

---

# LOCAL DIFFERENTIAL PRIVACY PRESERVING MECHANISMS FOR MULTI-AGENT REINFORCEMENT LEARNING

---

**Prashant Trivedi**

Industrial Engineering and Operations Research  
Indian Institute of Technology Bombay India  
trivedi.prashant15@iitb.ac.in

**Nandyala Hemachandra**

Industrial Engineering and Operations Research  
Indian Institute of Technology Bombay India  
nh@iitb.ac.in

## ABSTRACT

This work considers the local differential privacy (LDP) aspects of multi-agent reinforcement learning (MARL). We design a fully decentralized and generic multi-agent locally differential private (MA-LDP) algorithm that can handle any noise adding mechanisms. It takes the episodic form, where the data from each agent is anonymized using noise adding mechanism in each episode. MA-LDP uses the linear function approximations of the transition probabilities and the reward functions. We prove that the MA-LDP algorithm preserves the user's data privacy and attains the sub-linear regret for four noise mechanisms with different noise supports. Further, we compare the noise adding mechanisms with unbounded support to those with bounded support. A key observation is that for a suitably chosen bounded noise support, the regret of the MA-LDP algorithms is on-par or lower than the noise mechanism with unbounded support. We validate our theoretical findings on a network MDP with a large state and action spaces.

**Keywords** Differential privacy · Decentralized models · Noise mechanisms · Privacy loss · Sub-linear regret · Finite horizon MDP · Network MDP · Linear function approximations · Reinforcement learning · Bellman equations

## 1 Introduction

Multi-agent reinforcement learning (MARL) is one of the powerful tools used widely in many real-life applications. However, most of the algorithms rely on the user given data. Thus, the privacy of the users data is an utmost concern of any user. To this end, [Dwork et al. \(2006\)](#) introduced the notion of the differential privacy. The idea is to work with the anonymized user data, yet achieve the same performance and the user experience. Therefore, the notion of Locally Differentially Private (LDP) is introduced in [Dwork and Roth \(2014\)](#); [Kasiviswanathan et al. \(2011\)](#), which enhanced the performance while maintaining the privacy. However, to our knowledge the notion of LDP is not considered in MARL settings, where the privacy of user data is more serious. LDP in MARL can be used to protect the sensitive information that is shared or collected among the agents. Some of the applications include the health data analysis where a patients data is sensitive, interactions of agents within a network, say a social media platform etc.

In the MARL setup the agents usually have a common objective; however, the decisions are taken sequentially by each agent in a decentralised fashion. In this work, we introduce the notion of differential privacy for the MARL problem. In particular, we address the following questions. 1) Can we design a MARL algorithm that can preserve the user-data privacy while attaining the same outcome? We answer this question affirmatively and introduce a novel and generic Multi-Agent Local Differential Privacy (MA-LDP) MARL algorithm. This algorithm can work with any noise mechanism. In this work, we consider four noise adding mechanisms, Gaussian, Laplace, uniform and bounded Laplace (BL). Two of which have the unbounded support, whereas other two have bounded support. For each noise adding mechanism we investigate 3 aspects of our MA-LDP algorithm 1) its ability to preserve user sensitive information, 2) the effect on the regret, and 3) how its regret fares with that of the other mechanisms. The privacy from these noise mechanisms ranges from one extreme of no privacy loss ( $\epsilon = 0$ ) and less confidence ( $\delta > 0$ ) to some privacy loss ( $\epsilon > 0$ ) and high confidence ( $\delta = 0$ ).

Next, we compare the regret and the privacy guarantees of our MA-LDP algorithm across these noise mechanisms. This comparison is motivated by the following question: can the noise mechanisms with bounded support attain the same or better privacy and the regret guarantee as that of the noise mechanisms with unbounded support. Again we answer this question affirmatively. In particular, we show that the regret for the BL mechanism is of the same order as that of the Laplace mechanism, if the end points of the BL mechanism are chosen suitably. Our contributions are

(a) To address the LDP aspects in the MARL we introduce four novel noise adding mechanisms. Two of these mechanisms have unbounded noise support, whereas the remaining two have bounded support. We propose a generic algorithm MA-LDP that can work with any noise mechanism. Moreover, with these four noise mechanisms our MA-LDP achieves a sub-linear regret (Sections 5, 6).

(b) These four different noise adding mechanisms yield different privacy guarantees; Gaussian mechanism achieves  $(\epsilon, \delta)$  privacy, Laplace mechanism achieves  $(\epsilon, 0)$  privacy, uniform mechanism achieves  $(0, \delta)$ , and finally for bounded Laplace noise mechanism we get  $(\epsilon, 0)$  privacy. Thus, the privacy guarantees range from one extreme  $(\epsilon, 0)$  to  $(0, \delta)$  (Section 5, 6)

(c) While the regret of MA-LDP algorithm is sub-linear in total numbers of steps, it is super-linear in number of agents and the feature dimension for all the noise mechanisms. But this order with respect to number of agents and feature dimension is the same across all the four noise mechanisms.

(d) We compare the regret of MA-LDP standard (unbounded support) noise mechanisms with those of the bounded support mechanisms. If the end points of the support of bounded Laplace distribution is of the same order as that of the distribution parameter .. the regret of BL mechanism is on par with the Laplace mechanism. So, instead of injecting the potentially unbounded noise values, a suitably chosen bounded support can attain the lower regret (see Table 2, Section 7).

## 2 LDP for Decentralized MARL

In this Section, we first introduce the episodic decentralized MARL, and then provide the notion of LDP for the MARL.

### 2.1 Episodic Decentralized MARL

Let  $N = \{1, 2, \dots, n\}$  be the set of agents. An instance of multi-agent time inhomogeneous episodic Markov Decision Process is given by  $(N, \mathcal{S}, \{\mathcal{A}^i\}_{i \in N}, H, \{r_h^i\}_{i \in N, h \in H}, \{\mathbb{P}_h\}_{h \in H}, \{\mathcal{G}_t\}_{t \geq 0})$ . Here  $\mathcal{S}$  is the finite set of global state-space.  $\mathcal{A}^i$  is the finite set of local actions available to agent  $i \in N$ , i.e., the local action-space of agent  $i \in N$ . The global action at any global state  $s \in \mathcal{S}$  is  $\mathcal{A}(s) = \prod_{i \in N} \mathcal{A}^i(s)$ . We use product of the local action-space because the agents are independently taking the actions. A typical element in  $\mathcal{A}(s)$  is a vector of size  $n$ , one for each agent as  $(\mathbf{a}^1(s), \mathbf{a}^2(s), \dots, \mathbf{a}^n(s))$ , where  $\mathbf{a}^i(s) \in \mathcal{A}^i(s)$  represents the action taken by agent  $i$  when the global state is  $s$ . Let  $K$  be the total number of episodes; each episode consists of a fixed planning horizon  $H$ . We use  $h$  to denote an intermediate stage in planning horizon  $H$ . Moreover, let  $T = KH$  be the total number of interactions with the MDP. For a global state  $s$ , global action  $\mathbf{a}$ , each agent  $i \in N$  realizes a deterministic local reward  $r_h^i(s, \mathbf{a}) \in [0, 1]$  at the stage  $h$  of the planning horizon. It is important to note that the reward of agent  $i$  depends on the action taken by other agents and the global state. Moreover, this reward  $r_h^i$  is a private information of the agent  $i$  and hence not known to other agents. Once an action  $\mathbf{a}$  is taken in state  $s$  at stage  $h$ , the state change to  $s'$  with probability  $\mathbb{P}_h(s'|s, \mathbf{a})$ . In the MARL, this transition probability is generally unknown. Finally,  $\mathcal{G}_t$  is the time varying network that allows the sharing of some information across the agents at time  $t = kh$ . We describe more details about it later.

One of the goals in MARL is to learn a policy  $\pi = \{\pi_h\}_{h=1}^H$ , that is the collection of  $H$  functions, where each  $\pi_h(s)$  is the global action taken in the global state  $s$  at stage  $h$ . In particular, the aim is to learn an optimal policy, in a decentralized way, by interacting with the environment and observing the past information. To this end, we define the global state-action value function for the policy  $\pi$  at stage  $h$   $Q_h^\pi(s, \mathbf{a}) = \bar{r}_h(s, \mathbf{a}) + \mathbb{E}_\pi \left[ \sum_{h'=h+1}^H \bar{r}_{h'}(s_{h'}, \pi_{h'}(s_{h'})) \right]$ .

Here  $\bar{r}_h = \frac{1}{n} \sum_{i \in N} r_h^i$  and  $s_h = s$ ,  $\mathbf{a}_h = \mathbf{a}$  and  $s_{h'+1} \sim \mathbb{P}_{h'}(\cdot | s_{h'}, \mathbf{a}_{h'})$ . Note that this state-action value function represents the value of policy  $\pi$  when action  $\mathbf{a}$  is taken at stage  $h$  in the state  $s_h = s$ , and the policy  $\pi_h(s)$  is followed thereafter. Moreover, we also define the global state value function at stage  $h$  as  $V_h^\pi(s) = Q_h^\pi(s, \pi_h(s))$ . As opposed to the state action value function the state value function represents the value of a state  $s$  when policy  $\pi_h$  is followed starting for state  $s$  at stage  $h$ .

For simplicity of notation, given any function  $V : \mathcal{S} \rightarrow [0, H]$ , for all  $(s, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}$ , we define the cost-to-go function as  $\mathbb{P}_h V(s, \mathbf{a}) = \sum_{s' \in \mathcal{S}} \mathbb{P}_h(s'|s, \mathbf{a}) V(s')$ . Using above we write the Bellman equation for the policy  $\pi$  for all  $(s, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}$  as  $Q_h^\pi(s, \mathbf{a}) = \bar{r}_h(s, \mathbf{a}) + \mathbb{P}_h V_{h+1}^\pi(s, \mathbf{a})$ ;  $V_h^\pi(s) = Q_h^\pi(s, \pi_h(\mathbf{a}))$ ;  $V_{H+1}^\pi(s) = 0$ . Let the optimal state-action value be  $Q_h^*(s, \mathbf{a}) = \max_\pi Q_h^\pi(s, \mathbf{a})$ , and the optimal state value function be  $V_h^*(s) = \max_\pi V_h^\pi(s)$ .

The Bellman optimality equation for all  $(s, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}$  satisfies  $Q_h^*(s, \mathbf{a}) = \bar{r}_h(s, \mathbf{a}) + \mathbb{P}_h V_{h+1}^*(s, \mathbf{a})$ ;  $V_{h+1}^*(s) = \max_{\mathbf{a} \in \mathcal{A}} Q_h^*(s, \mathbf{a})$ ;  $V_{H+1}^*(s) = 0$ . Since  $\bar{r}_h(\cdot, \cdot)$  is a bounded function, for any policy  $\pi$ , both  $V^\pi(\cdot)$  and  $Q^\pi(\cdot, \cdot)$  are also bounded. This work assumes that the transition probability function  $\mathbb{P}_h$  is written as the linear mixture of given basis functions (Min et al. (2022); Vial et al. (2022); Liao et al. (2021)). In particular, we make the following standard assumption about the transition probability function.

**Assumption 1** (Transition probability approximation). *Suppose the feature mapping  $\phi : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}^{nd}$  is known and pre-given. There exists a  $\theta_h^* \in \mathbb{R}^{nd}$  with  $\|\theta_h^*\|_2 \leq \sqrt{nd}$  such that  $\mathbb{P}_h(s'|s, \mathbf{a}) = \langle \phi(s'|s, \mathbf{a}), \theta_h^* \rangle$  for any triplet  $(s', \mathbf{a}, s) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$  and stage  $h$ . Also, for a bounded function  $V : \mathcal{S} \mapsto [0, H]$ , it holds that  $\|\phi_V(s, \mathbf{a})\|_2 \leq H$ , where  $\phi_V(s, \mathbf{a}) = \sum_{s' \in \mathcal{S}} \phi(s'|s, \mathbf{a})V(s')$ .*

Using above we have,  $\mathbb{P}_h V(s, \mathbf{a}) = \sum_{s' \in \mathcal{S}} \langle \phi(s'|s, \mathbf{a}), \theta_h^* \rangle V(s') = \langle \phi_V(s, \mathbf{a}), \theta_h^* \rangle$ . Moreover, recall the Bellman optimality equation use the averaged reward  $\bar{r}_h(\cdot, \cdot)$ , however in our decentralized model this averaged reward is not known to any agent. To this end, each agent maintains an estimate the globally averaged reward function  $\bar{r}_h(\cdot, \cdot)$ .

