

Differential Privacy for Stochastic Matrices Using the Matrix Dirichlet Mechanism

Brandon Fallin*, Calvin Hawkins*, Bo Chen*, Parham Gohari[†],
Alexander Benvenuti*, Ufuk Topcu[†], and Matthew Hale*

Abstract—Stochastic matrices are commonly used to analyze Markov chains, but revealing them can leak sensitive information. Therefore, in this paper we introduce a technique to privatize stochastic matrices in a way that (i) conceals the probabilities they contain, and (ii) still allows for accurate analyses of Markov chains. Specifically, we use differential privacy, which is a statistical framework for protecting sensitive data. To implement it, we introduce the Matrix Dirichlet Mechanism, which is a probabilistic mapping that perturbs a stochastic matrix to provide privacy. We prove that this mechanism provides differential privacy, and we quantify the error induced in private stochastic matrices as a function of the strength of privacy being provided. We then bound the distance between the stationary distribution of the underlying, sensitive stochastic matrix and the stationary distribution of its privatized form. Numerical results show that, under typical conditions, privacy introduces error as low as 5.05% in the stationary distribution of a stochastic matrix.

I. INTRODUCTION

Control applications have become increasingly reliant on user data, e.g., in smart power grids, smart transit applications, and networks of robots [1]–[3]. There are often privacy concerns associated with sharing user data because of what it can reveal. For example, smart appliance usage data, driving routines, and other sensitive data streams can be revealing about a user’s past daily habits or locations, and allow inferences to be drawn about these behaviors in the future [4]–[7]. Data is still needed in many applications, so it is desirable to provide privacy to users while preserving the usefulness of their data.

In this paper, we provide privacy to Markov chain models of systems. Markov chains have been used to model smart power grids, users’ online behavior, and devices on the Internet of Things (IoT) [8]–[10]. Markov chains model a system with random transitions between a finite number of states, and the probabilities of these transitions are represented by a stochastic matrix [11], i.e., a matrix with non-negative entries

whose rows sum to 1. The entries of a stochastic matrix are sensitive because they can reveal how often a user engages in a certain behavior or how likely they are to engage in it, e.g., by revealing the probability of home occupancy at a given time of day, browsing and shopping patterns, or trends in smart device usage [12], [13]. These privacy threats pose a significant risk to individuals, and they motivate us to privatize stochastic matrices.

We do so using differential privacy, which originates in the computer science literature where it was originally used to protect sensitive data when databases are queried [14]. Differential privacy is appealing because (i) it is immune to post-processing, in the sense that arbitrary post-hoc computations on private data do not weaken its privacy protections, and (ii) it is robust to side information [15], in that gaining knowledge of some sensitive information does not weaken differential privacy by much [16]. Strong differential privacy protections can be attained while still providing accurate information [15]. In this paper, we consider identity queries of stochastic matrices, i.e., publishing a stochastic matrix.

There is a growing body of work on differential privacy in decision systems, including in multi-agent control, convex optimization, filtering and estimation, and symbolic systems [17]–[22]. These works generally implement differential privacy for numerical data using the Laplace or Gaussian mechanisms, which add noise to sensitive data before it is shared. However, these mechanisms are a poor fit for the privatization of stochastic matrices. Stochastic matrices have non-negative entries and row sums equal to 1, but the outputs of the Laplace and Gaussian mechanisms will not preserve these properties. Therefore, projection onto the allowable set of data would be required, but this has been shown to destroy the accuracy of private data in similar contexts [23]. Thus, new developments are needed.

In this paper, we (i) develop the Matrix Dirichlet Mechanism, the first differential privacy mechanism for stochastic matrices, (ii) bound the error induced by differential privacy between a privatized stochastic matrix and its non-private form in terms of the strength of privacy, (iii) quantify the utility of privatized stochastic matrices, (iv) bound the distance between the stationary distribution of a stochastic matrix and the stationary distribution of its privatized form, and (v) show in simulation that, under typical conditions, the errors induced by privacy are as low as 5.05%.

The rest of this paper is organized as follows. Section II provides background and problem statements. Section III implements differential privacy and quantifies the error it

*Brandon Fallin, Calvin Hawkins, Bo Chen, Alexander Benvenuti, and Matthew Hale are with the Department of Mechanical and Aerospace Engineering at the University of Florida, Gainesville, FL. Emails: {brandonfallin,calvin.hawkins,bo.chen,abenvenuti,matthewhale}@ufl.edu. BF, CH, BC, AB, and MH were supported in part by NSF under CAREER Grant 1943275, by AFOSR under Grant FA9550-19-1-0169, by ONR under Grant N00014-21-1-2502, and by AFRL under Grant FA8651-23-F-A008.

[†]Parham Gohari is with the Department of Electrical and Computer Engineering University of Texas at Austin, Austin, TX. Ufuk Topcu is with the Department of Aerospace Engineering and Engineering Mechanics at the University of Texas at Austin, Austin, TX. Emails: {pgohari, utopcu}@utexas.edu. PG and UT were supported in part by NSF under CAREER Grant 1652113 and by ONR under Grant N00014-21-1-2502.

induces. Section IV analyzes the trade-off between privacy and accuracy of the stationary distribution of a Markov chain. Section V provides simulations and Section VI concludes.

