(k(\eta+\frac{b}{2}),k(1-\bar{\eta}-\eta-\frac{b}{2}))}\right) \\ &\leq \log\left(\frac{k-k\bar{\eta}-1}{(2k\eta-1)(2k-2k\bar{\eta}-2k\eta-1)}\frac{1}{\exp(2-k+k\bar{\eta})}\right) \end{split}$$

$$\leq \log(k - k\bar{\eta} - 1) - \log(2k\eta - 1) - \log(2k - 2k\bar{\eta} - 2k\eta - 1) - 2 + k - k\bar{\eta}. \tag{13}$$

Then, using the fact that  $\log(x) < x$  in (13), the first term satisfies

$$\log(k - k\bar{\eta} - 1) \le k - k\bar{\eta} - 1.$$

Next, using the fact that

$$k \ge \max\left\{\frac{1}{\eta}, \frac{1}{1-\bar{\eta}-\eta}\right\}$$

from Assumption 1, the next two logarithm terms are both non-positive and an upper bound is furnished by eliminating them. Then we find

$$\log\left(\frac{\operatorname{beta}(k\eta,k(1-\bar{\eta}-\eta))}{\operatorname{beta}(k(\eta+\frac{b}{2}),k(1-\bar{\eta}-\eta-\frac{b}{2}))}\right)\leq 2k-3.$$

Combining this result with Theorem 1 completes the proof.

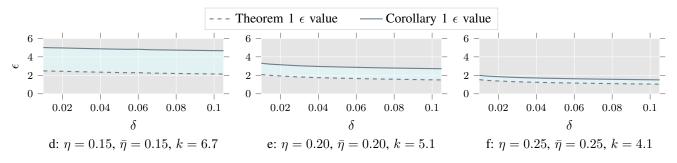


Fig. 2: An example where |W|=3 to compare the values of  $\epsilon$  using Theorem 1 with those computed using Corollary 1. For each value of  $\delta$ , first  $\gamma$  is optimized according to the optimization problem in (14), then the optimal  $\gamma$  is substituted in the expressions for values of  $\epsilon$ .

Because this bound on  $\epsilon$  is linear in k, it offers a more intuitive understanding of how changing k affects privacy.

In Figure 2, for three instances of  $(\eta, \bar{\eta}, k)$  and b = 0.1, we show how Theorem 1 and Corollary 1 capture the behavior of  $\epsilon$ . All three cases show that Corollary 1 overestimates the value of  $\epsilon$  relative to Theorem 1, and the size of overestimate decreases as  $\epsilon$  grows.

**Remark 1.** Note that if a mechanism is  $\epsilon_1$ -differentially private, it is also  $\epsilon_2$ -differentially private for all  $\epsilon_2 \geq \epsilon_1$ . Therefore, if the upper bound for  $\epsilon$  after simplification of beta functions is still within an acceptable range, e.g.,  $\delta \leq 0.05$  and  $\epsilon \leq 5$  [35], [36], [37], then using an over-approximation of  $\epsilon$  does not substantially harm our interpretation of the Dirichlet mechanism's protections.

Next, we point out that the parameter  $\gamma$ , which is used in the definition of  $\Omega_2$ , is not a parameter of the mechanism, in the sense that changing  $\gamma$  does not change the mechanism itself. Instead,  $\gamma$  balances the trade-off between privacy level and the probability of failing to guarantee that privacy level, *i.e.*, changing  $\gamma$  can decrease  $\epsilon$  in exchange for increasing  $\delta$  and vice versa.

In some cases, we are given the highest probability of privacy failure,  $\delta$ , that is acceptable, and one must maximize the level of privacy,  $\epsilon$ , subject to that upper bound. Let  $\hat{\delta}$  denote the maximum admissible value of  $\delta$ . Then we are interested in minimizing  $\epsilon$  while obeying  $\delta \leq \hat{\delta}$ . Using Theorem 1, we note that  $\epsilon$  is a strictly decreasing function of  $\gamma$ . Letting V be the set of vertices of  $\Delta_{n,W}^{(\eta,\bar{\eta})}$ , we then can minimize  $\epsilon$  by solving the problem

$$\max_{\gamma} \quad \gamma$$
 subject to  $\mathbb{P}[\mathcal{M}_{D}^{(k)}(p) \in \Omega_{1}] \geq 1 - \hat{\delta} \text{ for all } p \in V.$ 

Note that the feasible region of the optimization problem (14) is a convex set because the function  $\mathbb{P}[\mathcal{M}_D^{(k)}(p) \in \Omega_1]$  is a strictly decreasing function of  $\gamma$ . Therefore,  $\epsilon$  can be optimized for a given  $\hat{\delta}$  using off-the-shelf convex optimization toolboxes, and this will be done in Section VII. Next, we apply the Dirichlet mechanism to average queries.

# IV. DIRICHLET MECHANISM FOR DIFFERENTIAL PRIVACY OF AVERAGE QUERIES

In this section we consider a collection of N vectors indexed over  $i \in [N]$ , with the  $i^{th}$  denoted  $p^i \in \Delta_{n,W}^{(\eta,\bar{\eta})}$ . The goal is to compute the average of the collection  $\{p^i\}_{i\in[N]}$  while providing differential privacy. Accordingly, the space of sensitive data under consideration is now

$$\mathcal{S}:=\Big\{\{p^i\}_{i\in[N]}:N\in\mathbb{N}\text{ and }p^i\in\Delta_{n,W}^{(\eta,\bar{\eta})}\Big\}.$$

We next re-define the adjacency relationship for the average query setting.

**Definition 3.** Fix a scalar  $b \in (0,1]$ . Two collections in S, denoted  $\{p^i\}_{i \in [N]}$  and  $\{q^i\}_{i \in [N]}$ , are b-adjacent if there is some j such that

- 1)  $p^i = q^i$  for all  $j \neq i$ ,
- 2) there exist indices m and l such that  $p^j_{-(m,l)} = q^j_{-(m,l)}$  and  $||p^j q^j|| \le b$ .

With adjacency defined over collections, we have a corresponding definition of probabilistic differential privacy.

**Definition 4.** (Probabilistic differential privacy for collections) Let  $b \in (0,1]$  and  $W \subseteq [n-1]$  be given. Fix a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ . A mechanism  $\mathcal{M} : \left(\Delta_{n,W}^{(\eta,\bar{\eta})}\right)^N \times \Omega \to \Delta_n$  is said to be probabilistically  $(\epsilon, \delta)$ -differentially private if we can partition the output space  $\Delta_n$  into two disjoint sets  $\Omega_1, \Omega_2$ , such that, for all  $\mathcal{P} := \{p^i\}_{i \in [N]} \in \mathcal{S}$ ,

$$\mathbb{P}[\mathcal{M}(\mathcal{P}) \in \Omega_2] \le \delta,$$

and for all  $Q := \{q^i\}_{i \in [N]} \in \mathcal{S}$  b-adjacent to  $\mathcal{P}$  and for all  $x \in \Omega_1$ ,

$$\log\left(\frac{\mathbb{P}[\mathcal{M}(\mathcal{P}) = x]}{\mathbb{P}[\mathcal{M}(\mathcal{Q}) = x]}\right) \le \epsilon.$$

As with Definition 2, we use Definition 4 simply as a means to show that ordinary  $(\epsilon, \delta)$ -differential privacy holds.