Let  $\bar{r}_h(\cdot, \cdot; \mathbf{w}) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  be the class of parameterized functions where  $\mathbf{w} \in \mathbb{R}^p$  for some  $p \ll |\mathcal{S}||\mathcal{A}|$ . To obtain the estimate  $\bar{r}_h(\cdot, \cdot; \mathbf{w})$  we seek to minimize the following  $\min_{\mathbf{w}} \mathbb{E}_{s, \mathbf{a}} [\bar{r}_h(s, \mathbf{a}; \mathbf{w}) - \bar{r}_h(s, \mathbf{a})]^2$ . This optimization problem can be equivalently characterized as (both have the same stationary points)  $\min_{\mathbf{w}} \sum_{i=1}^n \mathbb{E}_{s, \mathbf{a}} [\bar{r}_h(s, \mathbf{a}; \mathbf{w}) - r_h^i(s, \mathbf{a})]^2$ . The details of this equivalence of the optimization problems is given in the Appendix G.1. This motivates the following updates for parameters of the global reward function parameters  $\mathbf{w}^i$  by agent  $i \in N$ ,  $\tilde{\mathbf{w}}_{k,h}^i \leftarrow \mathbf{w}_{k,h}^i + \gamma_{k,h} \cdot [r_h^i(\cdot, \cdot) - \bar{r}_h(\cdot, \cdot; \mathbf{w}_{k,h}^i)] \cdot \nabla_{\mathbf{w}} \bar{r}_h(\cdot, \cdot; \mathbf{w}_{k,h}^i)$ ;  $\mathbf{w}_{k+1,h}^i = \sum_{j \in N} l_{k,h}(i, j) \tilde{\mathbf{w}}_{k,h}^j$ . where  $l_{k,h}(i, j)$  is the  $(i, j)$ -th entry of the consensus matrix  $L_{k,h}$  obtained using communication network  $\mathcal{G}_{k,h}$  in the stage  $h$  of the  $k$ -th episode.  $\gamma_{k,h}$  is the step-size satisfying  $\sum_{k,h} \gamma_{k,h} = \infty$  and  $\sum_{k,h} \gamma_{k,h}^2 < \infty$ , and  $\bar{r}_h(\cdot, \cdot; \mathbf{w}_{k,h}^i)$  is the estimate of global reward function by agent  $i$  in the stage  $h$  of the episode  $k$ . We make following standard assumptions on  $\{L_{k,h}\}_{k,h \geq 0}$  and the features for the reward function approximation Zhang et al. (2018); Trivedi and Hemachandra (2022, 2023).

**Assumption 2** (Consensus matrix). *The consensus matrices  $\{L_t\}_{t \geq 0} \subseteq \mathbb{R}^{n \times n}$  satisfies (i)  $L_t$  is row stochastic, and  $\mathbb{E}(L_t)$  is column stochastic. Further, there exists a constant  $\kappa \in (0, 1)$  such that for any  $l_t(i, j) > 0$ , we have  $l_t(i, j) \geq \kappa$ ; (ii) Consensus matrix  $L_t$  respects  $\mathcal{G}_t$ , i.e.,  $l_t(i, j) = 0$ , if  $(i, j) \notin \mathcal{E}_t$ ; (iii) The spectral norm of  $\mathbb{E}[L_t^\top (I - \mathbb{1}\mathbb{1}^\top/n)L_t]$  is smaller than one.*

**Assumption 3** (Full rank). *For each agent  $i \in N$ , the reward function  $\bar{r}(s, \mathbf{a})$  is parameterized as  $\bar{r}(s, \mathbf{a}; \mathbf{w}) = \langle \psi(s, \mathbf{a}), \mathbf{w} \rangle$ . Here  $\psi(s, \mathbf{a}) = [\psi_1(s, \mathbf{a}), \dots, \psi_k(s, \mathbf{a})] \in \mathbb{R}^p$  are the features associated with pair  $(s, \mathbf{a})$ . We assume that features are uniformly bounded. Let the feature matrix  $\Psi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times p}$  have  $[\psi_m(s, \mathbf{a}), s \in \mathcal{S}, \mathbf{a} \in \mathcal{A}]^\top$  as its  $m$ -th column, then  $\Psi$  has full column rank.*

Since the global reward is unknown, each agent uses the parameterized reward and maintains its estimate of  $V_h(\cdot)$  and  $Q_h(\cdot, \cdot)$ . Let  $V_h^i(\cdot)$  and  $Q_h^i(\cdot, \cdot)$  be the estimate of these functions by agent  $i$  at stage  $h$ . So, the modified Bellman optimality equation for all  $(s, \mathbf{a})$  and for all agents  $i \in N$  is  $Q_h^{*,i}(s, \mathbf{a}; \mathbf{w}_{k,h}^i) = \bar{r}_h(s, \mathbf{a}; \mathbf{w}_{k,h}^i) + \mathbb{P}_h V_{h+1}^{*,i}(s, \mathbf{a}; \mathbf{w}_{k,h}^i)$ ;  $V_{h+1}^{*,i}(s; \mathbf{w}_{k,h}^i) = \max_{\mathbf{a} \in \mathcal{A}} Q_h^{*,i}(s, \mathbf{a}; \mathbf{w}_{k,h}^i)$ ;  $V_{H+1}^{*,i}(s; \mathbf{w}_{k,H+1}^i) = 0$ . We later show that as the number of episodes increases the reward function parameters converge, i.e., as  $k \rightarrow \infty$ ,  $\mathbf{w}_{k,h}^i \rightarrow \mathbf{w}^*$  a.s. for all  $i \in N$  and for all  $h \in [H]$ . Hence,  $\bar{r}_h(s, \mathbf{a}; \mathbf{w}_{k,h}^i) \rightarrow \bar{r}_h(s, \mathbf{a}; \mathbf{w}^*)$ ,  $Q_h^{*,i}(s, \mathbf{a}; \mathbf{w}_{k,h}^i) \rightarrow Q_h^{*,i}(s, \mathbf{a})$  and  $V_h^{*,i}(s; \mathbf{w}_{k,h}^i) \rightarrow V_h^{*,i}(s)$  as  $\bar{r}_h(s, \mathbf{a}; \mathbf{w}_{k,h}^i)$ ,  $Q_h^{*,i}(s, \mathbf{a}; \mathbf{w}_{k,h}^i)$  and  $V_h^{*,i}(s; \mathbf{w}_{k,h}^i)$  are continuous functions of  $\mathbf{w}^i$ , where  $Q_h^{*,i}(s, \mathbf{a})$  and  $V_h^{*,i}(s)$  are defined as  $Q^{*,i}(s, \mathbf{a}) = \bar{r}(s, \mathbf{a}; \mathbf{w}^*) + \mathbb{P}V^{*,i}(s, \mathbf{a})$ ;  $V^{*,i}(s) = \max_{\mathbf{a} \in \mathcal{A}} Q^{*,i}(s, \mathbf{a})$ .

At each episode  $k$ , with the initial state  $s_{k,1}$  agents choose a policy  $\pi_k$ . Moreover, in each stage  $h \in [H]$  of the episode  $k$ , for the observed state  $s_{k,h}$ , agents take an action according to the policy  $\pi_k$ , i.e.,  $\mathbf{a}_{k,h} = \pi_{k,h}(s_{k,h})$ . Further, they also observe a next state  $s_{k,h+1} \sim \mathbb{P}_h(\cdot | s_{k,h}, \mathbf{a}_{k,h})$ . The expected regret in the  $k$ -th episode is  $R_k = \frac{1}{n} \sum_{i \in N} \{V_1^{*,i}(s_1^k) - V_1^i(s_1^k)\}$ . Hence our objective is to design an algorithm with the sub-linear regret. The total expected regret in  $K$  episodes is defined as  $R_K = \sum_{k=1}^K \left( \frac{1}{n} \sum_{i \in N} \{V_1^{*,i}(s_1^k) - V_1^i(s_1^k)\} \right)$ . Note that we are using  $V_1^{*,i}(s_1^k)$  instead of  $V_1^*(s_1^k)$ . This is because  $\mathbf{w}_{k,h}^i \rightarrow \mathbf{w}^*$  and hence  $V_1^{*,i}(s_1^k) = V_1^*(s_1^k)$  for all  $i \in N$ . Such an independence of agents is desirable and is a sign of good decentralised algorithms.

## 2.2 LDP for MARL

In this Section, we introduce the notion of the multi-agent local differential privacy (MA-LDP). This definition is inspired by the single agent DP introduced in Dwork et al. (2006) and LDP introduced in Kasiviswanathan et al.

(2011); Duchi et al. (2013). Throughout, we use the following notations. Let  $\mathbf{D}_h = (D_h^1, D_h^2, \dots, D_h^n)$  and  $\mathbf{D}'_h = (D_h^{1'}, D_h^{2'}, \dots, D_h^{n'})$  are the different datasets collected by the server at stage  $h$ . Here for each agent  $i \in N$ ,  $D_h^i$  and  $D_h^{i'}$  differs at exactly one component. Let  $\mathbf{D}_{1:h-1}$  be the information collected from stage 1 to stage  $h$ , i.e.,  $\mathbf{D}_{1:h-1} = (\mathbf{D}_1, \mathbf{D}_2, \dots, \mathbf{D}_{h-1})$ .

In online RL, each episode  $k \in [K]$  is viewed as the trajectory associated with a specific user. Thus, there are  $K$  users and  $n$  agents. Note that in MA-LDP the agent is different from the user. So, for MA-LDP we guarantee that for any user the information of any agent  $i \in N$  send to the server is privatized. Therefore, the server is agnostic to the sensitive data.

**Definition 1** (MA-LDP). *For any  $\epsilon \geq 0$ , and  $\delta \geq 0$ , a randomized mechanism  $\mathcal{M}$  preserves  $(\epsilon, \delta)$  MA-LDP if for any two users  $u$  and  $u'$  and their corresponding data  $\mathbf{D}_u = (D_u^1, D_u^2, \dots, D_u^n) \in \mathcal{U}$  and  $\mathbf{D}_{u'} = (D_{u'}^1, D_{u'}^2, \dots, D_{u'}^n) \in \mathcal{U}$ , it satisfies  $\mathbb{P}(\mathcal{M}(\mathbf{D}_u) \in U) \leq e^\epsilon \mathbb{P}(\mathcal{M}(\mathbf{D}_{u'}) \in U) + \delta$ ,  $U \in \mathcal{U}$ . Here for each agent  $i \in N$ , the  $D_u^i$  and  $D_{u'}^i$  differs at exactly one component.*

### 3 MA-LDP algorithm

In this Section, we introduce the multi-agent locally differential private (MA-LDP) algorithm that achieves the sub-linear regret. Our MA-LDP algorithm is inspired from the UCRL-VTR algorithm of Jia et al. (2020) and also use some structure of UCRL-VTR-LDP algorithm for the single agent RL in Liao et al. (2021). Let  $K$  be the number of episodes, and each episode consists of a fixed planning horizon  $H$ . Initially, for each agent  $i \in N$ , the estimate of the global state value function  $V^i$  and the global state-action value function  $Q^i$  is taken as 0 at the  $H + 1$ -th stage. At every stage  $h$ , the optimistic estimator of the state action value function for each agent  $i \in N$  is obtained via the backward induction algorithm.

For every stage  $h \in [H]$  and agent  $i \in N$ , we initially receive  $\Lambda_{1,h}^i = \Sigma_{1,h}^i = \lambda \mathbf{I}$  and  $\hat{\theta}_{1,h}^i = \mathbf{0}_{nd}$  information from the server for the local user  $k = 1$ . For the local user  $k$  and the received information  $\Lambda_{k,h}^i, u_{k,h}^i$ , each agent  $i \in N$  uses the backward induction algorithm along with an additional UCB bonus term to get the optimistic estimator of the optimal state-action value function. The update is  $Q_{k,h}^i(\cdot, \cdot) \leftarrow \min\{H + 1 - h, \beta_{k,h}\} \|\Sigma_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(\cdot, \cdot)\|_2 + \bar{r}_h(\cdot, \cdot; \mathbf{w}_{k,h}^i) + \langle \hat{\theta}_{k,h}^i, \phi_{V_{k,h+1}^i}(\cdot, \cdot) \rangle$  here  $\beta_{k,h}$  is identified for each noise mechanism separately. We use  $\beta_{k,h}^G, \beta_{k,h}^L, \beta_{k,h}^{BU}, \beta_{k,h}^{BL}$  for each of the noise mechanisms respectively. These are identified in the next section.

In episode  $k$ , each agent  $i \in N$  realizes the initial global state  $s_{k,1}$  and then for all stages  $h \in [H]$  it takes action using the current optimistic estimator of the global state-action value  $Q_{k,h}^i$ . In particular, the action taken by agent  $i \in N$  is according to the max min criteria that captures its best action against the worst possible action by other agents (lines 11-17). Moreover, each agent maintains an intermediate reward function parameters for each stage  $h \in [H]$ . This intermediate reward function parameter is used to compute the reward function parameters for the next episode.