Notation: We use \mathbb{R} and \mathbb{N} to denote the real and natural numbers, respectively. The set \mathbb{R}_+ denotes the positive reals. We use $|S|$ to denote the cardinality of a finite set S . For $n \in \mathbb{N}$, let $[n] = \{1, \dots, n\}$. We use $\mathbf{1}_n$ to denote the vector of all ones in \mathbb{R}^n .

II. BACKGROUND AND PROBLEM STATEMENTS

This section provides background on Markov chains and differential privacy, and then it states the problems we solve.

A. Unit Simplex and Stochastic Matrices

Each row of a stochastic matrix is an element of the unit simplex. The unit simplex in \mathbb{R}^n is formally defined as the set $\Delta_n = \{x \in \mathbb{R}^n \mid \sum_{i=1}^n x_i = 1, x_i \geq 0 \forall i \in [n]\}$. Next, we define the bordered unit simplex, which is the set of vectors within the unit simplex whose components are a sufficient distance from 0 and 1.

Definition 1 (Bordered Unit Simplex). Let Δ_n° denote the interior of Δ_n . Fix $\eta, \bar{\eta} > 0$ and let $W \subseteq [n-1]$ satisfy $|W| \geq 2$. Then the bordered unit simplex is defined as $\Delta_{n,W}^{(\eta,\bar{\eta})} = \{x \in \Delta_n^\circ \mid \sum_{i \in W} x_i \leq 1 - \bar{\eta}, x_i \geq \eta \forall i \in W\}$. \diamond

We also establish mathematical notation for special functions used throughout this work. $\mathbb{P}[\cdot]$ denotes the probability of an event. For a random variable, $\mathbb{E}[\cdot]$ denotes its expectation and $\text{Var}[\cdot]$ denotes its variance. The notation $\|\cdot\|_1$ denotes the 1-norm of a vector or matrix. The space on which we use $\|\cdot\|_1$ will be clear from context. For $x, a, b \in \mathbb{R}_+$, we use the gamma function $\Gamma(x) = \int_0^\infty z^{x-1} \exp(-z) dz$, the psi function $\psi(x) = \frac{\Gamma'(x)}{\Gamma(x)}$, the bi-variate beta function $\text{beta}(a, b) = \int_0^1 t^{a-1} (1-t)^{b-1} dt = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$, and the regularized incomplete beta function $I_z(a, b) = \frac{\int_0^z x^{a-1} (1-x)^{b-1} dx}{\text{beta}(a, b)}$. Additionally, the gamma function satisfies $\Gamma(k+1) = k\Gamma(k)$.

In this work, we implement privacy by generating random matrices whose rows are in the unit simplex. A building block of this technique is the Dirichlet distribution on the unit simplex. For a parameter $k \in \mathbb{R}_+$ and a vector $p \in \Delta_n^\circ$, the Dirichlet distribution with mean p and parameterized by k is denoted $\mathcal{M}_D^{(k)}$ and defined as

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

where $B(kp) = \frac{\prod_{i=1}^n \Gamma(kp_i)}{\Gamma(k \sum_{i=1}^n p_i)}$.

A matrix whose rows all belong to the unit simplex is known as a stochastic matrix. Fix $n \in \mathbb{N}$. We define \mathcal{S}_n as the set of all $n \times n$ stochastic matrices. Formally, we have $\mathcal{S}_n = \{P \in \mathbb{R}^{n \times n} \mid P_{ij} \geq 0 \text{ for all } i, j \in [n], P\mathbf{1}_n = \mathbf{1}_n\}$.

B. Markov Chains

We now review the necessary background on Markov chains. See [11] for a more detailed exposition. A Markov chain is a sequence of random variables $X_1, X_2, X_3, \dots, X_k$ that possess the Markov property [11], namely the property

that $\mathbb{P}(X_{k+1} = x_{k+1} \mid X_1 = x_1, X_2 = x_2, \dots, X_k = x_k) = \mathbb{P}(X_{k+1} = x_{k+1} \mid X_k = x_k)$. In this work, we consider finite, irreducible, homogeneous Markov chains. That is, Markov chains where (i) its random variables take values in a finite set, (ii) transition is possible from one state to any other state using only transitions of positive probability, and (iii) the transition probabilities are independent of time [11]. We construct a transition probability matrix P where $P_{ij} = \mathbb{P}(X_{k+1} = x_j \mid X_k = x_i)$. Let P_i denote the i^{th} row of the transition probability matrix P .

At each time step, we compute the probability distribution of the states of a Markov chain. Let $\mu_k \in \Delta_n$ denote the probability distribution of states at time k . That is, $\mu_{k,i}$ is the probability that the Markov chain is in state i at time k . Let $\mu_0 \in \Delta_n$ denote the initial state distribution. Multiplying by the transition matrix P on the right updates the distribution by another time step, i.e., $\mu_k^T = \mu_{k-1}^T P$. In general, $\mu_k^T = \mu_0^T P^k$ for all $k \geq 1$. A common way to analyze Markov chains is through their steady-state behavior. Specifically, when P is finite, irreducible, and homogeneous, there exists a limit π as $k \rightarrow \infty$ that must satisfy $\pi^T = \pi^T P$, where π is the *stationary distribution* of the Markov chain.