The query we now consider is the average. Set  $\mathcal{P}=\{p^i\}_{i\in[N]}$  and  $\mathcal{Q}=\{q^i\}_{i\in[N]}$ . Mathematically we define the average operator  $\mathcal{A}$  via

$$\mathcal{A}(\mathcal{P}) := \frac{1}{N} \sum_{i=1}^{N} p^{i},$$

with  $\mathcal{A}(\mathcal{Q})$  defined analogously. The next theorem formalizes the privacy protections of the Dirichlet mechanism when applied to such averages.

**Theorem 2.** Fix  $\eta, \bar{\eta} \in (0, 1)$ , let  $b \in (0, 1]$ , let  $W \subseteq [n - 1]$ , and let Assumptions 1-3 hold. Let the adjacency relation in Definition 3 hold. Then the Dirichlet mechanism with parameter  $k \in \mathbb{R}_+$  defined as  $\mathcal{M}_D^{(k)}(\mathcal{P}) = \mathrm{Dir}_k(\mathcal{A}(\mathcal{P}))$  is  $(\epsilon, \delta)$ -differentially private, where

$$\epsilon = \log \left( \frac{\operatorname{beta}(k\eta, k(1 - \bar{\eta} - \eta))}{\operatorname{beta}(k(\eta + \frac{b}{2N}), k(1 - \bar{\eta} - \eta - \frac{b}{2N}))} \right) + \frac{kb}{2N} \log \left( \frac{1 - (|W| - 1)\gamma}{\gamma} \right), \quad (15)$$

and

$$\delta = 1 - \min_{p \in \Delta_{n,W}^{(\eta,\bar{\eta})}} \mathbb{P}[\mathcal{M}_D^{(k)}(\mathcal{A}(\mathcal{P})) \in \Omega_1].$$

*Proof.* We proceed by showing that Definition 4 is satisfied. Let  $x \in \Omega_1$ . Then, we are interested in the quantity

$$\frac{\mathbb{P}[\mathcal{M}_D^{(k)}(\mathcal{A}(\mathcal{P})) = x]}{\mathbb{P}[\mathcal{M}_D^{(k)}(\mathcal{A}(\mathcal{Q})) = x]} = \frac{B(k\mathcal{A}(\mathcal{Q})) \prod_{i=1}^n x_i^{kA_i(p)-1}}{B(k\mathcal{A}(\mathcal{P})) \prod_{i=1}^n x_i^{kA_i(q)-1}}.$$
 (16)

Based on the definition of the adjacency relationship for collections in Definition 3, A(p) and A(q) will differ only in their  $m^{th}$  and  $l^{th}$  entries. Taking the logarithm of both sides of (16) and using the same approach as in Theorem 1, we have that

$$\log \left( \frac{\mathbb{P}[\mathcal{M}_{D}^{(k)}(\mathcal{A}(\mathcal{P})) = x]}{\mathbb{P}[\mathcal{M}_{D}^{(k)}(\mathcal{A}(\mathcal{Q})) = x]} \right) \leq \frac{1}{\mathbb{P}[\mathcal{M}_{D}^{(k)}(\mathcal{A}(\mathcal{Q})) = x]}$$

$$= \max_{\mathcal{A}(\mathcal{P}), \mathcal{A}(\mathcal{Q})} \log \left( \frac{\mathbb{B}(k\mathcal{A}(\mathcal{Q}))}{\mathbb{B}(k\mathcal{A}(\mathcal{P}))} \right) + \frac{1}{\mathbb{E}[k]}$$

$$= \max_{\mathcal{A}(\mathcal{P}), \mathcal{A}(\mathcal{Q})} \log \left( \frac{1 - (|W| - 1)\gamma}{\gamma} \right)^{k|A_{m}(p) - A_{m}(q)|}.$$
(17)

Because  $\mathcal{P}$  and  $\mathcal{Q}$  are b-adjacent, and each entry of  $\mathcal{A}(\cdot)$  represents the average of a component, we have that

$$|\mathcal{A}_m(p) - \mathcal{A}_m(q)| \le \frac{b}{2N}.$$
(18)

Combining (17), (18) and Lemma 4 completes the proof for the value of  $\epsilon$ . For  $\delta$ , the same approach for calculating  $\delta$  in identity queries applies to average queries.

Remark 2. As seen in (15), the level of privacy increases with the number of vectors present in the collection. In particular,  $\epsilon \to 0$  as  $N \to \infty$ . This can be seen by taking the limit as  $N \to \infty$  in the expression for  $\epsilon$ . Noting that the second term is proportional to  $\frac{1}{N}$ , we observe that the first term is continuous over the positive reals and taking the limit drives the argument of the logarithm to one. This limiting behavior is consistent with the intuition that it should be harder to uncover the sensitive information of an individual in a population when their data is mixed together in an average.

As with Corollary 1, we provide simplified bounds on the value of  $\epsilon$  to ease the interpretation of each parameter's influence upon privacy.

**Corollary 2.** Let all conditions of Theorem 2 hold. Then, for average queries, the Dirichlet mechanism  $\mathcal{M}_D^{(k)}(\mathcal{P}) = \operatorname{Dir}_k(\mathcal{A}(\mathcal{P}))$  is  $(\epsilon, \delta)$ -differentially private with

$$\epsilon = 2k(1 - \bar{\eta}) - 3 + \frac{kb}{2N}\log\left(\frac{1 - (|W| - 1)\gamma}{\gamma}\right)$$

ana

$$\delta = 1 - \min_{p \in \Delta_{n,W}^{(\eta,\bar{\eta})}} \mathbb{P}[\mathcal{M}_D^{(k)}(\mathcal{A}(\mathcal{P})) \in \Omega_1].$$

*Proof:* The proof is similar to that of Corollary 1 and is therefore omitted.

# V. DIFFERENTIAL PRIVACY FOR GENERAL LINEAR OUERIES

In this section, we derive privacy guarantees for arbitrary linear queries over collections of vectors in the unit simplex. Examples of such queries are weighted averages of vectors of transition probabilities, e.g., in the smart power grid. In particular, with a variety of smart devices and smart buildings modeled as MDPs, one may wish to compute average behaviors, with the weights encoding the importance of a device or size of a building. We begin by establishing the class of queries to be considered, then we derive privacy guarantees provided by the Dirichlet mechanism for this class.