Once the action is taken, a new state is realized according to the unknown distribution  $\mathbb{P}_h(\cdot | s_h, \mathbf{a}_h)$  (line 21). To get the estimate of the true transition probability parameters each agent requires to send some information to the server. In MA-LDP algorithm data that is shared with the server from each agent is obtained via the ridge regression based minimization of the transition probability parameters. Thus, the server requires the true information  $\Delta \tilde{\Lambda}_{k,h}^i = \phi_{V_{k,h+1}^i}(s_{k,h}, \mathbf{a}_{k,h}) \phi_{V_{k,h+1}^i}(s_{k,h}, \mathbf{a}_{k,h})^\top$  and  $\Delta \tilde{u}_{k,h}^i = \phi_{V_{k,h+1}^i}(s_{k,h}, \mathbf{a}_{k,h}) V_{k,h+1}^i(s_{k,h+1})$  for each agent  $i \in N$  and stage  $h$  in episode  $k$ . However, to preserve the privacy, we privatize the above true information using different noise adding mechanisms. In particular, to the true information  $\Delta \tilde{\Lambda}_{k,h}^i$  we add a matrix  $\mathbf{W}_{k,h}^i$ , where each entry of the matrix  $\mathbf{W}_{k,h}^i$  is drawn according to a distribution corresponding the noise-adding mechanism used. Moreover, to the true information  $\Delta \tilde{u}_{k,h}^i$ , we add  $\xi_{k,h}^i$ . Again each entry of  $\xi_{k,h}^i$  is drawn from different distributions corresponding to different noise adding mechanisms. Let  $\Delta \Lambda_{k,h}^i$  and  $\Delta u_{k,h}^i$  be the anonymized information that is shared to the server (lines 26-32), i.e.,  $\Delta \Lambda_{k,h}^i \leftarrow \Delta \tilde{\Lambda}_{k,h}^i + \mathbf{W}_{k,h}^i$ ;  $\Delta u_{k,h}^i \leftarrow \Delta \tilde{u}_{k,h}^i + \xi_{k,h}^i$ . Server on the other hand collects this information and use it to updates  $\Sigma_{k,h}^i$  and  $u_{k,h}^i$  and hence gives the next estimate of the transition probability parameters. However, merely adding the noise might not preserve the PSD property, so we shift this matrix by adding an  $\eta \mathbf{I}$  to guarantee the PSD property. Finally, each agent sends  $\Lambda_{k,h}^i$  and  $u_{k,h}^i$  to the  $k + 1$ -th user. Apart from executing a policy that uses an optimistic estimator, each agent  $i \in N$  also updates the  $\Sigma_{k,h}$  and  $u_{k,h}$ . These updates are used to estimate the true model parameters  $\hat{\theta}_{k+1,h}^i$ . They are inspired from the minimizer to a regularized linear regression problem similar to Zhou et al. (2021). Lastly, to ensure the privacy of the reward function, the parameters  $\mathbf{w}^i$  of the reward functions are updated according to the consensus matrix. It is important to note that we use the stochastic

approximation based rule to update the reward function parameters  $\mathbf{w}^i$ . We want to emphasize that our algorithm uses the most recent reward function parameters. These parameters are updated in each episode, and stage of the episode via the consensus matrix in line 34 of the MA-LDP algorithm. The convergence of the reward function parameters ensures the convergence of state-action value function and these are used in regret analysis of MA-LDP algorithm.

---

**Algorithm 1** MA-Gaussian/Laplace-LDP
 

---

```

1: Require: Privacy parameters  $\epsilon, \delta$ ; failure probability  $\alpha$ ; parameter  $\eta, \mathbf{w}_{0,0}^i = 0$  for all  $i \in N$ 
2: Set  $\sigma = 4H^3 \sqrt{2 \log(2.5H/\delta)}/\epsilon$  for Gaussian; Set  $b = 4H^3 \sqrt{nd}/\epsilon$  for Laplace
3: for user  $k = 1, \dots, K$  do
4:   for  $i = 1, 2, \dots, n$  do
5:     For local user  $k$ :
6:       Receive  $\{\Sigma_{k,1}^i, \dots, \Sigma_{k,H}^i, \hat{\boldsymbol{\theta}}_{k,1}^i, \dots, \hat{\boldsymbol{\theta}}_{k,H}^i\}$ 
7:     end for
8:   for  $i = 1, \dots, n$  do
9:     for  $h = H, \dots, 1$  do
10:       $Q_{k,h}^i(\cdot, \cdot) \leftarrow \min \{H + 1 - h, \bar{r}_h(\cdot, \cdot; \mathbf{w}_{k,h}^i) + \langle \hat{\boldsymbol{\theta}}_{k,h}^i, \phi_{V_{k,h+1}^i}(\cdot, \cdot) \rangle + \beta_{k,h} \|\Sigma_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(\cdot, \cdot)\|_2\}$ 
11:       $V_{k,h}^i \leftarrow \max_{a^i \in \mathcal{A}^i} Q_{k,h}^i(\cdot, a^i, \mathbf{a}_{k,h}^{-i})$ 
12:    end for
13:  end for
14:  Receive the initial state  $s_{k,1}$ 
15:  for  $h = 1, \dots, H$  do
16:    for  $i = 1, \dots, n$  do
17:      Take action  $\mathbf{a}_{k,h}^i \leftarrow \arg \max_{a \in \mathcal{A}^i} \min_{a^{-i} \in \mathcal{A}^{-i}} Q_{k,h}^i(s_{k,h}, \mathbf{a}, \mathbf{a}^{-i})$ 
18:      Set  $\tilde{\mathbf{w}}_{k,h}^i \leftarrow \mathbf{w}_{k,h}^i + \gamma_{k,h} \cdot [r_h^i(\cdot, \cdot) - \bar{r}_h(\cdot, \cdot; \mathbf{w}_{k,h}^i)] \cdot \nabla_{\mathbf{w}} \bar{r}(\cdot, \cdot; \mathbf{w}_{k,h}^i)$ 
19:    end for
20:    Set  $\mathbf{a}_{k,h} = (\mathbf{a}_{k,h}^1, \mathbf{a}_{k,h}^2, \dots, \mathbf{a}_{k,h}^n)$ 
21:    Observe the next state  $s_{k,h+1}$ 
22:    for  $i = 1, \dots, n$  do
23:       $\Delta \tilde{\Lambda}_{k,h}^i = \phi_{V_{k,h+1}^i}(s_{k,h}, \mathbf{a}_{k,h}) \phi_{V_{k,h+1}^i}(s_{k,h}, \mathbf{a}_{k,h})^\top$ 
24:       $\Delta \tilde{u}_{k,h}^i = \phi_{V_{k,h+1}^i}(s_{k,h}, \mathbf{a}_{k,h}) V_{k,h+1}^i(s_{k,h+1})$ 
25:      Set  $\Delta \Lambda_{k,h}^i \leftarrow \Delta \tilde{\Lambda}_{k,h}^i + \mathbf{W}_{k,h}^i$ 
26:      Set  $\Delta u_{k,h}^i \leftarrow \Delta \tilde{u}_{k,h}^i + \boldsymbol{\xi}_{k,h}^i$ 
27:    end for
28:  end for
29:  for  $i = 1, \dots, n$  do
30:    Set  $D_k^i = \{\Delta \Lambda_{k,1}^i, \dots, \Delta \Lambda_{k,H}^i, \Delta u_{k,1}^i, \dots, \Delta u_{k,H}^i\}$ 
31:  end for
32:  Send  $\mathbf{D}_k = (D_k^1, D_k^2, \dots, D_k^n)$  to the server
33:  For the server:
34:  for  $h = 1, \dots, H$  do
35:    for  $i = 1, \dots, n$  do
36:       $\Lambda_{k+1,h}^i \leftarrow \Lambda_{k,h}^i + \Delta \Lambda_{k,h}^i$ 
37:       $u_{k+1,h}^i \leftarrow u_{k,h}^i + \Delta u_{k,h}^i$ 
38:       $\Sigma_{k+1,h}^i \leftarrow \Lambda_{k+1,h}^i + \eta I$ 
39:       $\hat{\boldsymbol{\theta}}_{k+1,h}^i \leftarrow (\Sigma_{k+1,h}^i)^{-1} u_{k+1,h}^i$ 
40:    end for
41:  end for
42:  Send  $\{\Sigma_{k+1,1}^i, \dots, \Sigma_{k+1,H}^i, \hat{\boldsymbol{\theta}}_{k+1,1}^i, \dots, \hat{\boldsymbol{\theta}}_{k+1,H}^i\}$  to the user  $k + 1$ 
43:  Update  $\mathbf{w}_{k+1,h}^i = \sum_{j \in N} l_{k,h}(i, j) \tilde{\mathbf{w}}_{k,h}^j, \forall h \in [H]$ 
44: end for
    
```

---

## 4 Preliminary Results

In this Section, we state the results and the definition that are useful in regret analysis and are common for all the noise adding mechanisms. To this end, we first define the notion of privacy loss.

**Definition 2** (Privacy loss [Dwork et al. \(2006\)](#)). *For any neighboring datasets  $d, d'$ , a mechanism  $\mathcal{M}$ , auxiliary input  $\mathbf{aux}$ , and an outcome  $o \in \mathbb{R}$ , the privacy loss at outcome  $o$  is defined as  $c(o; \mathcal{M}, \mathbf{aux}, d, d') := \log \frac{\mathbb{P}(\mathcal{M}(\mathbf{aux}, d) = o)}{\mathbb{P}(\mathcal{M}(\mathbf{aux}, d') = o)}$ .*

The privacy loss defined above represents the loss incurred by the outcome  $o$ , when the data  $d$  is replaced by  $d'$ . Note that this loss might be positive or negative. The  $(\epsilon, \delta)$  DP ensures that for neighbouring data  $d, d'$  the absolute value of the above privacy loss is bounded by  $\epsilon$  with probability at least  $1 - \delta$ . When  $\delta = 0$  we get  $(\epsilon, 0)$  DP, that is, the absolute privacy loss is at most  $\epsilon$  a.s. The  $(0, \delta)$  DP implies that the absolute privacy loss is 0 with confidence  $1 - \delta$ .

In our decentralized multi-agent LDP setting, the proof of the regret bounds for each of the noise adding mechanism relies on the asymptotic convergence of reward function parameters. In particular, we show that the reward function parameters converge independent of agent as  $k$  (number of episodes) goes to infinity. To this end, let  $d(s)$  be the stationary distribution of the Markov chain  $\{s_t\}_{t \geq 0}$  under policy  $\pi$ , and  $\pi(s, \mathbf{a})$  be the probability of taking action  $\mathbf{a}$  in state  $s$ . Moreover, let  $D^{s, \mathbf{a}} = \text{diag}[d(s)\pi(s, \mathbf{a})]$  be the diagonal matrix with  $d(s)\pi(s, \mathbf{a})$  as diagonal entries.

**Lemma 1.** *Under assumptions 2 and 3, for the sequence  $\{\mathbf{w}_{k, h}^i\}$ , we have  $\lim_k \mathbf{w}_{k, h}^i = \mathbf{w}^*$  a.s. for each agent  $i \in N$  and for all  $h \in [H]$ , where  $\mathbf{w}^*$  is unique solution to  $\Psi^\top D^{s, \mathbf{a}} (\Psi \mathbf{w}^* - \bar{r}) = 0$ .*

We would like to emphasize that this convergence is specifically needed for the decentralized multi-agent setting, and hence novel to our work. The proof is deferred to the Appendix A and it follows on the same lines as in [Trivedi and Hemachandra \(2023, 2022\)](#); [Zhang et al. \(2018\)](#) that uses stochastic approximations methods of [Borkar \(2022\)](#). Next, we show that in our MA-LDP algorithm each agent uses an optimistic estimator of the state-action value function in each episode. To this end, we have the following lemma common to all the noise adding mechanisms except a term  $\beta_{k, h}$  (see line 10 of the MA-LDP algorithm) that is separately identified in the Appendix C. The proof is via induction over  $h$ , see Appendix B.

**Lemma 2.** *Let  $Q_{k, h}^i$  and  $V_{k, h}^i$  be the estimate of the global state-action value and global state value functions respectively by agent  $i \in N$ . Then, for any pairs  $(s, \mathbf{a}, k, h) \in \mathcal{S} \times \mathcal{A} \times [K] \times [H]$ , we have  $Q_h^{i*}(s, \mathbf{a}) \leq Q_{k, h}^i(s, \mathbf{a})$  and  $V_h^{*, i}(s) \leq V_{k, h}^i(s)$ .*

In our regret analysis we also require the following bound on the estimated model parameters.