C. Differential Privacy

We briefly review differential privacy here. See [15] for a more complete exposition. Differential privacy is enforced by a randomized mapping, or *mechanism*, that outputs statistically “similar” private values for “close” pieces of sensitive data. An adjacency relation quantifies how “close” two pieces of data are. In the standard setup, two databases D and D' are adjacent if they differ in one entry [24]. Adjacency is a design choice we make that specifies what must be kept private. We state our adjacency relation for stochastic matrices in Section III-A. The condition that private outputs be statistically “similar” is formalized by the definition of differential privacy itself. In this work, we utilize probabilistic differential privacy, defined as follows.

Definition 2 (Probabilistic Differential Privacy [25]). Let P and Q be two adjacent data sets, let \mathcal{M} be a randomized privacy mechanism, and let \mathcal{S} be the set of possible outputs of the mechanism. The mechanism \mathcal{M} satisfies (ϵ, δ) -probabilistic differential privacy if we can partition the output space \mathcal{S} into two disjoint sets, Ω_1 and Ω_2 , such that for all P , we have $\mathbb{P}[\mathcal{M}(P) \in \Omega_2] \leq \delta$, and, for all Q adjacent to P and all $S \in \Omega_1$ we have $\log\left(\frac{\mathbb{P}[\mathcal{M}(P)=S]}{\mathbb{P}[\mathcal{M}(Q)=S]}\right) \leq \epsilon$. \diamond

The strength of differential privacy is quantified through two parameters: ϵ and δ . The value of ϵ controls the amount of information shared. In the literature, ϵ typically ranges from 0.01 to 10 [26]. The value of δ is the probability that more information is shared than should be allowed by ϵ , and this value typically ranges between 0 and 0.05 [27]. Smaller values of both ϵ and δ imply stronger privacy.

If a mechanism provides probabilistic (ϵ, δ) -differential privacy, then it provides conventional (ϵ, δ) -differential privacy [25]. For a privacy mechanism \mathcal{M} , conventional differential privacy states that for any measurable subset A of the

range of \mathcal{M} and all adjacent P and Q , we have the inequality $\mathbb{P}[\mathcal{M}(P) \in A] \leq e^\epsilon \mathbb{P}[\mathcal{M}(Q) \in A] + \delta$. It is shown that (ϵ, δ) -probabilistic differential privacy implies (ϵ, δ') -conventional differential privacy, where $\delta' < \delta$ [28].

The level of privacy of multiple queries of data can be calculated using methods of composition. While general sequences of private queries cause privacy to weaken [15], this weakening does not occur if the queries are of disjoint subsets of the sensitive data [29]. If the domain of the input to the mechanism is partitioned into disjoint sets and these disjoint sets are queried separately, then the ultimate privacy level is equal to the worst of the privacy guarantees of each query. This is formalized in the following lemma.

Lemma 1 (Parallel Composition [29]). Consider a database denoted by D that is partitioned into the disjoint subsets D_1, D_2, \dots, D_N , and suppose that there are privacy mechanisms $\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_N$, where \mathcal{M}_i is (ϵ_i, δ_i) -differentially private. Then, the release of the queries $\mathcal{M}_1(D_1), \mathcal{M}_2(D_2), \dots, \mathcal{M}_N(D_N)$ provides the database D with $(\max_{i \in [N]} \epsilon_i, \max_{i \in [N]} \delta_i)$ -differential privacy. \square

We will use Lemma 1 to develop a differential privacy mechanism for stochastic matrices in Section III.

D. Problem Statements

We now state the problems that we solve.

Problem 1. *Develop a mechanism that provides (ϵ, δ) -differential privacy to a stochastic matrix P .*

Problem 2. *Bound the expected difference of the entries of the sensitive input P and a private output produced by the mechanism developed in solving Problem 1.*

Problem 3. *Apply the developed mechanism to transition matrices of Markov chains and develop a bound on the distance between the private and non-private stationary distributions to quantify the trade-off between the level of privacy and accuracy.*

We next solve Problem 1 and develop and analyze a privacy mechanism to generate private stochastic matrices.

III. MATRIX DIRICHLET MECHANISM FOR DIFFERENTIAL PRIVACY OF IDENTITY QUERIES

We begin this section by establishing a formal adjacency definition for stochastic matrices and outlining the Matrix Dirichlet Mechanism in Section III-A. We then show the differential privacy guarantees provided by the Matrix Dirichlet Mechanism by computing δ in Section III-B and ϵ in Section III-C. This solves Problem 1. Lastly, the accuracy of private data is analyzed in Section III-D by bounding the expected difference between the entries of a sensitive stochastic matrix P and a private output produced by the Matrix Dirichlet Mechanism. This solves Problem 2.

A. Matrix Dirichlet Mechanism

For the use of the Matrix Dirichlet Mechanism, we require that row i of a sensitive stochastic matrix be in $\Delta_{n, W_i}^{(\eta, \bar{\eta})}$ for all $i \in [n]$. We let W in Definition 1 vary for each row, and we use W_i to denote the set of indices associated with row i . For convenience, we define a map

$$V : [n] \rightarrow \{W_1, \dots, W_n\}, \text{ where } V(i) = W_i. \quad (2)$$

We now define the set of sensitive matrices we consider.