# A. General Linear Queries over the Simplex

As above, we consider privacy over the set  $\mathcal{S}$ , which contains N-element collections of vectors in the unit simplex. We again consider  $\mathcal{P}=\{p^i\}_{i\in[N]}$ , with  $p^i\in\Delta_{n,W}^{(\eta,\bar{\eta})}$ . The collection  $\mathcal{P}$  can also be represented as an  $n\times N$  matrix, where column i is equal to  $p^i$ . With an abuse of notation, we also use  $\mathcal{P}$  to denote this matrix representation, and we note that  $\mathcal{P}_{ij}=p^j_i$ . With  $\mathcal{P}\in\mathbb{R}^{n\times N}$ , we can represent the linear queries of interest by vector multiplication on the right. Namely, a linear query  $L:\mathcal{S}\to\Delta_n$  can be represented via

$$L(\mathcal{P}) = \mathcal{P}\ell = \begin{pmatrix} \sum_{j=1}^{N} p_1^j \ell_j \\ \vdots \\ \sum_{j=1}^{N} p_n^j \ell_j \end{pmatrix}, \tag{19}$$

where  $\ell \in \mathbb{R}^N$ . The following lemma establishes the stronger statement that ensuring  $L(\mathcal{P}) \in \Delta_n$  for arbitrary  $\mathcal{P} \in \mathcal{S}$  in fact requires  $\ell \in \Delta_N$ .

**Lemma 6.** Let  $L: S \to \Delta_n$  be a linear query identified with the vector  $\ell \in \mathbb{R}^N$ . Then  $\ell \in \Delta_N$ .

*Proof:* Using (19), for  $L(\mathcal{P}) \in \Delta_n$ , we require that

$$\sum_{i=1}^{n} \sum_{j=1}^{N} p_i^j \ell_j = \sum_{j=1}^{N} \sum_{i=1}^{n} p_i^j \ell_j = \sum_{j=1}^{N} \ell_j \sum_{i=1}^{n} p_i^j = \sum_{j=1}^{N} \ell_j = 1,$$

where the third equality follows from  $p^j \in \Delta_n$ . We also must have  $\ell_j \geq 0$  for all  $j \in [N]$ .

We also incorporate additional boundedness of the entries of  $\ell$  in the following definition.

**Definition 5.** Fix  $\alpha \in (0,1]$ . A vector  $\ell \in \Delta_N$  is said to be  $\alpha$ -bounded if  $\|\ell\|_{\infty} \leq \alpha$ .

With a slight abuse of terminology, we will refer to a query L as  $\alpha$ -bounded if it is identified with an  $\alpha$ -bounded vector  $\ell$ . All vectors  $\ell \in \Delta_N$  are trivially 1-bounded, though the inclusion of  $\alpha$ -boundedness allows for additional bounds on the entries of  $\ell$  to be used in our privacy analysis. This inclusion in turn enables us to make stronger statements about the Dirichlet mechanism's guarantees, which we explore further in the next section.

## B. Privacy Guarantees for Linear Queries

With this characterization of  $\ell$  in hand, we now quantify the privacy guarantees afforded to such queries by the Dirichlet mechanism. As above, determining these privacy guarantees will require bounding ratios of Dirichlet distributions, a component of which is a term involving gamma functions. The next lemma provides a bound on this term.

**Lemma 7.** Fix  $\eta, \bar{\eta} \in (0,1)$  and  $W \subseteq [n-1]$ , and let Assumptions 1-3 hold. Fix  $b \in (0,1]$  and let  $\mathcal{P}, \mathcal{Q} \in \mathcal{S}$  be adjacent according to Definition 3. Let  $L: \mathcal{S} \to \Delta_n$  be an  $\alpha$ -bounded linear query over  $\mathcal{S}$  identified with  $\ell \in \Delta_n$ . Then

$$\frac{\mathrm{B}\!\left(k\mathcal{Q}\ell\right)}{\mathrm{B}\!\left(k\mathcal{P}\ell\right)} \leq \frac{\mathrm{beta}\!\left(k\eta,k(1-\bar{\eta}-\eta)\right.}{\mathrm{beta}\!\left(k(\eta+\frac{b\alpha}{2}),k(1-\bar{\eta}-\eta-\frac{b\alpha}{2})\right)}.$$

*Proof:* The proof is similar to those of Lemmas 3 and 4 and is therefore omitted.

It is using this lemma that we next state the privacy guarantees afforded to arbitrary linear queries over S.

**Theorem 3.** Let  $L: S \to \Delta_n$  be  $\alpha$ -bounded. Fix an adjacency parameter  $b \in (0,1]$  and let Assumptions 1-3 hold. Then, for adjacency as defined in Definition 3, the Dirichlet mechanism applied to L, denoted  $\mathcal{M}_D(L(\cdot)): S \to \Delta_n$ , is  $(\epsilon, \delta)$ -differentially private, where

$$\begin{split} \epsilon &= \log \left( \frac{\mathrm{beta} \left( k \eta, k (1 - \bar{\eta} - \eta) \right)}{\mathrm{beta} \left( k (\eta + \frac{b \alpha}{2}), k (1 - \bar{\eta} - \eta - \frac{b \alpha}{2}) \right)} \right) \\ &+ \frac{k b \alpha}{2} \log \left( \frac{1 - (|W| - 1) \gamma}{\gamma} \right) \end{split}$$

and

$$\delta = 1 - \min_{p \in \Delta_{n,W}^{(\eta,\bar{\eta})}} \mathbb{P} \big[ \mathcal{M}_D^{(k)}(p) \in \Omega_1 \big].$$

*Proof:* This proof is similar to that of Theorem 2, but with an arbitrary  $\alpha$ -bounded query instead of the average.

**Remark 3.** Setting  $\alpha = \frac{1}{N}$  in Theorem 3 recovers Theorem 2 because the average is a  $\frac{1}{N}$ -bounded query.

As above, we provide a simplified bound on  $\epsilon$  that offers a straightforward dependence of  $\epsilon$  upon other parameters in the problem.

**Corollary 3.** Let all conditions of Theorem 3 hold. Then, for an arbitrary  $\alpha$ -bounded linear query L, the Dirichlet

mechanism  $\mathcal{M}_D(L(\cdot))$  :  $\mathcal{S} \to \Delta_n$  is  $(\epsilon, \delta)$ -differentially private with

$$\epsilon = 2k(1 - \bar{\eta}) - 3 + \frac{kb\alpha}{2}\log\left(\frac{1 - (|W| - 1)\gamma}{\gamma}\right)$$

and

$$\delta = 1 - \min_{p \in \Delta_{n,W}^{(\eta,\bar{\eta})}} \mathbb{P} \big[ \mathcal{M}_D^{(k)}(p) \in \Omega_1 \big].$$

*Proof:* This proof is similar to Corollaries 1 and 2 and is therefore omitted.