**Lemma 3.** *If  $\eta = 1$ , then for any fixed policy  $\pi$  and all pairs  $(s, a, h, k) \in \mathcal{S} \times \mathcal{A} \times [H] \times [K]$ , with probability at least  $1 - \alpha/2$  for all  $i \in N$ , we have  $\|(\Sigma_{k, h}^i)^{1/2}(\hat{\boldsymbol{\theta}}_{k, h}^i - \boldsymbol{\theta}_h^*)\| \leq \beta_k$ .*

The proof of this Lemma uses a decomposition of  $\|(\Sigma_{k, h}^i)^{1/2}(\hat{\boldsymbol{\theta}}_{k, h}^i - \boldsymbol{\theta}_h^*)\|$  in three different terms, we call them  $\mathbf{q}_1, \mathbf{q}_2$ , and  $\mathbf{q}_3$ . We give this decomposition in Appendix C. Each of these terms are bounded differently for every noise mechanism that helps in identifying the respective  $\beta_k$ .

## 5 Gaussian and Laplace noise adding mechanisms (unbounded support)

In this Section, we consider two most popular noise adding mechanisms, Gaussian and Laplace. To our knowledge, the multi-agent version of these are not studied so far. We first show that our MA-LDP algorithm preserves  $(\epsilon, \delta)$  and  $(\epsilon, 0)$  LDP for Gaussian and Laplace mechanisms respectively. We also show that MA-LDP algorithm for these mechanisms achieves the sub-linear regret.

### 5.1 Gaussian Noise Adding Mechanism

For the Gaussian mechanism, we add  $\mathbf{W}_{k, h}^i$  to  $\Lambda_{k, h}^i$  in line (21) of MA-LDP algorithm. Recall,  $\mathbf{W}_{k, h}^i$  is a symmetric matrix, so for  $l \leq m$ , each entry  $(l, m)$  of  $\mathbf{W}_{k, h}^i$  is sampled from the Gaussian distribution  $\mathcal{N}(0, \sigma^2)$  distribution. Moreover, to  $u_{k, h}^i$  we add  $\boldsymbol{\xi}_{k, h}^i$  sampled from  $\mathcal{N}(\mathbf{0}_{nd}, \sigma^2 \mathbf{I}_{nd \times nd})$  distribution.

**Theorem 1.** *If we choose the parameter  $\sigma$  of the Gaussian distribution  $\mathcal{N}(0, \sigma^2)$  such that  $\sigma = 4H^3 \sqrt{2 \log(2.5H/\delta)}/\epsilon$ , then, the Gaussian mechanism  $\mathcal{M}_G$  will satisfy  $(\epsilon, \delta)$  MA-LDP property.*

The proof of this Theorem is referred to Appendix D and D.1. and use Theorem A.2 of Liao et al. (2021) (for details see Appendix G.2). Next, we show that the above mentioned Gaussian mechanism with privacy parameters  $\epsilon, \delta$  achieves a sub-linear regret.

**Theorem 2.** *Let  $\alpha \in (0, 1)$ ,  $\eta = 1$ , privacy parameter  $\epsilon > 0, \delta > 0$ , and  $\beta_{k,h}^G = c_g(nd)^{3/4}(H - h + 1)^{3/2}k^{1/4} \log(ndT/\alpha)((\log(H - h + 1)/\delta))^{1/4}\sqrt{1/\epsilon}$  with  $c_g$  being an absolute constant. Consider the Gaussian noise mechanism with parameter  $\sigma$  as in Theorem 1. Then, for any user  $k$ , with probability at least  $1 - \alpha$ , the total regret of MA-LDP algorithm in the first  $T = KH$  steps with BL noise mechanism is at most  $\tilde{O}(n^{5/4}d^{5/4}H^{7/4}T^{3/4} \log(ndT/\alpha)(\log(H/\delta))^{1/4}\sqrt{1/\epsilon})$ .*

*Proof.* (Sketch Only). Here we provide the outline of the proof. The proof involve four major steps: 1) to show that the model parameters obtained from the MA-LDP algorithm converges to the true model parameters  $\theta^*$  within the confidence radius  $\beta_k^g$  (Lemma 3 in previous Section). 2) Using induction argument to show that in each step MA-LDP algorithm uses an optimistic estimator of the state-action value function (Lemma 2 in the previous Section). 3) The convergence of the reward function parameters that is shown in the previous subsection in Lemma 1. Finally, 4) bounding the regret via concentration inequalities. The details of this step are given in Appendix E.1.  $\square$

## 5.2 Laplace Noise Adding Mechanism

Unlike Gaussian mechanism for the Laplace noise adding mechanisms, we add  $\mathbf{W}_{k,h}^i$  to  $\Lambda_{k,h}^i$  in line (21) of the MA-LDP algorithm. Recall,  $\mathbf{W}_{k,h}^i$  is a symmetric matrix, so for  $l \leq m$ , each entry  $(l, m)$  of  $\mathbf{W}_{k,h}^i$  is sampled from the Laplace distribution  $\mathcal{L}(b)$  distribution. Moreover, to each entry of  $u_{k,h}^i$  we add  $\xi_{k,h}^i$  sampled from  $\mathcal{L}(b; \cdot)$  distribution.

**Theorem 3.** *If we choose the parameter  $b$  of the Laplace distribution  $\mathcal{L}(b)$  such that  $b = \frac{4H^3\sqrt{nd}}{\epsilon}$ , then, the Laplace noise adding mechanism  $\mathcal{M}_L$  will satisfy  $(\epsilon, 0)$  MA-LDP property.*

The proof of this theorem is deferred to Appendix D and D.2. Next, we provide the upper bound on the regret incurred by MA-LDP algorithm with Laplace mechanism.

**Theorem 4.** *Let  $\alpha \in (0, 1)$ ,  $\eta = 1$ , privacy parameter  $\epsilon > 0$ , and  $\beta_{k,h}^L = c_l(nd)^{3/4}(H - h + 1)^{3/2}k^{1/4} \log(ndT/\alpha)\sqrt{1/\epsilon}$  with  $c_l$  being an absolute constant. Consider the Laplace noise mechanism with parameter  $b$  as in Theorem 3. Then, for any user  $k$ , with probability at least  $1 - \alpha$ , the total regret of MA-LDP algorithm in the first  $T = KH$  steps with Laplace noise mechanism is at most  $\tilde{O}(n^{5/4}d^{5/4}H^{7/4}T^{3/4} \log(ndT/\alpha)\sqrt{1/\epsilon})$ .*

## 6 Uniform and bounded Laplace noise adding mechanisms

In this Section, we consider the noise adding mechanism with the bounded support. The motivation comes from the fact that noise mechanisms with unbounded support allows all the noise values to attain the LDP with given privacy parameters  $\epsilon, \delta$ ; some of these noise values can be arbitrarily large. So, the central question is: does restricting the noise values to a bounded support change the privacy property or the regret or both? If so, how? If it doesn't affect the privacy property, then is it that adding a bounded noise attains the same regret with the same privacy parameters. In this later case, it means that we need not inject huge noises to the system. In the this section, we show that this is indeed possible, and for a certain relation between the distribution parameters and the support of the Laplace mechanism, our MA-LDP algorithm attains the same regret up to the constants; this is one of our main result. Moreover, we compare these regrets with the most natural choice of the uniform noise mechanism. We show that the privacy guarantees of the uniform and Laplace mechanisms with bounded support are two extremes,  $(0, \delta)$  and  $(\epsilon, 0)$  LDP respectively.

### 6.1 Uniform Noise Adding Mechanism

The uniform (U) mechanism with MA-LDP algorithm works as follows; for a given  $a > 0$  we denote by  $\mathcal{U}[-a, a]$  the uniform distribution over the interval  $[-a, a]$ . For the uniform mechanism, we add  $\mathbf{W}_{k,h}^i$  to  $\Lambda_{k,h}^i$  in line (21) of the MA-LDP algorithm. Recall,  $\mathbf{W}_{k,h}^i$  is a symmetric matrix, so for  $l \leq m$ , each entry  $(l, m)$  of  $\mathbf{W}_{k,h}^i$  is sampled from the uniform distribution  $\mathcal{U}[-a, a]$  distribution. Moreover, to  $u_{k,h}^i$  we add  $\xi_{k,h}^i$  whose each entry is sampled from  $\mathcal{U}[-a, a]$  distribution.

**Theorem 5.** *If we choose the parameter  $a$  of the uniform distribution  $\mathcal{U}[-a, a]$  such that  $a = 4H^3\sqrt{\log(2H/\delta)}$ , then MA-LDP algorithm with mechanism  $\mathcal{M}_U$  as preserves  $(0, \delta)$  LDP property.*

The proof of this theorem is deferred to the Appendix D and D.3 The regret of MA-LDP algorithm with BU mechanism is given in the following theorem. The proof is given in Appendix E and E.3

**Theorem 6.** Let  $\alpha \in (0, 1)$ ,  $\eta = 1$ , privacy parameter  $\delta > 0$ , and  $\beta_{k,h}^U = c_u(nd)^{3/4}(H - h + 1)^{3/2}k^{1/4} \log(ndT/\alpha)(\log(H - h + 1/\delta))^{1/4}$ , with  $c_u$  being an absolute constant. Consider the uniform noise mechanism with parameter  $a$  as in Theorem 5. Then, for any user  $k$ , with probability at least  $1 - \alpha$ , the total regret of MA-LDP algorithm in the first  $T = KH$  steps with uniform noise mechanism is at most  $\tilde{O}(n^{5/4}d^{5/4}H^{7/4}T^{3/4} \log(ndT/\alpha)(\log(H/\delta))^{1/4})$ .

## 6.2 Bounded Laplace Noise Adding Mechanism

In this Section, we design another noise adding mechanism with bounded support. We call it bounded Laplace (BL) mechanism. The BL mechanism is obtained by restricting the support of the Laplace distribution to  $[-B, B]$  for a given  $B > 0$ . The probability density of the bounded Laplace distribution,  $\mathcal{BL}(x; b)$ , with parameter  $b$  is  $f_{\mathcal{BL}}(x; b) = (2b(1 - \exp(-\frac{B}{b}))^{-1} \exp(-\frac{|x|}{b}))^{-1}$ ,  $\forall x \in [-B, B]$ , and 0 otherwise. The variance of this distribution is  $\zeta = 2b^2(1 - \exp(-\frac{B}{b}))^{-1} - \kappa$  where  $\kappa = ((B + b)^2 + b^2) \times \exp(-\frac{B}{b}) \times (1 - \exp(-\frac{B}{b}))^{-1}$ . For the BL mechanism, we add  $\mathbf{W}_{k,h}^i$  to  $\Lambda_{k,h}^i$  in line (21) of the MA-LDP algorithm. Recall,  $\mathbf{W}_{k,h}^i$  is a symmetric matrix, so for  $l \leq m$ , each entry  $(l, m)$  of  $\mathbf{W}_{k,h}^i$  is sampled from  $\mathcal{BL}(b; B)$  distribution. Moreover, to  $u_{k,h}^i$  we add  $\xi_{k,h}^i$  whose each entry is sampled from  $\mathcal{BL}(b, B)$  distribution.

**Theorem 7.** If we choose the parameter  $b$  of the bounded Laplace distribution  $\mathcal{BL}(b; B)$  such that  $b = \frac{4H^3\sqrt{nd}}{\epsilon}$ , then, the BL mechanism  $\mathcal{M}_{BL}$  will satisfy  $(\epsilon, 0)$  MA-LDP property.

The proof of the above theorem is given to the Appendix D and D.4 Thus, the BL mechanism preserves  $(\epsilon, 0)$  LDP similar to the unbounded case. The following theorem bounds the regret of MA-LDP algorithm with BL mechanism. The proof of is deferred to Appendix E and E.4.