Definition 3 (Stochastic Matrices). Fix $n \in \mathbb{N}$ and, for each $i \in [n]$, fix a collection of indices $W_i \subseteq [n-1]$. Fix $\eta, \bar{\eta} > 0$. Let $\mathcal{S}_{n, V}^{(\eta, \bar{\eta})}$ be the set of all stochastic matrices whose i^{th} row is in $\Delta_{n, W_i}^{(\eta, \bar{\eta})}$ from Definition 1 for all $i \in [n]$. Then $\mathcal{S}_{n, V}^{(\eta, \bar{\eta})} = \{P \in \mathcal{S}_n \mid P_i \in \Delta_{n, V(i)}^{(\eta, \bar{\eta})} \forall i \in [n]\}$, where V is from (2) and P_i is the i^{th} row of P . \diamond

We impose the following assumption on η and $\bar{\eta}$ to ensure that ratios of probability distributions over stochastic matrices are bounded when showing that Matrix Dirichlet Mechanism provides differential privacy.

Assumption 1. For $\mathcal{S}_{n, V}^{(\eta, \bar{\eta})}$ in Definition 3, it holds that $\eta > 0$, $\bar{\eta} > 0$, and $\eta + \bar{\eta} < \frac{1}{2}$. It also holds that $W_i \subseteq [n-1]$ and $|V(i)| = |W_i| \geq 2$ for all $i \in [n]$. \diamond

Next, we state our adjacency relation.

Definition 4 (Adjacency). Fix $n \in \mathbb{N}$, $\eta > 0$, and $\bar{\eta} > 0$. Let Assumption 1 hold. Let $P, Q \in \mathcal{S}_{n, V}^{(\eta, \bar{\eta})}$ be $n \times n$ stochastic matrices, let P_i be the i^{th} row of P , and let P_{ij} be the $i^{\text{th}} j^{\text{th}}$ entry of P ; Q_i and Q_{ij} are defined analogously. For an adjacency parameter $b \in (0, 1]$, P and Q are b -adjacent if, for all $i \in [n]$, there exist $j, k \in W_i$ such that $P_{ij} \neq 0$, $P_{ik} \neq 0$, $P_{i\ell} = Q_{i\ell}$ for all $\ell \neq j, k$, and $\|P_i - Q_i\|_1 \leq b$. \diamond

Section II-C showed that the conventional definition of adjacency allows for databases to differ in one entry. Here, we must consider each row differing in two entries because rows must sum to 1, and it is not possible to change only a single entry. We now formalize the notion of differential privacy for stochastic matrices.

Definition 5 (Differential Privacy for Stochastic Matrices). Fix $n \in \mathbb{N}$, $\eta > 0$, $\bar{\eta} > 0$, and $b \in (0, 1]$. Let Assumption 1 hold. Let $P, Q \in \mathcal{S}_{n, V}^{(\eta, \bar{\eta})}$ be two b -adjacent $n \times n$ stochastic matrices. Let V be defined as in (2). A mechanism $\mathcal{M} : \mathcal{S}_{n, V}^{(\eta, \bar{\eta})} \rightarrow \mathcal{S}_n$ is (ϵ, δ) -differentially private, if, for any measurable subset A of the range of \mathcal{M} and all b -adjacent P, Q , we have $\mathbb{P}[\mathcal{M}(P) \in A] \leq e^\epsilon \mathbb{P}[\mathcal{M}(Q) \in A] + \delta$. \diamond

The query we privatize is the identity query of a stochastic matrix. We now introduce our Matrix Dirichlet Mechanism that randomly maps elements of $\mathcal{S}_{n, V}^{(\eta, \bar{\eta})}$ to \mathcal{S}_n .

Definition 6 (Matrix Dirichlet Mechanism). The Matrix Dirichlet Mechanism Dir_M with parameter $k \in \mathbb{R}_+$ takes as input a stochastic matrix $P \in \mathcal{S}_{n, V}^{(\eta, \bar{\eta})}$, and outputs $\tilde{P} \in \mathcal{S}_n$

via $\tilde{P} \sim \text{Dir}_M(kP)$, where $\text{Dir}_M(kP)$ equals

$$\left(\frac{1}{B(kP_1)} \prod_{j=1}^n X_{1j}^{kP_{1j}-1}, \dots, \frac{1}{B(kP_n)} \prod_{j=1}^n X_{nj}^{kP_{nj}-1} \right)^T$$

Definition 6 shows that the Matrix Dirichlet Mechanism outputs \tilde{P} by applying (1) to each row of the matrix P . The parameter k in Definition 6 can be tuned to adjust the level of privacy provided. Given η and $\bar{\eta}$, we apply the following assumption to the privacy parameter k .

Assumption 2. The privacy parameter k for the Matrix Dirichlet Mechanism satisfies $k \geq \max\{\frac{1}{\eta}, \frac{1}{1-\eta-\bar{\eta}}\}$. \diamond

Sections III-B and III-C will show that the Matrix Dirichlet Mechanism from Definition 6 satisfies (ϵ, δ) -differential privacy in Definition 5. We interpret the rows of the sensitive matrix P as a disjoint partition of the entire sensitive matrix. The Matrix Dirichlet Mechanism privatizes the elements of this disjoint partition independently. This is parallel composition as in Lemma 1. We show that the Matrix Dirichlet Mechanism provides conventional (ϵ, δ) -differential privacy for stochastic matrices by (i) computing the (ϵ_i, δ_i) -probabilistic differential privacy guarantees from Definition 2 for P_i , for all $i \in [n]$, (ii) using the fact that (ϵ_i, δ_i) -probabilistic differential privacy for row P_i implies conventional differential privacy for P_i , for all $i \in [n]$, and (iii) using parallel composition from Lemma 1.