#### VI. ACCURACY ANALYSIS

We analyze the accuracy of the Dirichlet mechanism in two ways: first in terms of its moments and concentration about its mean, and second by comparison to the existing Gaussian mechanism for differential privacy.

## A. Analytical Accuracy

**Proposition 1.** Let  $x \in \Delta_n$  be the output of a Dirichlet mechanism with input  $p \in \Delta_{n,W}^{(\eta,\bar{\eta})}$  and parameter  $k \in \mathbb{R}_+$ . Then we have that  $\mathbb{E}[x_i] = p_i$  and

$$Var[x_i] = \frac{p_i(1 - p_i)}{k + 1}.$$
 (20)

*Proof.* Let  $\bar{p} = \sum_{r=1}^{n} k p_r$ . Equation (49.9) in [38] gives

$$\mathbb{E}[x_i] = \frac{kp_i}{\bar{p}} = p_i \quad \text{ and } \quad \operatorname{Var}[x_i] = \frac{kp_i(\bar{p} - kp_i)}{\bar{p}^2(\bar{p} + 1)}.$$

Since the input p belongs to the unit simplex, we have that  $\bar{p} = k$ . Substituting  $\bar{p}$  with k concludes the proof.

**Remark 4.** As seen in (20) the variance of the output depends on the input data  $p_i$ . However, we can find the worst-case variance by maximizing the expression for the variance which occurs at  $p_i = 0.5$ . Hence, we have that

$$Var[x_i] \le \frac{1}{4(k+1)}.$$

It is this form of upper bound that we use to establish the concentration of the Dirichlet mechanism's output about its input. In particular, we bound the probability of a large deviation of the private output from the sensitive input.

**Theorem 4.** Fix  $\eta, \bar{\eta} \in (0,1)$  and  $W \subseteq [n-1]$ , and let Assumptions 1 and 2 hold. Let  $\mathcal{M}_D^{(k)}$  denote a Dirichlet mechanism with parameter  $k \in \mathbb{R}_+$ . Let  $\mu \in (0,1)$  and  $\theta \in (0,e^{-2\mu^2})$  be given. Then a sufficient condition to ensure

$$\mathbb{P}\Big[ \big\| \mathcal{M}_D^{(k)}(p) - p \big\|_{\infty} \le \mu \Big] \ge 1 - \theta$$

is to select  $k = -\frac{\log(\theta)}{2\mu^2} - 1$ .

*Proof:* Lemma 2 in [28] implies that, for any  $\beta > 0$ ,

$$\mathbb{P}\left[\|\mathcal{M}_{D}^{(k)}(p) - p\|_{\infty} \ge \sqrt{\frac{\log(1/\beta)}{2(k+1)}}\right] \ge 1 - \beta.$$

The result follows by setting  $\beta = \theta$ , setting  $\mu = \sqrt{\log(1/\beta)/2(k+1)}$  and solving for k.

Theorem 4 provides both the means to assess accuracy of the Dirichlet mechanism for a given k, as well as a prescriptive tool for selecting k based on desired accuracy.

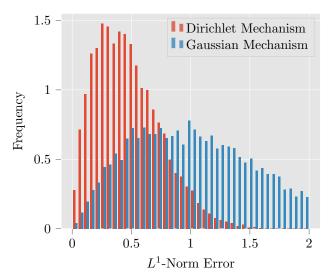


Fig. 3: Histogram of  $L^1$ -norm error when simulating (2.30, 0.05)—differentially private identity queries using the Dirichlet and Gaussian mechanisms over 10,000 random input vectors. This plot shows that on average, the Dirichlet mechanism produces a much smaller error compared to the Gaussian mechanism while providing identical privacy guarantees.

# B. Numerical Comparison to the Gaussian Mechanism

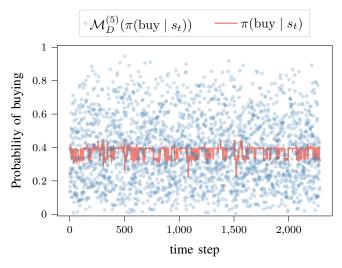
In this section we compare the accuracy of the Dirichlet mechanism to the Gaussian mechanism. The Dirichlet mechanism is designed to provide  $(\epsilon, \delta)$ -differential privacy for data on the unit simplex, and a state-of-the-art method for  $(\epsilon, \delta)$ -differential privacy is the Gaussian mechanism. We therefore use the Gaussian Mechanism and project the outputs back onto the unit simplex [16], [39] to provide a benchmark for the accuracy of the Dirichlet mechanism.

To compare the two mechanisms, we consider identity queries over inputs in  $\Delta_3$ . Let  $\hat{\delta}=0.05$  and k=3. Then Theorem 1 implies that the queries are (2.30,0.05)—differentially private. For an  $(\epsilon,\delta)$ -differentially private Gaussian mechanism for identity queries, added noise must satisfy  $\sigma>\sqrt{2\ln(1.25/\delta)}/\epsilon$  [16]. Thus for (2.30,0.05)-differential privacy we require  $\sigma>1.103$  and we use  $\sigma=1.120$ .

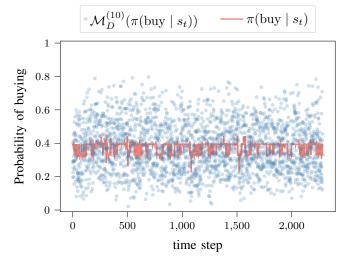
We have simulated the two mechanisms over 10,000 random vectors in the unit simplex and measured the mechanisms'  $L^1$ -norm error, defined as  $\|e(p)\|_1 = \|\mathcal{M}(p) - p\|_1$ . The Dirichlet mechanism had a mean  $L^1$ -norm error of 0.478 and the Gaussian mechanism's mean error was 0.981, which is more than double the average error of the Dirichlet mechanism. A histogram of the results is shown in Figure 3, which illustrates that the Dirichlet mechanism is significantly more accurate than a comparable Gaussian mechanism while providing the same privacy guarantees.

## VII. SIMULATION RESULTS

In this section, we simulate the output of the Dirichlet mechanism for identity and average queries. We consider a reinforcement-learning agent that trades stocks in the AnyTrading environment from Open AI Gym [40]. In the simulations we use Google's stock data that the dataset in



(a) Distribution of private outputs for identity queries of stock trading policies with k=5.



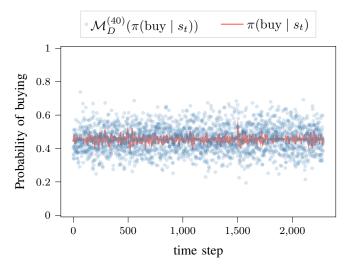
(b) Distribution of private outputs for identity queries of stock trading policies with k=10.