**Theorem 8.** Let  $\alpha \in (0, 1)$ ,  $\eta = 1$ , privacy parameter  $\epsilon > 0$ , and  $\beta_{k,h}^{BL} = c_{bl}(nd)^{3/4}(H - h + 1)^{3/2}k^{1/4} \log(ndT/\alpha) \sqrt{1/\epsilon}$ , with  $c_{bl}$  being an absolute constant. Consider the BL noise mechanism with parameter  $b$  as in Theorem 7. Then, for any user  $k$ , with probability at least  $1 - \alpha$ , the total regret of MA-LDP algorithm in the first  $T = KH$  steps with BL noise mechanism is at most  $\tilde{O}(n^{5/4}d^{5/4}\zeta^{1/4}H^{1/4}T^{3/4} \log(ndT/\alpha))$

Mechanism	Type	Privacy	Order of Regret
Gaussian	Unbounded	$(\epsilon, \delta)$	$\tilde{O}((nd)^{5/4}H^{7/4}T^{3/4} \log(ndT/\alpha)(\log(H/\delta))^{1/4} \sqrt{1/\epsilon})$
Laplace	Unbounded	$(\epsilon, 0)$	$\tilde{O}((nd)^{5/4}H^{7/4}T^{3/4} \log(ndT/\alpha) \sqrt{1/\epsilon})$
Uniform	Bounded	$(0, \delta)$	$\tilde{O}((nd)^{5/4}H^{7/4}T^{3/4} \log(ndT/\alpha)(\log(H/\delta))^{1/4})$
Bounded Laplace	Bounded	$(\epsilon, 0)$	$\tilde{O}((nd)^{5/4}\zeta^{1/4}H^{1/4}T^{3/4} \log(ndT/\alpha))$

Table 1: Privacy guarantees and the order of regret for different noise adding mechanisms.

## 7 Comparison of regret for different noise mechanisms

First, from the regret expressions of Gaussian and Laplace mechanism in Theorems 2 and 4 it follows that  $R_K(\epsilon_1)/R_K(\epsilon_2) = \sqrt{\epsilon_2/\epsilon_1}$ . Thus, we have:

**Theorem 9.** If privacy parameters  $\epsilon_1$  and  $\epsilon_2$  are such that  $\epsilon_1 > \epsilon_2$ . Then for both the Gaussian and Laplace mechanisms we have that  $R_K(\epsilon_1) < R_K(\epsilon_2)$ .

Moreover, we have the following Theorem for the cumulative regret between the Gaussian and Laplace mechanism for the same privacy parameter  $\epsilon$ .

**Theorem 10.** Let  $R_K^G(\epsilon), R_K^L(\epsilon)$  be the cumulative regret of the Gaussian and Laplace mechanism respectively with privacy parameters  $\epsilon, \delta$ , and  $H > 2$ . Then,  $R_K^G(\epsilon) > R_K^L(\epsilon)$ .

The proof follows from the fact that the Gaussian noise mechanism gives the approximate LDP, i.e.,  $\delta > 0$ , so  $\log(H/\delta) > 1$  for all  $\delta \in (0, 1)$  and  $H > 2$ . Therefore,  $R_K^G(\epsilon)/R_K^L(\epsilon) = \log(H/\delta) > 1$ .

Apart from the above, recall  $\zeta$ , the variance of the bounded Laplace distribution and is a function of  $B$  and  $b$ . So, depending on whether  $B$  has the same order as that of  $b$  or not, we get different expressions of  $\zeta$ , thus different



order of the regret (Theorem 8). So, we compare the regret of the BL mechanism with the Laplace mechanism when  $B = O(b^\gamma)$  for different possible values of  $\gamma$ . For the case of  $\gamma > 1$ , we have that  $R_K^{BL}/R_K^L = (H^3/\epsilon)^{\frac{2}{3}}$ . If

$B$	$R_K^{BL}$
$O(b^\gamma), 0 \leq \gamma \leq 1$	$\tilde{O}((nd)^{5/4}H^{7/4}T^{3/4}\log(ndT/\alpha))\sqrt{1/\epsilon}$
$O(b^\gamma), \gamma > 1$	$\tilde{O}((nd)^{5/4}H^{7/4}H^{3\gamma/2}T^{3/4}\log(ndT/\alpha))\sqrt{1/\epsilon^{\gamma+1}}$

Table 2: Regret bound for the Bounded Laplace (BL) mechanism. MA-LDP algorithm with BL mechanism offers the same order of regret as that of the Laplace mechanism when  $B = O(b^\gamma)$  for  $\gamma \in [0, 1]$ .

$(H^3/\epsilon)^{\frac{2}{3}} > 1$  and  $\gamma > 1$ , then the regret of MA-LDP with BL mechanism is more than that of Laplace mechanism. So, for the given problem instance, the regret of the BL mechanism will be the same as that of the Laplace mechanism if  $\gamma \in [0, 1]$ . Further, if  $\gamma > 1$  and  $(H^3/\epsilon)^{\frac{2}{3}} < 1$ , then BL will have lower regret than that of Laplace. Though the BL mechanism injects noise from a bounded support, the regret of MA-LDP algorithm with BL mechanism is either on par or lower with that of the Laplace mechanism in most of the cases, i.e., when  $B = O(b^\gamma)$  with  $\gamma \in [0, 1]$ . In a very restrictive setting where  $\gamma > 1$  and  $(H^3/\epsilon)^{\frac{2}{3}} > 1$ , the regret from BL mechanism is more than that of the Laplace mechanism.

**Theorem 11.** *If the end point  $B$  of the support of bounded Laplace,  $BL$ , distribution (with parameter  $b$ ) is of order  $O(b^\gamma)$ , where  $\gamma \leq 1$ , then the order of the regret of the MA-LDP algorithm with BL mechanism is the same as that of the Laplace mechanism.*

## 8 Computational Experiments

In this section, we give computational results to validate the usefulness of our MA-LDP algorithm. To this end, we consider a network with  $q + 2$  nodes shown in Figure 1 below,  $q \geq 1$ , i.e.,  $\{s_{in}, 1, 2, \dots, q, g\}$ , where  $s_{in}, g$  are the initial and goal nodes respectively. The number of global states are  $(q + 2)^n$ .

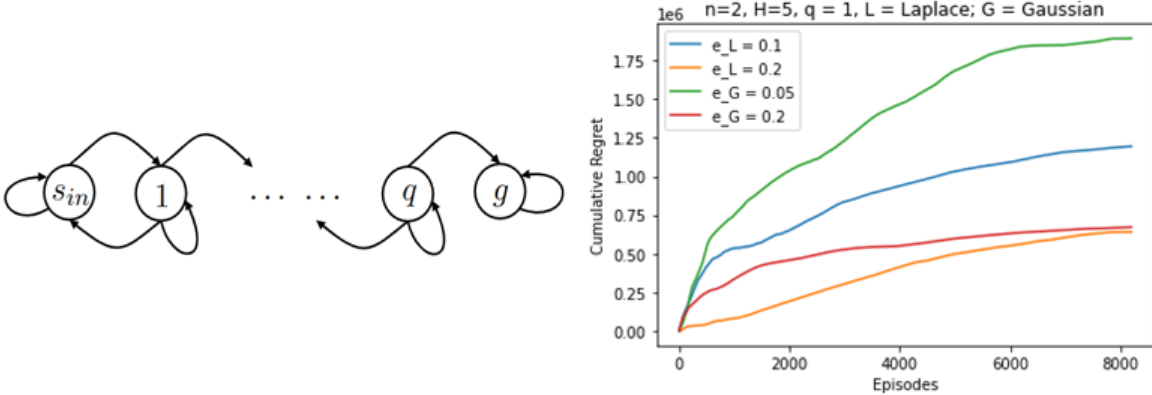


Figure 1: (Left) The MDP problem instance that we consider. (Right) Cumulative regret with number of episodes for the Laplace and Gaussian mechanism.

In each state the actions available to each agent are  $\mathcal{A}^i = \{-1, 1\}^{d-1}$ , for given  $d \geq 2$ . So, the total number of actions is  $2^{n(d-1)}$ . Each agent  $i \in N$  receives a reward of 5/1000 units for taking any action in  $s_{in}$ , a reward of 1000 for action in  $g$ , and the reward of 0 unit for action in other nodes. The collective objective of the agents is to reach the goal node in a decentralized way while maximizing the overall reward. To address humongous state, action space, we parameterize the transition probability as  $\mathbb{P}_\theta(s'|s, \mathbf{a}) = \langle \phi(s'|s, \mathbf{a}), \boldsymbol{\theta}(s) \rangle$  for each  $(s', \mathbf{a}, s) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$ . The feature design for this transition probability is novel to our work and is given below. Let  $S(s)$  is the set all feasible states from state  $s$ . We define the global features as:  $\phi(s'|s, \mathbf{a}) = (\phi(s'|s^1, \mathbf{a}^1), \dots, \phi(s'|s^n, \mathbf{a}^n))$ , if  $s \neq g$ ,  $s' \in S(s)$ ;  $\mathbf{0}_{nd}$ , if  $s \neq g, s' \notin S(s)$ ,  $\mathbf{0}_{nd}$ , if  $s = g, s' \neq g$ , and  $(\mathbf{0}_{nd-1}, \alpha(s))$ , if  $s = g, s' = g$ . We identify  $\alpha(s)$  as  $\alpha(s) = \frac{|S(s)|}{n} \left\{ \frac{x_0}{2} + x_{q+1} + \sum_{j=1}^q \frac{x_j}{3} \right\}$ . Here  $x_0, x_1, \dots, x_q, x_{q+1}$  are the number of agents at the nodes  $s_{in}, 1, \dots, q, g$  respectively in the state  $s$ . The local features  $\phi(s'^i | s^i, \mathbf{a}^i)$  are defined in Appendix F. Moreover, the

transition probability parameters for any state  $s$  are taken as  $\theta(s) = \left(\theta^1, \frac{1}{\alpha(s)}, \theta^2, \frac{1}{\alpha(s)}, \dots, \theta^n, \frac{1}{\alpha(s)}\right)$  where  $\theta^i \in \left\{-\frac{\Delta}{n(d-1)}, \frac{\Delta}{n(d-1)}\right\}^{d-1}$ , and  $\Delta < \delta$ . More details of the experiments and the MDP involved are given in Appendix F.1. The result below shows that the above choice of feature design  $\phi(s'|s, \mathbf{a})$  and model parameters  $\theta(s)$  yields a valid MDP.

**Lemma 4.** *For every  $\theta(s)$ , features  $\phi(s'|s, \mathbf{a})$  satisfies the following: (a)  $\sum_{s'} \langle \phi(s'|s, \mathbf{a}), \theta(s) \rangle = 1, \forall s, \mathbf{a}$ ; (b)  $\langle \phi(s' = \mathbf{g}|s = \mathbf{g}, \mathbf{a}), \theta(s) \rangle = 1, \forall \mathbf{a}$ ; (c)  $\langle \phi(s' \neq \mathbf{g}|s = \mathbf{g}, \mathbf{a}), \theta(s) \rangle = 0, \forall \mathbf{a}$ .*

The proof of this Lemma is deferred to the Appendix F.1. We implement our MA-LDP algorithm with various privacy parameters on the Gaussian and Laplace mechanisms with unbounded support. Figure 1 shows the cumulative regret with  $n = 2$ , number of nodes in the network as 3, planning horizon  $H = 5$ . All the values are averaged over 10 runs; each run has  $K = 8500$  episodes. Here are some observations from the experiments. Firstly, for a given mechanism, the cumulative regret with lower privacy losses are higher than the higher privacy losses. That is  $R_K(\epsilon_L = 0.1) > R_K(\epsilon_L = 0.2)$ ; and  $R_K(\epsilon_G = 0.05) > R_K(\epsilon_G = 0.2)$ . This illustrates the results of Theorems 9 and 10 that compare cumulative regrets of various privacy losses.

## 9 Related Work

The idea of DP is first introduced in Dwork et al. (2006). DP is motivated by the designing the algorithms that preserve the user’s sensitive data. It indicates that changing or removing a data point has little influence on any observable output. However, DP has risk of data leakage, and is vulnerable to membership influence attacks Shokri et al. (2017). Hence a stronger notion ‘Locally DP’ is introduced Kasiviswanathan et al. (2011); Duchi et al. (2013). Under LDP, users send privatized data to the server and each individual user maintains its own sensitive data. The server, on the other hand, is totally agnostic about the sensitive data. The LDP problem in the multi-agent distributed optimization is studied in Dobbe et al. (2018). The amplification of the privacy in machine learning is studied in Cyffers and Bellet (2022). The optimal noise adding mechanism are introduced in Geng and Viswanath (2015). Apart from the learning and decision making the notion of DP is also common in other fields of science including Hu and Fang (2022); Duchi et al. (2013). In particular, Hu and Fang (2022) study the  $K$  armed bandit problem with distributionally trust model of the DP that guarantees the privacy without trustworthy server. Recently, Jia et al. (2020); Jin et al. (2020) propose single agent RL scheme with linear function approximation of the transition probability function. Recently, Liao et al. (2021) use the UCRL-VTR algorithm of Jia et al. (2020) and incorporate the notion of DP into it. However, these consider a single agent taking the decisions in the finite horizon models. In this work, we introduce MA-LDP algorithm in a fully decentralized MARL framework that use different noise mechanisms and identify conditions when bounded support mechanism attains regret that is comparable to that is offered by conventional mechanisms.