B. Computing δ

In this section, we compute δ_i for each row P_i . We begin by analyzing the rows P_i for $i \in [n]$, which form a disjoint partition of the input P . We choose W_i in Definition 1 such that it satisfies Assumption 1. We partition the output space of the Matrix Dirichlet Mechanism applied to a row P_i into two disjoint sets, Ω_1^i and Ω_2^i . For all $i \in [n]$, fix $\gamma_i \in (0, 1)$ and define the sets Ω_1^i and Ω_2^i as $\Omega_1^i = \{x \in \Delta_n \mid x_j \geq \gamma_i \text{ for all } j \in W_i\}$, and $\Omega_2^i = \Delta_n \setminus \Omega_1^i$. We use the following assumption for γ_i .

Assumption 3. Fix $W_i \subseteq [n-1]$. Then $\gamma_i \leq \frac{1}{|W_i|}$. \diamond

Since Ω_1^i and Ω_2^i are disjoint sets, we have

$$\mathbb{P}[\mathcal{M}_D^{(k)}(P_i) \in \Omega_2^i] = 1 - \mathbb{P}[\mathcal{M}_D^{(k)}(P_i) \in \Omega_1^i]. \quad (3)$$

Using Definition 2, we compute the row-wise probabilities $\mathbb{P}[\mathcal{M}_D^{(k)}(P_i) \in \Omega_2^i]$, which give δ_i for each row P_i using probabilistic differential privacy. From (3), we first compute $\mathbb{P}[\mathcal{M}_D^{(k)}(P_i) \in \Omega_1^i]$ for all $i \in [n]$.

Lemma 2. Fix $n \in \mathbb{N}$, $\eta > 0$ and $\bar{\eta} > 0$. For all $i \in [n]$, fix $W_i \subseteq [n-1]$, let $\mathcal{S}_{n,V}^{(\eta, \bar{\eta})}$ be defined as in Definition 3, and let Assumptions 1, 2, and 3 hold. For all $i \in [n]$, define $\mathcal{A}_{r_i} = \{X_i \in \mathbb{R}^{r_i-1} \mid \sum_{j \in [r_i-1]} X_{ij} \leq 1, X_{ij} \geq \gamma_i \forall j \in W_i\}$ for all $r_i \geq |W_i| + 1$. Then, for the Matrix Dirichlet Mechanism Dir_M with parameter $k \in \mathbb{R}_+$, we have that

$$\mathbb{P}[\mathcal{M}_D^{(k)}(P_i) \in \Omega_1^i] = \frac{1}{B(k\tilde{P}_{W_i})} \int_{\mathcal{A}_{|W_i|+1}} \prod_{j \in W_i} X_{ij}^{kP_{ij}-1} \cdot \left(1 - \sum_{j \in W_i} X_{ij}\right)^{k(1-\sum_{j \in W_i} P_{ij})-1} \prod_{j \in W_i} dX_{ij}.$$

We note that $\tilde{P}_{W_i} \in \Delta_{|W_i|+1}$ is equal to P_i after removing entries with indices outside W_i and an entry equal to $1 - \sum_{j \in W_i} P_{ij}$ is appended as its final entry.

Proof. Apply [23, Lemma 1] to each row of P . \square

Lemma 2 shows that instead of an $(n-1)$ -fold integral of the Dirichlet PDF in (1), the computation of $\mathbb{P}[\mathcal{M}_D^{(k)}(P_i) \in \Omega_1^i]$ can be reduced to a $|W_i|$ -fold integral. From (3), an upper bound for $\mathbb{P}[\mathcal{M}_D^{(k)}(P_i) \in \Omega_2^i]$ can be found by minimizing $\mathbb{P}[\mathcal{M}_D^{(k)}(P_i) \in \Omega_1^i]$, where $\mathbb{P}[\mathcal{M}_D^{(k)}(P_i) \in \Omega_2^i] \leq 1 - \min_{P_i \in \Delta_{n, W_i}^{(\eta, \bar{\eta})}} \mathbb{P}[\mathcal{M}_D^{(k)}(P_i) \in \Omega_1^i] =: \delta_i$.

This value of δ_i is for probabilistic differential privacy from Definition 2. This implies that the same δ_i can be used in conventional differential privacy from Definition 5. We know that each row P_i is a disjoint partition of P . Using Lemma 1 we see that, for privacy of the entire stochastic matrix as given in Definition 5, we have $\delta = \max_{i \in [n]} \delta_i$. Next, we will compute ϵ .

C. Computing ϵ

As above, we begin by analyzing a row P_i for some $i \in [n]$. Fix $\eta, \bar{\eta} \in (0, 1]$ satisfying Assumption 1, $b \in (0, 1]$, and $W_i \subseteq [n-1]$ for all $i \in [n]$. For a given $k \in \mathbb{R}_+$, we must bound the following term to compute ϵ_i : $\log(\mathbb{P}[\mathcal{M}_D^{(k)}(P_i) = X_i] / \mathbb{P}[\mathcal{M}_D^{(k)}(Q_i) = X_i])$, where $X_i \in \Omega_1^i$ and P_i and Q_i are b -adjacent in the sense of Definition 4. We now establish the differential privacy guarantees of the Matrix Dirichlet distribution and solve Problem 1.