Fig. 4: Distributions of outputs for identity queries on 0.2-adjacent inputs in  $\Delta_2$ . The top figure uses k=5 and the bottom figure uses k=10.

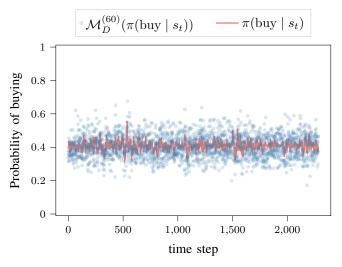
[41] provides. The agent learns a policy to either buy or sell a stock, and its policy at each point in time is therefore an element of  $\Delta_2$ . Privatizing both entries gives  $W = \{1,2\}$ . For identity queries, we train a single agent using the A2C baseline in [42], and we plot the agent's policy at each point in time when privatized with the Dirichlet mechanism with two levels of privacy. For average queries, we repeat this training for four agents and release the average of their policies under two different privacy levels.

### A. Simulation of Identity Queries

In Figure 4, we show an example of privatizing an identity query on a single agent's trading policy. We plot the probability of buying the stock, and the complementary probability is one minus the shown probability. The actual learned trading



(a) Private averages of the stock trading policies learned by four agents with k=40.



(b) Private averages of the stock trading policies learned by four agents with k=60.

Fig. 5: Private averages of four agents' learned stock trading policies. The top plot shows the case of k=40 and the bottom plot shows k=60. As expected, a larger value of k provides weaker privacy protections in exchange for outputs more closely concentrated about the sensitive inputs that produced them.

policy is shown in the solid orange line in the center, while its private form is shown in the blue circles around it. We use the values k=5 and k=10 to illustrate different levels of privacy, and both instances use the adjacency parameter b=0.2. Theorem 1 implies that these identity queries are provided with (1.45,0.04)-differential privacy with k=5 and (2,0.04)-differential privacy with k=10.

As expected, increasing k shows a more concentrated distribution of private outputs about the sensitive inputs to the mechanism. This illustrates that one must tradeoff stronger privacy for reduced accuracy and vice versa.

## B. Simulation of Average Queries

In Figure 5, we show an example of privatizing the average trading policy learned by four agents trading stocks. We again use the adjacency parameter b=0.2. Here, we choose k=40 and k=60. Using Theorem 2, we find that k=40 provides (0.83,0.04)-differential privacy, while k=60 provides (1.02,0.04)-differential privacy.

As seen in Figure 5, the output when k=40 is less concentrated around the average. We note here as well that general linear queries exhibit the same behavior and that is why results specific to that case have been omitted: simulation results for average queries show the concentration of the Dirichlet mechanism's outputs about its inputs when those inputs are functions of collections of vectors, and this concentration is the same for general linear queries.

#### VIII. CONCLUSION

In this work we introduced a mechanism used for privatizing data inputs that belong to the unit simplex. We used the Dirichlet distribution to probabilistically map a vector within the unit simplex to itself. We proved that the Dirichlet mechanism is differentially private with high probability in identity, average, and linear queries. Our simulation results validated that the privacy bounds and the accuracy of the mechanism are within ranges typically considered in the differential privacy literature. As an extension to this work, we are interested in applying the Dirichlet mechanism to privatizing a policy in a Markov decision process. In particular, we are interested in showing how accurate the Dirichlet mechanism is in terms of the total accumulated rewards.

## REFERENCES

- [1] G. W. Hart, "Nonintrusive appliance load monitoring," *Proceedings of the IEEE*, vol. 80, no. 12, pp. 1870–1891, 1992.
- [2] Y. Wang, Z. Huang, S. Mitra, and G. E. Dullerud, "Differential privacy in linear distributed control systems: Entropy minimizing mechanisms and performance tradeoffs," *IEEE Transactions on Control of Network* Systems, vol. 4, no. 1, pp. 118–130, March 2017.
- [3] E. Nozari, P. Tallapragada, and J. Cortés, "Differentially private distributed convex optimization via functional perturbation," *IEEE Transactions on Control of Network Systems*, vol. 5, no. 1, pp. 395–408, March 2018.
- [4] M. Hale, A. Jones, and K. Leahy, "Privacy in feedback: The differentially private lqg," in 2018 Annual American Control Conference (ACC), June 2018, pp. 3386–3391.
- [5] A. Jones, K. Leahy, and M. Hale, "Towards differential privacy for symbolic systems," in 2019 American Control Conference (ACC), July 2019, pp. 372–377.
- [6] M. L. Puterman, Markov Decision Processes.: Discrete Stochastic Dynamic Programming. John Wiley & Sons, 2014.
- [7] R. S. Sutton, A. G. Barto et al., Introduction to reinforcement learning. MIT press Cambridge, 1998, vol. 2, no. 4, ch. 3.
- [8] Y. Savas, M. Ornik, M. Cubuktepe, M. O. Karabag, and U. Topcu, "Entropy maximization for markov decision processes under temporal logic constraints," *IEEE Transactions on Automatic Control*, 2019.
- [9] B. Wu, M. Cubuktepe, and U. Topcu, "Switched linear systems meet markov decision processes: Stability guaranteed policy synthesis," in 2019 IEEE 58th Annual Conference on Decision and Control (CDC). IEEE, 2019, to appear, preprint arXiv:1904.11456.
- [10] K. Chatterjee, R. Majumdar, and T. A. Henzinger, "Markov decision processes with multiple objectives," in *Annual Symposium on Theoretical Aspects of Computer Science*. Springer, 2006, pp. 325–336.
- [11] S. Brechtel, T. Gindele, and R. Dillmann, "Probabilistic mdp-behavior planning for cars," in 2011 14th International IEEE Conference on Intelligent Transportation Systems (ITSC), Oct 2011, pp. 1537–1542.