## 10 Discussion

In this work, we consider the notion of the LDP in the fully decentralized MARL setting. We first define LDP for MARL and then propose a generic MA-LDP algorithm which can handle any noise adding mechanism. We show that the MA-LDP algorithm preserves the privacy for four different noise adding mechanisms, then prove that it also achieves the sub-linear regret. Next, we compare the noise mechanisms with bounded support with that of unbounded support. Our key observation is that if the support of bounded noise distribution is picked appropriately, the regret is lower than the unbounded support noise mechanism. Thus, injecting a bounded noise is often sufficient for LDP without substantially affecting the nature of the regret. We illustrate our results on a networked MDP with many states and actions.

The work we consider offers a rich set of further possibilities. We mention some of these here. Firstly, our regret bound is super-linear in the number of agents and feature dimension; towards this, a nice update rule for the optimistic estimators of the state-action value function can be attempted. Secondly, to our knowledge, the regret bound we show are of its first kind, so an attempt to get a better sub-linear regret bound is possible. Moreover, a matching lower bound can also be tried. A careful study of bounded support noise mechanism that leads to the lower regret bounds with low noise values would be interesting.

## References

Abadi, M., Chu, A., Goodfellow, I., McMahan, H. B., Mironov, I., Talwar, K., and Zhang, L. (2016). Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pages 308–318.

- Abbasi-Yadkori, Y., Pál, D., and Szepesvári, C. (2011). Improved algorithms for linear stochastic bandits. *Advances in neural information processing systems*, 24.
- Borkar, V. S. (2022). *Stochastic approximation: a dynamical systems viewpoint. Second Edition*, volume 48. Springer.
- Cyffers, E. and Bellet, A. (2022). Privacy amplification by decentralization. In *International Conference on Artificial Intelligence and Statistics*, pages 5334–5353. PMLR.
- Dobbe, R., Pu, Y., Zhu, J., Ramchandran, K., and Tomlin, C. (2018). Customized local differential privacy for multi-agent distributed optimization. *arXiv preprint arXiv:1806.06035*.
- Duchi, J. C., Jordan, M. I., and Wainwright, M. J. (2013). Local privacy, data processing inequalities, and statistical minimax rates. *arXiv preprint arXiv:1302.3203*.
- Dwork, C., McSherry, F., Nissim, K., and Smith, A. (2006). Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, pages 265–284. Springer.
- Dwork, C. and Roth, A. (2014). The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407.
- Geng, Q. and Viswanath, P. (2015). Optimal noise adding mechanisms for approximate differential privacy. *IEEE Transactions on Information Theory*, 62(2):952–969.
- Hu, W. and Fang, H. (2022). Decentralized matrix factorization with heterogeneous differential privacy. *arXiv preprint arXiv:2212.00306*.
- Jia, Z., Yang, L., Szepesvari, C., and Wang, M. (2020). Model-based reinforcement learning with value-targeted regression. In *Learning for Dynamics and Control*, pages 666–686. PMLR.
- Jin, C., Yang, Z., Wang, Z., and Jordan, M. I. (2020). Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143. PMLR.
- Kasiviswanathan, S. P., Lee, H. K., Nissim, K., Raskhodnikova, S., and Smith, A. (2011). What can we learn privately? *SIAM Journal on Computing*, 40(3):793–826.
- Kushner, H. and Yin, G. G. (2003). *Stochastic approximation and recursive algorithms and applications*, volume 35. Springer Science & Business Media.
- Liao, C., He, J., and Gu, Q. (2021). Locally differentially private reinforcement learning for linear mixture markov decision processes. *arXiv preprint arXiv:2110.10133*.
- Metivier, M. and Priouret, P. (1984). Applications of a Kushner and Clark lemma to general classes of stochastic algorithms. *IEEE Transactions on Information Theory*, 30(2):140–151.
- Min, Y., He, J., Wang, T., and Gu, Q. (2022). Learning stochastic shortest path with linear function approximation. In *International Conference on Machine Learning*, pages 15584–15629. PMLR.
- Ross, S. M. (2022). *Simulation*. Academic Press.
- Shokri, R., Stronati, M., Song, C., and Shmatikov, V. (2017). Membership inference attacks against machine learning models. In *2017 IEEE symposium on security and privacy (SP)*, pages 3–18. IEEE.
- Tao, T. (2012). *Topics in random matrix theory*, volume 132. American Mathematical Soc.
- Trivedi, P. and Hemachandra, N. (2022). Multi-agent natural actor-critic reinforcement learning algorithms. *Dynamic Games and Applications*, pages 1–31.
- Trivedi, P. and Hemachandra, N. (2023). Multi-agent congestion cost minimization with linear function approximation. *arXiv preprint arXiv:2301.10993*.
- Vial, D., Parulekar, A., Shakkottai, S., and Srikant, R. (2022). Regret bounds for stochastic shortest path problems with linear function approximation. In *International Conference on Machine Learning*, pages 22203–22233. PMLR.
- Zhang, K., Yang, Z., Liu, H., Zhang, T., and Basar, T. (2018). Fully decentralized multi-agent reinforcement learning with networked agents. In *International Conference on Machine Learning*, pages 5872–5881. PMLR.
- Zhou, D., Gu, Q., and Szepesvari, C. (2021). Nearly minimax optimal reinforcement learning for linear mixture markov decision processes. In *Conference on Learning Theory*, pages 4532–4576. PMLR.

## A Proofs of convergence of reward function parameters (Lemma 1)

*Proof.* Let  $t = kh$ , therefore, as  $k \rightarrow \infty$  we have  $t \rightarrow \infty$ . The proof of this result is on the same lines to [Zhang et al. \(2018\)](#); [Trivedi and Hemachandra \(2023, 2022\)](#). We briefly give the proof details here.

To prove the convergence of the reward function parameters, we use the following Proposition to give bounds on  $\mathbf{w}_t^i$  for all  $i \in N$ . For proof, we refer to [Zhang et al. \(2018\)](#).

**Proposition 1.** *Under assumptions 2, and 3 the sequence  $\{\mathbf{w}_t^i\}$  satisfy  $\sup_t \|\mathbf{w}_t^i\| < \infty$  a.s., for all  $i \in N$ .*

Let  $\mathcal{F}_t = \sigma(r_\tau, \mathbf{w}_\tau, \mathbf{s}_\tau, \mathbf{a}_\tau, L_{\tau-1}, \tau \leq t)$  be the filtration which is an increasing  $\sigma$ -algebra over time  $t$ . Define the following for notation convenience. Let  $r_t = [r_t^1, \dots, r_t^n]^\top \in \mathbb{R}^n$ , and  $\mathbf{w}_t = [(\mathbf{w}_t^1)^\top, \dots, (\mathbf{w}_t^n)^\top]^\top \in \mathbb{R}^{np}$ . Moreover, let  $A \otimes B$  represent the Kronecker product of any two matrices  $A$  and  $B$ . Let  $y_t = [(y_t^1)^\top, \dots, (y_t^n)^\top]^\top$ , where  $y_{t+1}^i = [(r_{t+1}^i - \psi_t^\top \mathbf{w}_t^i) \psi_t^\top]^\top$ . Recall,  $\psi_t = \psi(\mathbf{s}_t, \mathbf{a}_t)$ . Let  $I$  be the identity matrix of the dimension  $p \times p$ . Then update of  $\mathbf{w}_t$  can be written as

$$\mathbf{w}_{t+1} = (L_t \otimes I)(\mathbf{w}_t + \gamma_t \cdot y_{t+1}). \quad (1)$$

Let  $\mathbb{1} = (1, \dots, 1)$  represents the vector of all 1's. We define the operator  $\langle \mathbf{w} \rangle = \frac{1}{n}(\mathbb{1}^\top \otimes I)\mathbf{w} = \frac{1}{n} \sum_{i \in N} \mathbf{w}^i$ . This  $\langle \mathbf{w} \rangle \in \mathbb{R}^p$  represents the average of the vectors in  $\{\mathbf{w}^1, \mathbf{w}^2, \dots, \mathbf{w}^n\}$ . Moreover, let  $\mathcal{J} = (\frac{1}{n} \mathbb{1} \mathbb{1}^\top) \otimes I \in \mathbb{R}^{np \times np}$  is the projection operator that projects a vector into the consensus subspace  $\{\mathbb{1} \otimes u : u \in \mathbb{R}^p\}$ . Thus  $\mathcal{J}\mathbf{w} = \mathbb{1} \otimes \langle \mathbf{w} \rangle$ . Now define the disagreement vector  $\mathbf{w}_\perp = \mathcal{J}_\perp \mathbf{w} = \mathbf{w} - \mathbb{1} \otimes \langle \mathbf{w} \rangle$ , where  $\mathcal{J}_\perp = I - \mathcal{J}$ . Here  $I$  is  $np \times np$  dimensional identity matrix. The iteration  $\mathbf{w}_t$  can be decomposed as the sum of a vector in disagreement space and a vector in consensus space, i.e.,  $\mathbf{w}_t = \mathbf{w}_{\perp,t} + \mathbb{1} \otimes \langle \mathbf{w}_t \rangle$ . The proof of convergence consists of two steps.

**Step 01:** To show  $\lim_t \mathbf{w}_{\perp,t} = 0$  a.s. From Proposition 1 we have  $\mathbb{P}[\sup_t \|\mathbf{w}_t\| < \infty] = 1$ , i.e.,  $\mathbb{P}[\cup_{K_1 \in \mathbb{Z}^+} \{\sup_t \|\mathbf{w}_t\| < K_1\}] = 1$ . It suffices to show that  $\lim_t \mathbf{w}_{\perp,t} \mathbb{1}_{\{\sup_t \|\mathbf{w}_t\| < K_1\}} = 0$  for any  $K_1 \in \mathbb{Z}^+$ . Lemma 5.5 in [Zhang et al. \(2018\)](#) proves the boundedness of  $\mathbb{E}[\|\beta_t^{-1} \mathbf{w}_{\perp,t}\|^2]$  over the set  $\{\sup_t \|\mathbf{w}_t\| \leq K_1\}$ , for any  $K_1 > 0$ . We state the lemma here.

**Proposition 2** (Lemma 5.5 in [Zhang et al. \(2018\)](#)). *Under assumptions 2, and 3 for any  $K_1 > 0$ , we have*

$$\sup_t \mathbb{E}[\|\beta_t^{-1} \mathbf{w}_{\perp,t}\|^2 \mathbb{1}_{\{\sup_t \|\mathbf{w}_t\| \leq K_1\}}] < \infty.$$

From Proposition 2 we obtain that for any  $K_1 > 0$ ,  $\exists K_2 < \infty$  such that for any  $t \geq 0$ ,  $\mathbb{E}[\|\mathbf{w}_{\perp,t}\|^2] < K_2 \gamma_t^2$  over the set  $\{\sup_t \|\mathbf{w}_t\| < K_1\}$ . Since  $\sum_t \gamma_t^2 < \infty$ , by Fubini's theorem we have  $\sum_t \mathbb{E}(\|\mathbf{w}_{\perp,t}\|^2 \mathbb{1}_{\{\sup_t \|\mathbf{w}_t\| < K_1\}}) < \infty$ . Thus,  $\sum_t \|\mathbf{w}_{\perp,t}\|^2 \mathbb{1}_{\{\sup_t \|\mathbf{w}_t\| < K_1\}} < \infty$  a.s. Therefore,  $\lim_t \mathbf{w}_{\perp,t} \mathbb{1}_{\{\sup_t \|\mathbf{w}_t\| < K_1\}} = 0$  a.s. Since  $\{\sup_t \|\mathbf{w}_t\| < \infty\}$  with probability 1, thus  $\lim_t \mathbf{w}_{\perp,t} = 0$  a.s. This ends the proof of Step 01.