Theorem 1. Fix $n \in \mathbb{N}$, $\eta > 0$, $\bar{\eta} > 0$, $b \in (0, 1]$, and $W_i \subseteq [n-1]$ for all $i \in [n]$. Let Assumptions 1, 2, and 3 hold. Let the adjacency relation in Definition 4 hold. Then, the Matrix Dirichlet Mechanism with parameter $k \in \mathbb{R}_+$, defined in Definition 6 and denoted as $\text{Dir}_M(kP)$, is (ϵ, δ) -differentially private, where

$$\epsilon = \log\left(\frac{\text{beta}(k\eta, k(1-\bar{\eta}-\eta))}{\text{beta}(k(\eta + \frac{b}{2}), k(1-\bar{\eta}-\eta - \frac{b}{2}))}\right) + \max_{i \in [n]} \frac{kb}{2} \log\left(\frac{1 - (|W_i| - 1)\gamma_i}{\gamma_i}\right),$$

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

Proof. See [30, Theorem 1]. \square

We next quantify the trade-off between privacy and accuracy for our mechanism.

D. Accuracy

Through providing (ϵ, δ) -differential privacy, the Matrix Dirichlet Mechanism randomizes the entries of the matrix P . Here, we solve Problem 2 and provide an upper bound on how much privacy perturbs the individual entries of P .

Theorem 2. Fix $n \in \mathbb{N}$, $\eta > 0$, $\bar{\eta} > 0$, $b \in (0, 1]$, and $W_i \subseteq [n-1]$ for all $i \in [n]$. Let Assumptions 1, 2, and 3 hold. Let the adjacency relation in Definition 4 hold. Fix

a sensitive stochastic matrix $P \in \mathcal{S}_{n,V}^{(\eta,\bar{\eta})}$ and a parameter $k \in \mathbb{R}_+$. Let $\tilde{P} \sim \text{Dir}_M(kP)$, where the Matrix Dirichlet Mechanism Dir_M is given in Definition 6. Then,

$$\mathbb{E}[|P_{ij} - \tilde{P}_{ij}|] \leq \frac{\Gamma(k)2^{1-k}}{\Gamma(k/2)^2k},$$

and

$$\mathbb{E}[|P_{ij} - \tilde{P}_{ij}|^2] \leq \frac{k}{4(k^2 + k)}.$$

Proof. See [30, Theorem 2]. \square

IV. STATIONARY DISTRIBUTION PERTURBATION BOUND

In this section, we solve Problem 3 by quantifying how the implementation of privacy alters the stationary distribution of a finite, irreducible, homogeneous Markov chain. The change between the private and non-private stationary distribution of a transition probability matrix can be bounded using perturbation theory. We first define the matrix A as $A = I - P$, where I is the identity matrix and P is the transition probability matrix.

Using A , we solve for the fundamental matrix of the Markov chain, Z , which is defined as $Z = (A - \mathbb{1}_n \pi^T)^{-1}$, where π is the stationary distribution of the Markov chain. The fundamental matrix Z always exists for a finite, irreducible, homogeneous Markov chain. We denote the stationary distribution corresponding to a privatized transition matrix as $\tilde{\pi}$. We use the following bound.

Lemma 3 (Stationary Distribution Perturbation Bound [31]). The norm-wise perturbation bound of the stationary distribution of a finite, homogeneous, irreducible Markov chain is of the form: $\|\pi - \tilde{\pi}\|_1 \leq \|Z\|_1 \|P - \tilde{P}\|_1$. \square

The privatized transition probability matrix \tilde{P} is generated by randomizing the original transition probability matrix using the Matrix Dirichlet Mechanism from Definition 6. We therefore bound $\mathbb{E}[\|\pi - \tilde{\pi}\|_1]$.

Theorem 3. Fix $n \in \mathbb{N}$, $\eta > 0$, $\bar{\eta} > 0$, $b \in (0, 1]$, and $W_i \subseteq [n-1]$ for all $i \in [n]$. Let Assumptions 1, 2, and 3 hold. Fix $n \in \mathbb{N}$ and consider a stochastic matrix $P \in \mathcal{S}_{n,V}^{(\eta,\bar{\eta})}$ that corresponds to a finite, homogeneous, irreducible Markov chain. Let π denote its stationary distribution. Let $\tilde{P} \sim \text{Dir}_M(kP)$ be its privatized form with stationary distribution $\tilde{\pi}$. Then the expected distance between the private and non-private stationary distribution can be upper bounded as

$$\mathbb{E}[\|\pi - \tilde{\pi}\|_1] \leq \|Z\|_1 \left\{ \frac{n\Gamma(k)2^{1-k}}{\Gamma(k/2)^2k} + \left(\frac{n-1}{n}\right)^{\frac{1}{2}} \cdot \left\{ \frac{n^2k}{4(k^2+k)} - \max_i \left[4 \frac{\eta^{2k[1-\vartheta_i]} \cdot \vartheta_i^{2k(1-\eta)}}{k^2 \cdot \text{beta}(k\eta, k\vartheta_i)^2} \right] \right\}^{\frac{1}{2}} \right\},$$

where $\vartheta_i = \bar{\eta} + (|W_i| - 1)\eta$, $Z = (A - \mathbb{1}_n \pi^T)^{-1}$, and $|W_i|$ is the number of indices in W_i for a given row P_i .

Proof. See [30, Theorem 3]. \square

Next, we demonstrate the accuracy of the perturbed stationary distribution using simulated results.