- [12] S. Misra, A. Mondal, S. Banik, M. Khatua, S. Bera, and M. S. Obaidat, "Residential energy management in smart grid: A markov decision process-based approach," in 2013 IEEE International Conference on Green Computing and Communications and IEEE Internet of things and IEEE Cyber, Physical and Social Computing. IEEE, 2013, pp. 1152–1157.
- [13] C. Yu, J. Liu, and S. Nemati, "Reinforcement learning in healthcare: A survey," arXiv preprint arXiv:1908.08796, 2019.
- [14] X. Pan, W. Wang, X. Zhang, B. Li, J. Yi, and D. Song, "How you act tells a lot: Privacy-leaking attack on deep reinforcement learning," in Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems, ser. AAMAS '19. Richland, SC: International Foundation for Autonomous Agents and Multiagent Systems, 2019, p. 368–376.
- [15] C. Dwork, F. McSherry, K. Nissim, and A. Smith, "Calibrating noise to sensitivity in private data analysis," in *Theory of cryptography* conference. Springer, 2006, pp. 265–284.
- [16] C. Dwork, A. Roth et al., "The algorithmic foundations of differential privacy," Foundations and Trends® in Theoretical Computer Science, vol. 9, no. 3–4, pp. 211–407, 2014.
- [17] S. P. Kasiviswanathan and A. Smith, "On the'semantics' of differential privacy: A bayesian formulation," *Journal of Privacy and Confidential*ity, vol. 6, no. 1, 2014.
- [18] K. Nissim and A. Wood, "Is privacy privacy?" Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences, vol. 376, no. 2128, p. 20170358, 2018.
- [19] J. Cortés, G. E. Dullerud, S. Han, J. Le Ny, S. Mitra, and G. J. Pappas, "Differential privacy in control and network systems," in 2016 IEEE 55th Conference on Decision and Control (CDC). IEEE, 2016, pp. 4252–4272.
- [20] S. Han and G. J. Pappas, "Privacy in control and dynamical systems," Annual Review of Control, Robotics, and Autonomous Systems, vol. 1, pp. 309–332, 2018.
- [21] Z. Xu, K. Yazdani, M. T. Hale, and U. Topcu, "Differentially private controller synthesis with metric temporal logic specifications," 2019.
- [22] C. Hawkins and M. Hale, "Differentially private formation control," in 59th IEEE Conference on Decision and Control (CDC), 2020, Accepted; to be presented. Preprint available at: https://arxiv.org/abs/2004.02744.
- [23] M. Gotz, A. Machanavajjhala, G. Wang, X. Xiao, and J. Gehrke, "Publishing search logs—a comparative study of privacy guarantees," *IEEE Transactions on Knowledge and Data Engineering*, vol. 24, no. 3, pp. 520–532, 2011.
- [24] D. B. Karp, "Normalized incomplete beta function: log-concavity in parameters and other properties," *Journal of Mathematical Sciences*, vol. 217, no. 1, pp. 91–107, 2016.
- [25] N. Holohan, S. Antonatos, S. Braghin, and P. Mac Aonghusa, "The bounded laplace mechanism in differential privacy," *Journal of Privacy* and Confidentiality, vol. 10, no. 1, 2020.
- [26] S. Vadhan, "The complexity of differential privacy," in *Tutorials on the Foundations of Cryptography*. Springer, 2017, pp. 347–450.
- [27] P. Gohari, B. Wu, M. T. Hale, and U. Topcu, "The dirichlet mechanism for differential privacy on the unit simplex," in *Proceedings of the 2020 American Control Conference (ACC)*, 2020, accepted. To be presented. Preprint at: https://arxiv.org/abs/1910.00043.
- [28] P. Gohari, M. Hale, and U. Topcu, "Privacy-preserving policy synthesis in Markov decision processes," in 59th IEEE Conference on Decision and Control (CDC), 2020, Accepted; to be presented. Preprint available at: https://arxiv.org/abs/2004.07778.
- [29] A. Machanavajjhala, D. Kifer, J. Abowd, J. Gehrke, and L. Vilhuber, "Privacy: Theory meets practice on the map," in *Proceedings of the* 2008 IEEE 24th International Conference on Data Engineering. IEEE Computer Society, 2008, pp. 277–286.
- [30] M. Gotz, A. Machanavajjhala, G. Wang, X. Xiao, and J. Gehrke, "Publishing search logs—a comparative study of privacy guarantees," *IEEE Transactions on Knowledge and Data Engineering*, vol. 24, no. 3, pp. 520–532, 2011.
- [31] Y. Chen and X. Ye, "Projection Onto A Simplex," arXiv e-prints, p. arXiv:1101.6081, Jan. 2011.
- [32] J. Le Ny and G. J. Pappas, "Differentially private filtering," *IEEE Transactions on Automatic Control*, vol. 59, no. 2, pp. 341–354, 2013.
- [33] J. Rao and M. Sobel, "Incomplete dirichlet integrals with applications to ordered uniform spacings," *Journal of Multivariate Analysis*, vol. 10, no. 4, pp. 603–610, 1980.
- [34] S. Boyd and L. Vandenberghe, Convex Optimization. New York, NY, USA: Cambridge University Press, 2004.

- [35] F. McSherry and R. Mahajan, "Differentially-private network trace analysis," ACM SIGCOMM Computer Communication Review, vol. 41, no. 4, pp. 123–134, 2011.
- [36] L. Bonomi, L. Xiong, R. Chen, and B. Fung, "Privacy preserving record linkage via grams projections," arXiv preprint arXiv:1208.2773, 2012.
- [37] J. Hsu, M. Gaboardi, A. Haeberlen, S. Khanna, A. Narayan, B. C. Pierce, and A. Roth, "Differential privacy: An economic method for choosing epsilon," in 2014 IEEE 27th Computer Security Foundations Symposium. IEEE, 2014, pp. 398–410.
- [38] S. Kotz, N. Balakrishnan, and N. L. Johnson, Continuous multivariate distributions, Volume 1: Models and applications. John Wiley & Sons, 2004, vol. 1.
- [39] Y. Chen and X. Ye, "Projection onto a simplex," arXiv preprint arXiv:1101.6081, 2011.
- [40] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "Openai gym," 2016.
  [41] M. A. Haghpanah, "gym-anytrading," https://github.com/AminHP/
- [41] M. A. Haghpanah, "gym-anytrading," https://github.com/AminHP/ gym-anytrading/blob/master/gym\_anytrading/datasets/data/STOCKS\_ GOOGL.csv, 2019.
- [42] A. Hill, A. Raffin, M. Ernestus, A. Gleave, A. Kanervisto, R. Traore, P. Dhariwal, C. Hesse, O. Klimov, A. Nichol, M. Plappert, A. Radford, J. Schulman, S. Sidor, and Y. Wu, "Stable baselines," https://github.com/ hill-a/stable-baselines, 2018.
- [43] A. Prékopa, "On logarithmic concave measures and functions," Acta Scientiarum Mathematicarum, vol. 34, pp. 335–343, 1973.
- [44] A. Prékopa, "Logarithmic concave measures with application to stochastic programming," *Acta Scientiarum Mathematicarum*, vol. 32, pp. 301–316, 1971.

#### APPENDIX A: PROOF OF LEMMA 2

We first state a result from [43].