**Step 02:** To show the convergence of the consensus vector  $\mathbb{1} \otimes \langle \mathbf{w}_t \rangle$ , first note that the iteration of  $\langle \mathbf{w}_t \rangle$  (Equation (1)) can be written as

$$\begin{aligned} \langle \mathbf{w}_{t+1} \rangle &= \frac{1}{N}(\mathbb{1}^\top \otimes I)(L_t \otimes I)(\mathbb{1} \otimes \langle \mathbf{w}_t \rangle + \mathbf{w}_{\perp,t} + \gamma_t y_{t+1}) \\ &= \langle \mathbf{w}_t \rangle + \gamma_t \langle (L_t \otimes I)(y_{t+1} + \gamma_t^{-1} \mathbf{w}_{\perp,t}) \rangle \\ &= \langle \mathbf{w}_t \rangle + \gamma_t \mathbb{E}(\langle y_{t+1} \rangle | \mathcal{F}_t) + \beta_t \xi_{t+1}, \end{aligned} \quad (2)$$

where

$$\begin{aligned} \xi_{t+1} &= \langle (L_t \otimes I)(y_{t+1} + \gamma_t^{-1} \mathbf{w}_{\perp,t}) \rangle - \mathbb{E}(\langle y_{t+1} \rangle | \mathcal{F}_t), \text{ and} \\ \langle y_{t+1} \rangle &= [(\bar{r}_{t+1} - \psi_t^\top \langle \mathbf{w}_t \rangle) \psi_t^\top]^\top. \end{aligned}$$

Note that  $\mathbb{E}(\langle y_{t+1} \rangle | \mathcal{F}_t)$  is Lipschitz continuous in  $\langle \mathbf{w}_t \rangle$ . Moreover,  $\xi_{t+1}$  is a martingale difference sequence and satisfies

$$\mathbb{E}[\|\xi_{t+1}\|^2 | \mathcal{F}_t] \leq \mathbb{E}[\|y_{t+1} + \gamma_t^{-1} \mathbf{w}_{\perp,t}\|_{R_t}^2 | \mathcal{F}_t] + \|\mathbb{E}(\langle y_{t+1} \rangle | \mathcal{F}_t)\|^2, \quad (3)$$

where  $R_t = \frac{L_t^\top \mathbb{1} \mathbb{1}^\top L_t \otimes I}{n^2}$  has bounded spectral norm. Bounding first and second terms in RHS of Equation (3), we have, for any  $K_1 > 0$

$$\mathbb{E}(\|\xi_{t+1}\|^2 | \mathcal{F}_t) \leq K_3(1 + \|\langle \mathbf{w}_t \rangle\|^2),$$

over the set  $\{\sup_t \|\mathbf{w}_t\| \leq K_1\}$  for some  $K_3 < \infty$ . Thus condition (3) of assumption 4 is satisfied. The ODE associated with the Equation (2) has the form

$$\langle \dot{\mathbf{w}} \rangle = -\Psi^\top D^{s,a} \Psi \langle \mathbf{w} \rangle + \Psi^\top D^{s,a} \bar{r} \quad (4)$$

Let the RHS of Equation (4) be  $h(\langle \mathbf{w} \rangle)$ . Note that  $h(\langle \mathbf{w} \rangle)$  is Lipschitz continuous in  $\langle \mathbf{w} \rangle$ . Also, recall that  $D^{s,\mathbf{a}} = \text{diag}[d(s) \cdot \pi(s, \mathbf{a}), s \in \mathcal{S}, \mathbf{a} \in \mathcal{A}]$ . Hence the ODE given in Equation (4) has unique globally asymptotically stable equilibrium  $\mathbf{w}^*$  satisfying

$$\Psi^\top D^{s,\mathbf{a}}(\bar{r} - \Psi \mathbf{w}^*) = 0.$$

Moreover, from Propositions 1, and 2, the sequence  $\{\mathbf{w}_t\}$  is bounded almost surely, so is the sequence  $\{\langle \mathbf{w}_t \rangle\}$ . Specializing Corollary 8.1 and Theorem 8.3 on page 114-115 in Borkar (2022) we have  $\lim_t \langle \mathbf{w}_t \rangle = \mathbf{w}^*$  a.s. over the set  $\{\text{sup}_t \|\mathbf{w}_t\| \leq K_1\}$  for any  $K_1 > 0$ . This concludes the proof of Step 02.

The proof of the Theorem follows from Proposition 1 and results from Step 01. Thus, we have  $\lim_t \mathbf{w}_t^i = \mathbf{w}^*$  a.s. for each  $i \in N$ . This implies,  $\lim_k \mathbf{w}_{kh}^i = \mathbf{w}^*$  a.s. for each  $i \in N$  and for all  $h \in [H]$ .  $\square$

## B Proof of Lemma 2

*Proof.* The proof of this lemma is by induction over  $h$ . Consider the basic case  $h = H+1$ . By assumption we have that  $Q_{k,H+1}^i(\cdot, \cdot) = 0 = Q_{H+1}^{*,i}(\cdot, \cdot)$ , and  $V_{k,H+1}^i(\cdot) = 0 = V_{H+1}^{*,i}(\cdot)$ . Now suppose that this statement is true for all  $h+1$ , so we have  $Q_{k,h+1}^i(\cdot, \cdot) \geq Q_{k,h+1}^{*,i}(\cdot, \cdot)$ , and  $V_{k,h+1}^i(\cdot) \geq V_{k,h+1}^{*,i}(\cdot)$ . For any stage  $h$  and  $(s, \mathbf{a})$ , if  $Q_{k,h}^i(s, \mathbf{a}) \geq H$ , then the statement is also true for stage  $h$ , i.e.,  $Q_{k,h}^i(s, \mathbf{a}) \geq H \geq Q_h^{*,i}(s, \mathbf{a})$ . However, if  $Q_{k,h}^i(s, \mathbf{a}) \leq H$ , then consider the following:

$$\begin{aligned} Q_{k,h}^i(s, \mathbf{a}) - Q_h^{*,i}(s, \mathbf{a}) &\stackrel{(i)}{=} \bar{r}_h(s, \mathbf{a}; \mathbf{w}_{k,h}^i) + \left\langle \hat{\boldsymbol{\theta}}_{k,h}^i, \phi_{V_{k,h+1}^i}(s, \mathbf{a}) \right\rangle + \beta_k \|\Sigma_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(s, \mathbf{a})\|_2 \\ &\quad - \bar{r}_h(s, \mathbf{a}; \mathbf{w}_{k,h}^i) - \mathbb{P}_h V_{h+1}^{*,i}(s, \mathbf{a}) \\ &= \left\langle \hat{\boldsymbol{\theta}}_{k,h}^i, \phi_{V_{k,h+1}^i}(s, \mathbf{a}) \right\rangle + \beta_k \|\Sigma_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(s, \mathbf{a})\|_2 - \mathbb{P}_h V_{h+1}^{*,i}(s, \mathbf{a}) \\ &\stackrel{(ii)}{=} \left\langle \hat{\boldsymbol{\theta}}_{k,h}^i, \phi_{V_{k,h+1}^i}(s, \mathbf{a}) \right\rangle + \beta_k \|\Sigma_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(s, \mathbf{a})\|_2 - \mathbb{P}_h V_{h+1}^{*,i}(s, \mathbf{a}) \\ &\quad - \left\langle \boldsymbol{\theta}_h^*, \phi_{V_{k,h+1}^i}(s, \mathbf{a}) \right\rangle + \left\langle \boldsymbol{\theta}_h^*, \phi_{V_{k,h+1}^i}(s, \mathbf{a}) \right\rangle \\ &= \beta_k \|\Sigma_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(s, \mathbf{a})\|_2 - \left\langle \hat{\boldsymbol{\theta}}_{k,h}^i - \boldsymbol{\theta}_h^*, \phi_{V_{k,h+1}^i}(s, \mathbf{a}) \right\rangle \\ &\quad - \mathbb{P}_h V_{h+1}^{*,i}(s, \mathbf{a}) + \mathbb{P}_h V_{k,h+1}^i(s, \mathbf{a}) \\ &\stackrel{(iii)}{\geq} \beta_k \|\Sigma_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(s, \mathbf{a})\|_2 - \|\Sigma_{k,h}^{i1/2} (\hat{\boldsymbol{\theta}}_{k,h}^i - \boldsymbol{\theta}_h^*)\|_2 \|\Sigma_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(s, \mathbf{a})\|_2 \\ &\quad - \mathbb{P}_h V_{h+1}^{*,i}(s, \mathbf{a}) + \mathbb{P}_h V_{k,h+1}^i(s, \mathbf{a}) \\ &\stackrel{(iv)}{\geq} -\mathbb{P}_h V_{h+1}^{*,i}(s, \mathbf{a}) + \mathbb{P}_h V_{k,h+1}^i(s, \mathbf{a}) \\ &\stackrel{(v)}{\geq} 0 \end{aligned}$$

In (i) we use the update of  $Q_{k,h}^i(s, \mathbf{a})$  as in line 7 of the MA-LDP algorithm 1. In (ii) we add and subtract a inner product term. Inequality (iii) follows from the Cauchy-Schwartz inequality. The inequality (iii) follows from Lemma 3, and (iv) is by induction assumption. Finally, the last inequality (v) uses the monotone property of  $\mathbb{P}_h$  with respect to the partial ordering of the function.  $\square$

## C Finding $\beta_k$ (Proof of Lemma 3)

In this Section, we give the details of  $\beta_k$  for each of the noise adding mechanism. To this end, we first decompose  $\|(\Sigma_{k,h}^i)^{1/2} (\hat{\boldsymbol{\theta}}_{k,h}^i - \boldsymbol{\theta}_h^*)\| = \|(\Sigma_{k,h}^i)^{-1/2} (\mathbf{q}_1 + \mathbf{q}_2 + \mathbf{q}_3)\|$  into three terms. Note that this decomposition is common to all the noise adding mechanism. However, to get the exact expression we bound each term differently for different noise mechanisms. Consider the difference

$$\hat{\boldsymbol{\theta}}_{k,h}^i - \boldsymbol{\theta}_h^* \stackrel{(i)}{=} (\Sigma_{k,h}^i)^{-1} \sum_{\tau=1}^{k-1} \{\phi_{V_{\tau,h+1}^i}(s_{\tau,h}, \mathbf{a}_{\tau,h}) V_{\tau,h+1}^i(s_{\tau,h+1}) + \boldsymbol{\xi}_{\tau,h}^i\} - \boldsymbol{\theta}_h^*$$

$$\begin{aligned}
 &= (\Sigma_{k,h}^i)^{-1} \left\{ -\Sigma_{k,h}^i \boldsymbol{\theta}_h^* + \sum_{\tau=1}^{k-1} \{ \phi_{V_{\tau,h+1}^i} V_{\tau,h+1}^i + \boldsymbol{\xi}_{\tau,h}^i \} \right\} \\
 &\stackrel{(ii)}{=} (\Sigma_{k,h}^i)^{-1} \left\{ -(\lambda \mathbf{I} + \sum_{\tau=1}^{k-1} \phi_{V_{\tau,h+1}^i} \phi_{V_{\tau,h+1}^i}^\top + \mathbf{W}_h^i) \boldsymbol{\theta}_h^* + \sum_{\tau=1}^{k-1} \{ \phi_{V_{\tau,h+1}^i} V_{\tau,h+1}^i + \boldsymbol{\xi}_{\tau,h}^i \} \right\} \\
 &= (\Sigma_{k,h}^i)^{-1} \left\{ -\lambda \boldsymbol{\theta}_h^* - \sum_{\tau=1}^{k-1} \phi_{V_{\tau,h+1}^i} \phi_{V_{\tau,h+1}^i}^\top \boldsymbol{\theta}_h^* - \mathbf{W}_h^i \boldsymbol{\theta}_h^* + \sum_{\tau=1}^{k-1} \{ \phi_{V_{\tau,h+1}^i} V_{\tau,h+1}^i + \boldsymbol{\xi}_{\tau,h}^i \} \right\} \\
 &\stackrel{(iii)}{=} (\Sigma_{k,h}^i)^{-1} \left\{ (-\lambda \mathbf{I} - \mathbf{W}_h^i) \boldsymbol{\theta}_h^* + \sum_{\tau=1}^{k-1} \phi_{V_{\tau,h+1}^i} [V_{\tau,h+1}^i - \mathbb{P}_h V_{\tau,h+1}^i] + \sum_{\tau=1}^{k-1} \boldsymbol{\xi}_{\tau,h}^i \right\}
 \end{aligned}$$

where  $\mathbf{W}_h^i = \sum_{\tau=1}^{k-1} \mathbf{W}_{\tau,h}^i$ . Here (i) uses the definition of  $\hat{\boldsymbol{\theta}}_{k,h}^i$  given in line 30 of the algorithm 1. The (ii) uses the update definition of  $\Sigma_{k,h}^i$ . In (iii) we combine some terms and use the linear function approximation of the transition probability. So, from the above equation we can write the following:

$$\|(\Sigma_{k,h}^i)^{1/2}(\hat{\boldsymbol{\theta}}_{k,h}^i - \boldsymbol{\theta}_h^*)\| = \|(\Sigma_{k,h}^i)^{-1/2}(\mathbf{q}_1 + \mathbf{q}_2 + \mathbf{q}_3)\| \quad (5)$$

where  $\mathbf{q}_1 = (-\lambda \mathbf{I} - \mathbf{W}_h^i) \boldsymbol{\theta}_h^*$ ;  $\mathbf{q}_2 = \sum_{\tau=1}^{k-1} \phi_{V_{\tau,h+1}^i} [V_{\tau,h+1}^i - \mathbb{P}_h V_{\tau,h+1}^i]$ ; and  $\mathbf{q}_3 = \sum_{\tau=1}^{k-1} \boldsymbol{\xi}_{\tau,h}^i$ .