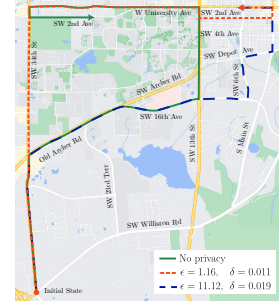


Fig. 1. A random walk over 10 states of the Markov chain model of Gainesville, Florida conducted for varying levels of privacy. As ϵ increases, the sampled path becomes more similar to the random walk without privacy.

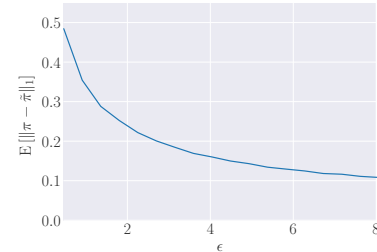


Fig. 2. Using an adjacency parameter of $b = 0.025$ and $\gamma_i = 0.001$ for all $i \in [n]$, the expected value of the error between the private and non-private stationary distribution was found over 10,000 applications of the Matrix Dirichlet Mechanism at each value of $\epsilon \in [0.5, 9]$, with $0 < \delta \leq 0.011$. This result shows that weaker privacy provides greater accuracy for the stationary distribution, and vice versa.

V. SIMULATION RESULTS

This section presents simulation results. We consider a Markov chain model of a traffic system generated from the Annual Average Daily Traffic (AADT) of some of the major streets in Gainesville, Florida from 2021 [32]. The transition probabilities between feasible states of the Markov chain model were found using frequency analysis. In this example, the support of the Markov chain model is assumed to be public knowledge. To this end, the entries in P equal to 0 are not perturbed using the Matrix Dirichlet Mechanism.

As ϵ shrinks and privacy strengthens, the entries of P are perturbed more by privacy, which will affect predictions of routes taken by users. The first transition matrix P was derived directly from the AADT traffic statistics and is treated as the sensitive data. The second transition probability matrix \tilde{P}_1 was generated by privatizing P using Dir_M with a parameter $k = 9.87$ corresponding to $\epsilon = 1.16$ and $\delta = 0.011$. The final transition matrix \tilde{P}_2 was generated by privatizing P using Dir_M with a parameter $k = 98.7$ corresponding to $\epsilon = 11.12$ and $\delta = 0.019$. For both applications of Dir_M , an adjacency parameter $b = 0.025$ was selected. We fix $\eta = 0.10$ and $\bar{\eta} = 0.051$. The set W_i was selected as the indices of the $n - 1$ largest non-zero entries in each row of P , and we set $\gamma_i = 0.001$ for all $i \in [n]$. A random walk was then performed for 10 steps for each of the 3 transition probability matrices. The results illustrated in Figure 1 show that as ϵ increases, the random walk tends toward state transitions similar to those in the sensitive matrix P .

Figure 2 shows the relationship between the strength of privacy and the 1-norm of the difference between the private and non-private stationary distributions, and it was generated by considering $k \in [10, 200]$. For a given k , the values of ϵ and δ were calculated, 10,000 private matrices were generated, and the stationary distribution $\tilde{\pi}$ was computed for each. The average stationary distribution over the 10,000 private responses was calculated and subtracted from the original stationary distribution, π .

Figure 2 shows that the simulated error monotonically decreases as privacy increases. Figure 2 shows a maximum expected error of 0.485 corresponding to $\epsilon = 0.5$ and a minimum error of 0.101 for $\epsilon = 8$. With respect to the maximum possible error, 2, the private stationary distribution error ranges from 5.05% to 24.25% under typical privacy conditions. This error is spread over the 32 entries of the stationary distribution, thus the error associated with a single state is quite small, even for strong privacy guarantees.

VI. CONCLUSION

This paper introduced the Matrix Dirichlet Mechanism to provide differential privacy to stochastic matrices. We proved that this mechanism satisfies differential privacy guarantees, and quantified the error induced in private stochastic matrices as a function of the strength of privacy. The mechanism was then applied to Markov chains and we quantified how the stationary distribution of the sensitive stochastic matrix is changed by privacy. Future work includes an extension of the Matrix Dirichlet Mechanism to doubly stochastic matrices.