**Lemma 8** (Theorem 3 in [43]). Let  $f_1, \ldots, f_k$  be non-negative and Borel measurable functions defined on  $\mathbb{R}^n$  and let

Borel measurable functions defined on 
$$\mathbb{R}^n$$
 and let  $r(t) = \sup_{\substack{(x_1, \dots, x_k) \ \lambda_1 x_1 + \dots + \lambda_k x_k = t}} f_1(x_1) \cdots f_k(x_k), \quad t \in \mathbb{R}^n,$ 

where  $\lambda_1, \ldots, \lambda_k$  are positive constants satisfying the equality  $\lambda_1 + \cdots + \lambda_k = 1$ . Then, the function r(t) is also Borel measurable and

$$\int_{\mathbb{R}^n} r(t)dt \ge \left(\int_{\mathbb{R}^n} f_1(x_1)^{\frac{1}{\lambda_1}} dx_1\right)^{\lambda_1} \cdots \left(\int_{\mathbb{R}^n} f_k(x_k)^{\frac{1}{\lambda_k}} dx_k\right)^{\lambda_k}.$$

We next review the definition of log-concave functions. A function  $g: \mathbb{R}^n \to \mathbb{R}$  is said to be log-concave if for all  $x_1, x_2 \in \mathbb{R}^n$  and  $\theta \in [0, 1]$ , we have that

$$g(\theta x_1 + (1 - \theta)x_2) \ge (g(x_1))^{\theta} (g(x_2))^{1 - \theta}.$$

This condition is equivalent to

$$g(t) \ge \sup_{\theta u + (1-\theta)v = t} g(u)^{\theta} g(v)^{1-\theta}. \tag{21}$$

Note that g is log-concave if and only if  $\log g$  is concave. Next, for  $x \in \mathbb{R}^n$  and  $p \in \Delta_{n,W}^{(\eta,\bar{\eta})}$  let  $f: \mathbb{R}^n \times \Delta_{n,W}^{(\eta,\bar{\eta})} \to [0,1]$  be defined as

$$f(x,p) = \frac{\prod_{i \in W} x_i^{kp_i - 1} \left( 1 - \sum_{i \in W} x_i \right)^{k(1 - \sum_{i \in W} p_i) - 1}}{B(k\tilde{p}_W)}.$$

For a fixed  $p \in \Delta_{n,W}^{(\eta,\bar{\eta})}$ , let

$$f_1(x) := f(x, p).$$

The function  $f_1(x)$  is the Dirichlet probability distribution function with parameter  $\alpha \in \mathbb{R}^W$ , where  $\alpha := k\tilde{p}_W$ . Since

 $p\in \Delta_{n,W}^{(\eta,\bar{\eta})}$ , we have that  $\alpha_i\geq 1$ , for all  $i\in [|W|]$ . Therefore  $f_1$  is a log-concave function [44, Equation (4.4)]. Then, by (21),

$$f(t_x, p) \ge \sup_{\beta u_x + (1-\beta)v_x = t_x} f(u_x, p)^{\beta} f(v_x, p)^{1-\beta}$$
 (22)

for all  $p \in \Delta_{n,W}^{(\eta,\bar{\eta})}$ , all  $t_x, u_x, v_x \in \mathbb{R}^n$ , and  $\beta \in [0,1]$ . Similarly, for a fixed  $x \in \mathbb{R}^n$ , let

$$f_2(p) := f(x, p).$$

Toward evaluating the Hessian of  $\log f_2(p)$ , recall the trigamma function  $\psi^{(1)}(x) = \frac{d^2}{dx^2} \log (\Gamma(x))$ . Then

$$\left[\nabla^2 \log f_2(p)\right]_{i,j} = \begin{cases} -k^2 \psi^{(1)}(kx) & i = j \text{ and } i, j \in W \\ 0 & \text{otherwise} \end{cases}.$$

The trigamma function is positive on the interval  $(0, \infty)$ . Therefore, the Hessian matrix of  $\log(f_2(p))$  is a diagonal matrix whose diagonal entries are either zero or negative. This provides log-concavity of  $f_2(p)$ . Therefore, using (21),

$$f(x,t_p) \ge \sup_{\beta u_p + (1-\beta)v_p = t_p} f(x,u_p)^{\beta} f(x,v_p)^{1-\beta}$$
 (23)

for all  $x \in \mathbb{R}^n$ , all  $t_p, u_p, v_p \in \Delta_{n,W}^{(\eta,\bar{\eta})}$ , and all  $\beta \in [0,1]$ . Next, fix a choice of  $\lambda \in [0,1]$ , choose  $\tilde{u}_x, \tilde{v}_x, \tilde{u}_p, \tilde{v}_p \in \mathbb{R}^n$  such that

$$\lambda \tilde{u}_x + (1 - \lambda)\tilde{v}_x = t_x$$
 and  $\lambda \tilde{u}_p + (1 - \lambda)\tilde{v}_p = p$ .

Assigning  $u_x$  to x in (23), we find

$$f(u_x, p) \ge \sup_{\beta u_p + (1-\beta)v_p = p} f(u_x, u_p)^{\beta} f(u_x, v_p)^{1-\beta}$$
  
 
$$\ge f(u_x, \tilde{u}_p)^{\lambda} f(u_x, \tilde{v}_p)^{1-\lambda}, \tag{24}$$

where the second inequality follows by setting  $\beta = \lambda$ . Similarly, we can write

$$f(v_{x}, p) \ge \sup_{\beta u_{p} + (1-\beta)v_{p} = p} f(v_{x}, u_{p})^{\beta} f(v_{x}, v_{p})^{1-\beta}$$
  

$$\ge f(v_{x}, \tilde{u}_{p})^{\lambda} f(v_{x}, \tilde{v}_{p})^{1-\lambda}. \tag{25}$$

Revisiting (22), using (24) and (25), we can write

$$f(t_{x}, p) \geq \sup_{\beta u_{x} + (1-\beta)v_{x} = t_{x}} f(u_{x}, p)^{\beta} f(v_{x}, p)^{1-\beta}$$

$$= \sup_{\lambda \tilde{u}_{x} + (1-\lambda)\tilde{v}_{x} = t_{x}} f(\tilde{u}_{x}, p)^{\lambda} f(\tilde{v}_{x}, p)^{1-\lambda}$$

$$\geq \sup_{\lambda \tilde{u}_{x} + (1-\lambda)\tilde{v}_{x} = t_{x}} \left( f(\tilde{u}_{x}, \tilde{u}_{p})^{\lambda^{2}} f(\tilde{v}_{x}, \tilde{v}_{p})^{\lambda(1-\lambda)} \right)$$

$$\cdot f(\tilde{v}_{x}, \tilde{u}_{p})^{(1-\lambda)\lambda} f(\tilde{v}_{x}, \tilde{v}_{p})^{(1-\lambda)^{2}} \right).$$