REFERENCES

- [1] F. Fioretto, T. W. Mak, and P. Van Hentenryck, "Differential privacy for power grid obfuscation," *IEEE Transactions on Smart Grid*, vol. 11, no. 2, pp. 1356–1366, 2019.
- [2] Y. Li, D. Yang, and X. Hu, "A differential privacy-based privacy-preserving data publishing algorithm for transit smart card data," *Transportation Research Part C: Emerging Technologies*, vol. 115, p. 102634, 2020.
- [3] A. Prorok and V. Kumar, "A macroscopic model for differential privacy in dynamic robotic networks," *arXiv preprint arXiv:1703.04797*, 2017.
- [4] M. R. Alam, M. B. I. Reaz, and M. M. Ali, "Speed: An inhabitant activity prediction algorithm for smart homes," *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans*, vol. 42, no. 4, pp. 985–990, 2011.
- [5] Q. Gong, S. Midlam-Mohler, V. Marano, and G. Rizzoni, "An iterative markov chain approach for generating vehicle driving cycles," *SAE International Journal of Engines*, vol. 4, no. 1, pp. 1035–1045, 2011.
- [6] A. Asahara, K. Maruyama, A. Sato, and K. Seto, "Pedestrian-movement prediction based on mixed markov-chain model," in *Proceedings of the 19th ACM SIGSPATIAL international conference on advances in geographic information systems*, 2011, pp. 25–33.
- [7] M. Hale, P. Barooah, K. Parker, and K. Yazdani, "Differentially private smart metering: Implementation, analytics, and billing," in *Proceedings of the 1st ACM International Workshop on Urban Building Energy Sensing, Controls, Big Data Analysis, and Visualization*. Association for Computing Machinery, 2019, p. 33–42.
- [8] J. Widén, A. M. Nilsson, and E. Wäckelgård, "A combined markov-chain and bottom-up approach to modelling of domestic lighting demand," *Energy and Buildings*, vol. 41, no. 10, pp. 1001–1012, 2009.
- [9] S. Vermeer and D. Trilling, "Toward a better understanding of news user journeys: A markov chain approach," *Journalism Studies*, vol. 21, no. 7, pp. 879–894, 2020.
- [10] J. Dong, G. Li, W. Ma, and J. Liu, "Personalized recommendation system based on social tags in the era of internet of things," *Journal of Intelligent Systems*, vol. 31, no. 1, pp. 681–689, 2022.
- [11] D. A. Levin and Y. Peres, *Markov chains and mixing times*. American Mathematical Soc., 2017, vol. 107.
- [12] J. Lundström, E. Järpe, and A. Verikas, "Detecting and exploring deviating behaviour of smart home residents," *Expert Systems with Applications*, vol. 55, pp. 429–440, 2016.
- [13] S. Rendle, C. Freudenthaler, and L. Schmidt-Thieme, "Factorizing personalized markov chains for next-basket recommendation," in *Proceedings of the 19th international conference on World wide web*, 2010, pp. 811–820.
- [14] C. Dwork, F. McSherry, K. Nissim, and A. Smith, "Calibrating noise to sensitivity in private data analysis," in *Theory of Cryptography: Third Theory of Cryptography Conference, TCC 2006, New York, NY, USA, March 4-7, 2006. Proceedings 3*. Springer, 2006, pp. 265–284.
- [15] 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.
- [16] S. P. Kasiviswanathan and A. Smith, "On the 'semantics' of differential privacy: A bayesian formulation," *Journal of Privacy and Confidentiality*, vol. 6, no. 1, 2014.
- [17] C. Hawkins and M. Hale, "Differentially private formation control," in *2020 59th IEEE Conference on Decision and Control (CDC)*. IEEE, 2020, pp. 6260–6265.
- [18] G. Yuan, Y. Yang, Z. Zhang, and Z. Hao, "Convex optimization for linear query processing under approximate differential privacy," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016, pp. 2005–2014.
- [19] J. Le Ny, "Differentially private kalman filtering," *Differential Privacy for Dynamic Data*, pp. 55–75, 2020.
- [20] B. Chen, K. Leahy, A. Jones, and M. Hale, "Differential privacy for symbolic systems with application to markov chains," *Automatica*, vol. 152, p. 110908, 2023.
- [21] Y. Wang, M. Hale, M. Egerstedt, and G. E. Dullerud, "Differentially private objective functions in distributed cloud-based optimization," in *2016 IEEE 55th Conference on Decision and Control (CDC)*, 2016, pp. 3688–3694.
- [22] P. Gohari, M. Hale, and U. Topcu, "Privacy-preserving policy synthesis in markov decision processes," in *2020 59th IEEE Conference on Decision and Control (CDC)*, 2020, pp. 6266–6271.
- [23] P. Gohari, B. Wu, C. Hawkins, M. Hale, and U. Topcu, "Differential privacy on the unit simplex via the dirichlet mechanism," *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 2326–2340, 2021.
- [24] C. Dwork and J. Lei, "Differential privacy and robust statistics," in *Proceedings of the forty-first annual ACM symposium on Theory of computing*, 2009, pp. 371–380.
- [25] A. Machanavajjhala, D. Kifer, J. Abowd, J. Gehrke, and L. Vilhuber, "Privacy: Theory meets practice on the map," in *2008 IEEE 24th international conference on data engineering*. IEEE, 2008, pp. 277–286.
- [26] 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.
- [27] C. Hawkins, B. Chen, K. Yazdani, and M. Hale, "Node and edge differential privacy for graph laplacian spectra: Mechanisms and scaling laws," *arXiv preprint arXiv:2211.15366*, 2022.
- [28] 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.
- [29] N. Ponomareva, H. Hazimeh, A. Kurakin, Z. Xu, C. Denison, H. B. McMahan, S. Vassilvitskii, S. Chien, and A. Thakurta, "How to dp-fy ml: A practical guide to machine learning with differential privacy," *arXiv preprint arXiv:2303.00654*, 2023.
- [30] B. Fallin, C. Hawkins, B. Chen, P. Gohari, A. Benvenuti, U. Topcu, and M. Hale, "Technical report: differentially privacy for stochastic matrices." [Online]. Available: <http://corelab.mae.ufl.edu/papers/Fallin-CDC23-TR.pdf>
- [31] E. Seneta, "Sensitivity to perturbation of the stationary distribution: some refinements," *Linear Algebra and its Applications*, vol. 108, pp. 121–126, 1988.
- [32] "Florida traffic online." [Online]. Available: <https://tdaappsprod.dot.state.fl.us/fto/>