The first line holds for all  $\beta \in [0, 1]$ , while the second follows from setting  $\beta = \lambda$  for the choice of  $\lambda$  above. Note that

$$\begin{split} \lambda \tilde{u}_x + (1-\lambda)\tilde{v}_x &= \\ \lambda^2 \tilde{u}_x + \lambda (1-\lambda)\tilde{u}_x + \lambda (1-\lambda)\tilde{v}_x + (1-\lambda)^2 \tilde{v}_x. \end{split}$$

Since  $\lambda^2 + \lambda(1-\lambda) + \lambda(1-\lambda) + (1-\lambda)^2 = 1$ , Theorem 8 applies. Therefore, we can write

$$\begin{split} \int_{\mathcal{A}_n} f(t_x, p) dt_x &\geq \\ & \left( \int_{\mathcal{A}_n} f(u_x, \tilde{u}_p) du_x \right)^{\lambda^2} \left( \int_{\mathcal{A}_n} f(u_x, \tilde{v}_p) du_x \right)^{\lambda(1-\lambda)} \\ & \left( \int_{\mathcal{A}_n} f(v_x, \tilde{u}_p) dv_x \right)^{(1-\lambda)\lambda} \left( \int_{\mathcal{A}_n} f(v_x, \tilde{v}_p) dv_x \right)^{(1-\lambda)^2}. \end{split}$$

By renaming the variables  $t_x$ ,  $u_x$  and  $v_x$  to x inside the integrals and merging the similar terms into one, we find

$$\int_{\mathcal{A}_n} f(x, p) dx \ge \left( \int_{\mathcal{A}_n} f(x, \tilde{u}_p) dx \right)^{\lambda} \left( \int_{\mathcal{A}_n} f(x, \tilde{v}_p) dx \right)^{(1-\lambda)},$$

where  $\lambda \tilde{u}_p + (1 - \lambda)\tilde{v}_p = p$ . Therefore,  $\int_{\mathcal{A}_n} f(x, p) dx$  is log-concave in p.

## APPENDIX B: PROOF OF LEMMA 3

Let  $c = p_i + p_j = q_i + q_j$ . Then using (6), we have that

$$\frac{\det(kq_i, kq_j)}{\det(kp_i, kp_j)} = \frac{\Gamma(kq_i)\Gamma(k(c - q_i))}{\Gamma(kp_i)\Gamma(k(c - p_i))}$$

$$= \frac{\Gamma(kq_j)\Gamma(k(c - q_j))}{\Gamma(kp_j)\Gamma(k(c - p_j))}.$$
(26)

Using the definition of the digamma function, we have

$$\frac{\partial}{\partial x} \left[ \frac{\Gamma(x-a)}{\Gamma(x-b)} \right] = \frac{\Gamma(x-a)[\psi^{(0)}(x-a) - \psi^{(0)}(x-b)]}{\Gamma(x-b)}.$$
(28)

Because the digamma function is strictly increasing on the interval  $(0,\infty)$ , the derivative in (28) is positive if and only if x-b < x-a, which is true if and only if a < b. Returning to (26) and (27), we see that (26) is increasing in c if  $q_i < p_i$  and that (27) is increasing in c if  $q_j < p_j$ . Therefore, we will construct an upper bound using (26) if  $q_i < p_i$  and we will construct an upper bound using (27) if  $q_j < p_j$ . For concreteness, suppose  $q_i < p_i$ . Then (26) is an increasing function of c. By definition,  $c = q_i + q_j$  and  $i, j \in W$ . Then  $c \le 1 - \bar{\eta}$  and we find

$$\begin{split} \frac{\operatorname{beta}(kp_i, kp_j)}{\operatorname{beta}(kq_i, kq_j)} &= \frac{\Gamma(kq_i)\Gamma(k(c-q_i))}{\Gamma(kp_i)\Gamma(k(c-p_i))} \\ &\leq \frac{\operatorname{beta}(kq_i, k(1-\bar{\eta}-q_i))}{\operatorname{beta}(kp_i, k(1-\bar{\eta}-p_i))}. \end{split}$$

The other case proceeds identically

## APPENDIX C: PROOF OF LEMMA 5

Convexity of  $\exp(\cdot)$  and Jensen's inequality imply

$$\text{beta}(a,b) \ge \exp\left(\int_0^1 \! \log\left(x^{a-1}(1-x)^{b-1}\right) dx\right) = 2 - (a+b).$$

The upper bound follows from

$$2\alpha\beta \leq \alpha^2 + \beta^2$$
, for all  $\alpha, \beta \in \mathbb{R}$ .

Substituting  $\alpha, \beta$  with  $x^{a-1}$  and  $y^{b-1}$  in the integral representation of the beta function completes the proof.

#### **AUTHORS**



Parham Gohari joined the Department of Electrical and Computer Engineering at the University of Texas at Austin as a Ph.D. student in Fall 2018. Prior to his graduate studies, he attended Sharif University of Technology where he received his B.Sc. in electrical engineering. His research interests include developing theory and algorithms for decision-making systems operating in safety- and privacy-critical environments.



Bo Wu received his B.E. degree from Harbin Institute of Technology, China, in 2008, an M.S. degree from Lund University, Sweden, in 2011 and Ph.D. degree from the University of Notre Dame, USA, in 2018, all in electrical engineering. He is currently a postdoctoral researcher at the Oden Institute for Computational Engineering and Sciences at the University of Texas at Austin. His research interest is to apply formal methods, learning, and control in autonomous systems, such as robotic systems, communication systems, and human-in-the-loop systems,

to provide privacy, security, and performance guarantees.



Calvin Hawkins is a PhD student in Mechanical Engineering at the University of Florida. He received his bachelor's degree in Mechanical Engineering from Wayne State University in May, 2019. His current research interests are broadly in the area of privacy in control, with an emphasis on bringing differential privacy to new classes of data typically used by control systems and quantifying the effects of privacy upon feedback.



Matthew Hale is an Assistant Professor of Mechanical and Aerospace Engineering at the University of Florida. He received his BSE in Electrical Engineering summa cum laude from the University of Pennsylvania in 2012, and his MS and PhD in Electrical and Computer Engineering from the Georgia Institute of Technology in 2015 and 2017, respectively. He directs the Control, Optimization, and Robotics Engineering (CORE) Lab at the University of Florida, and his research interests include multi-agent systems, privacy in control, distributed

optimization, and mobile robotics. He was the Teacher of the Year in the UF Department of Mechanical and Aerospace Engineering for the 2018-2019 school year, and he received an NSF CAREER Award in 2020 for his work on privacy in control systems.



Ufuk Topcu joined the Department of Aerospace Engineering at the University of Texas at Austin as an assistant professor in Fall 2015. He received his Ph.D. degree from the University of California at Berkeley in 2008. He held research positions at the University of Pennsylvania and California Institute of Technology. His research focuses on the theoretical, algorithmic and computational aspects of design and verification of autonomous systems through novel connections between formal methods, learning theory and controls.