To complete the proof we need to bound each  $\mathbf{q}_1, \mathbf{q}_2, \mathbf{q}_3$ . Here, we give the general framework for bounding each of these terms, and later we will specialize them to each of these noise adding mechanism.

### Bounding $\mathbf{q}_1$

To bound  $\mathbf{q}_1$ , we need to give the upper and the lower bound on the eigenvalues of the symmetric matrix  $\mathbf{W}_h^i$ . Recall, each entry of matrix  $\mathbf{W}_{j,h}^i$  is sampled from the a distribution (Gaussian, Laplace, uniform, and bounded Laplace). So, the variance of the matrix  $\mathbf{W}_h^i$  is  $(k-1)\sigma^2$ , where  $\sigma^2$  is the variance of the corresponding noise distribution. So, from the known concentration results of Tao (2012), we have that

$$\mathbb{P} \left( \left\| \sum_{\tau=1}^{k-1} \mathbf{W}_{\tau,h}^i \right\| \geq \sigma \sqrt{k-1} (\sqrt{4nd} + 2 \log(6H/\alpha)) \right) \leq \frac{\alpha}{6H} \quad (6)$$

that is

$$\mathbb{P} (\| \mathbf{W}_h^i \| \geq \Gamma) \leq \frac{\alpha}{6H} \quad (7)$$

where  $\Gamma = \sigma \sqrt{k-1} (\sqrt{4nd} + 2 \log(6H/\alpha))$ .

For the symmetric matrix  $\mathbf{W}_h^i$  the PSD property might not be preserved, so we add a basic matrix  $2\Gamma \mathbf{I}$  to the matrix  $\mathbf{W}_h^i = \sum_{\tau=1}^{k-1} \mathbf{W}_{\tau,h}^i$  for each stage  $h \in [H]$ . Thus, the eigenvalues of the shifted matrix are bounded in the interval  $[\Gamma, 3\Gamma]$  with probability  $1 - \alpha/6$ . Define the following event

$$\mathcal{E}_1 := \{ \forall h \in [H], \forall j \in [nd], \Gamma \leq \sigma_j \leq 3\Gamma \}, \quad (8)$$

where  $\sigma_j$ 's are the eigenvalues of the matrix  $\mathbf{W}_h^i$ , and we have  $\mathbb{P}(\mathcal{E}_1) \geq 1 - \alpha/6$ . Let  $\rho_{\max} = 3\Gamma + \lambda$ , and  $\rho_{\min} = \Gamma + \lambda$ . Then for the term  $\mathbf{q}_1$ , we have

$$\begin{aligned}
 \|(\Sigma_{k,h}^i)^{-1/2} \mathbf{q}_1\| &\stackrel{(i)}{\leq} \|(\mathbf{W}_h^i + \lambda \mathbf{I})^{-1/2} \mathbf{q}_1\| \\
 &= \|(\mathbf{W}_h^i + \lambda \mathbf{I})^{-1/2} (-\mathbf{W}_h^i - \lambda \mathbf{I}) \boldsymbol{\theta}_h^*\| \\
 &= \|(\mathbf{W}_h^i + \lambda \mathbf{I})^{1/2} \boldsymbol{\theta}_h^*\| \\
 &\stackrel{(ii)}{\leq} \sqrt{\rho_{\max}} \|\boldsymbol{\theta}_h^*\| \\
 &\stackrel{(iii)}{\leq} \sqrt{\rho_{\max} nd} \\
 &= \sqrt{\{3\sigma \sqrt{k-1} (\sqrt{4nd} + 2 \log(6H/\alpha)) + \lambda\} nd}
 \end{aligned}$$

$$= \sqrt{3nd\sigma\sqrt{k-1}(\sqrt{4nd} + 2\log(6H/\alpha))} + \lambda nd.$$

Since  $\rho_{\max} = 3\Gamma + \lambda = 3\sigma\sqrt{k-1}(\sqrt{4nd} + 2\log(6H/\alpha)) + \lambda$  and  $\|\boldsymbol{\theta}_h^*\| \leq \sqrt{nd}$ . The first inequality holds because  $(\Sigma_{k,h}^i)^{-1/2} \succeq (\mathbf{W}_h^i + \lambda\mathbf{I})^{-1/2}$ . The second inequality holds due to event  $\mathcal{E}_1$ . And the last inequality (iii) holds because of the Assumption 1.

### Bounding $\mathbf{q}_2$

For the term  $\mathbf{q}_2$ , we have the following

$$\|(\Sigma_{k,h}^i)^{-1/2}\mathbf{q}_2\| = \left\| \sum_{\tau=1}^{k-1} \phi_{V_{\tau,h+1}^i} [V_{\tau,h+1}^i - \mathbb{P}_h V_{\tau,h+1}^i] \right\|_{(\Sigma_{k,h}^i)^{-1}} \quad (9)$$

$$\leq \left\| \sum_{\tau=1}^{k-1} \phi_{V_{\tau,h+1}^i} [V_{\tau,h+1}^i - \mathbb{P}_h V_{\tau,h+1}^i] \right\|_{\mathbf{Z}^{-1}}, \quad (10)$$

where  $\mathbf{Z} = \lambda\mathbf{I} + \sum_{\tau=1}^{k-1} \phi_{V_{\tau,h+1}^i} \phi_{V_{\tau,h+1}^i}^\top$ . The inequality holds because,  $\Sigma_{k,h}^i \succeq \mathbf{Z} = \lambda\mathbf{I} + \sum_{\tau=1}^{k-1} \phi_{V_{\tau,h+1}^i} \phi_{V_{\tau,h+1}^i}^\top$ . Let  $\eta_{\tau,h+1}^i = V_{\tau,h+1}^i - \mathbb{P}_h V_{\tau,h+1}^i = V_{\tau,h+1}^i - \phi_{V_{\tau,h+1}^i}^\top \boldsymbol{\theta}_h^*$ . Moreover, let  $\{\mathcal{G}_t\}_{t=1}^\infty$  be a filtration,  $\{\phi_{V_{\tau,t}^i}, \eta_{\tau,t}^i\}_{t=1}^\infty$  a stochastic process so that  $\phi_{V_{\tau,t}^i}$  is  $\mathcal{G}_t$ -measurable and  $\eta_{\tau,t}^i$  is  $\mathcal{G}_{t+1}$ -measurable. With above notations, we have

$$|\eta_{\tau,h}^i| = |V_{\tau,h}^i - \mathbb{P}_h V_{\tau,h}^i| \leq H. \quad (11)$$

The above is true because  $V_{\tau,h}^i \leq H$ . Further, we have

$$\mathbb{E}[(\eta_{\tau,h}^i)^2 | \mathcal{G}_h] \leq \mathbb{E}[(V_{\tau,h}^i)^2 | \mathcal{G}_h] \leq H^2. \quad (12)$$

Moreover, we define the following event

$$\mathcal{E}_2 = \left\{ \forall h \in [H], \|\mathbf{q}_2\|_{\mathbf{Z}^{-1}} \leq 4H \left( 2\sqrt{nd \log \left( 1 + \frac{(k-1)H^2}{nd\lambda} \right)} \log \left( \frac{24(k-1)^2}{\alpha} \right) + \log \left( \frac{24(k-1)^2}{\alpha} \right) \right) \right\}. \quad (13)$$

From Theorem 2 of Zhou et al. (2021) we have that the probability of the above event is at least  $1 - \alpha/6$ .

### Bounding $\mathbf{q}_3$

The term  $\mathbf{q}_3$  can be bounded as

$$\left\| \sum_{\tau=1}^{k-1} \boldsymbol{\xi}_{\tau,h}^i \right\|_{(\Sigma_{k,h}^i)^{-1}} \leq \left\| \sum_{\tau=1}^{k-1} \boldsymbol{\xi}_{\tau,h}^i \right\|_{(\mathbf{W}_h^i + \lambda\mathbf{I})^{-1}} \leq \frac{1}{\sqrt{\rho_{\min}}} \left\| \sum_{\tau=1}^{k-1} \boldsymbol{\xi}_{\tau,h}^i \right\|_2 \quad (14)$$

where the first inequality holds due to the fact that  $\Sigma_{k,h}^i \succeq \mathbf{W}_h^i + \lambda\mathbf{I}$ , and the second inequality holds due to the definition of event  $\mathcal{E}_1$ . So, with probability  $1 - \alpha/6H$ , we have

$$\left\| \sum_{\tau=1}^{k-1} \boldsymbol{\xi}_{\tau,h}^i \right\|_2 \leq \sigma \sqrt{(k-1)nd \log \frac{12ndH}{\alpha}}. \quad (15)$$

We define the event  $\mathcal{E}_3$  as

$$\mathcal{E}_3 = \left\{ \forall h \in [H] : \left\| \sum_{\tau=1}^{k-1} \boldsymbol{\xi}_{\tau,h}^i \right\|_2 \leq \sigma \sqrt{(k-1)nd \log \frac{12ndH}{\alpha}} \right\}. \quad (16)$$

Taking the union bound on all the stage  $h \in [H]$ , we have

$$\left\| \sum_{\tau=1}^{k-1} \boldsymbol{\xi}_{\tau,h}^i \right\|_{(\Sigma_{k,h}^i)^{-1}} \stackrel{(i)}{\leq} \frac{1}{\sqrt{\rho_{\min}}} \left\| \sum_{\tau=1}^{k-1} \boldsymbol{\xi}_{\tau,h}^i \right\|_2$$

$$\begin{aligned}
 & \stackrel{(ii)}{\leq} \frac{\sigma \sqrt{(k-1)nd \log \frac{12ndH}{\alpha}}}{\sqrt{\rho_{\min}}} \\
 & \stackrel{(iii)}{=} \frac{\sigma \sqrt{(k-1)nd \log \frac{12ndH}{\alpha}}}{\sqrt{\sigma \sqrt{k-1}(\sqrt{4nd} + 2 \log(6H/\alpha)) + \lambda}} \\
 & \stackrel{(iv)}{\leq} \frac{\sqrt{\sigma}(k-1)^{1/4} \sqrt{nd \log \frac{12ndH}{\alpha}}}{\sqrt{\sqrt{4nd} + 2 \log(6H/\alpha)}} \\
 & \stackrel{(v)}{\leq} (nd)^{1/4} \sqrt{\sigma}(k-1)^{1/4} \sqrt{\log(12ndH/\alpha)}. \tag{17}
 \end{aligned}$$

The inequality (i) follows from the relation between the  $l_2$  norm and the norm on  $(\Sigma_{k,h}^i)^{-1}$  and hence bounded by the  $\frac{1}{\sqrt{\rho_{\min}}}$ . The inequality (ii) follows from definition of event  $\mathcal{E}_3$ . In (iii) we use the definition of  $\rho_{\min}$ . Inequality (iv) follows after dropping  $\lambda \geq 0$  in the previous expression. Finally (v) follows by combining the common terms.

Combining the above bounds for  $\mathbf{q}_1, \mathbf{q}_2, \mathbf{q}_3$  from Equations (9), (13), and (17), we have that with probability at least  $1 - \alpha/2$ , we have for each  $i \in N$  and for each  $h \in [H]$ ,

$$\|(\Sigma_{k,h}^i)^{1/2}(\hat{\theta}_{k,h}^i - \theta_h^*)\| \leq \beta_k \tag{18}$$

here,

$$\beta_k = c(nd)^{3/4} \sqrt{\sigma} k^{1/4} \log(ndT/\alpha) \tag{19}$$

where  $c$  is an absolute constant (different for different mechanisms.)

In the next Section, we identify the  $\beta_k$  for each noise adding mechanism.

### C.1 $\beta_k^G$ for the Gaussian mechanism

For the  $(\epsilon, \delta)$  LDP for the Gaussian noise adding mechanism, the  $\sigma$  is taken as  $\sigma = \frac{2H\sqrt{2\log(2.5H/\delta)}\Delta f}{\epsilon}$ . This is obtained using the Lemma 5 given in Appendix G.2. Moreover,  $\Delta f$  is the  $l_2$  sensitivity which is identified as  $2H^2$ . Using this  $\sigma = \frac{4H^3\sqrt{2\log(2.5H/\delta)}}{\epsilon}$  in Equation (19), we have  $\beta_k^G$  for the Gaussian mechanism is

$$\beta_k^G = c_g(nd)^{3/4} H^{3/2} k^{1/4} \log(ndT/\alpha) (\log(H/\delta))^{1/4} \sqrt{1/\epsilon}. \tag{20}$$

### C.2 $\beta_k^L$ for the Laplace mechanism

The variance of the Laplace mechanism is  $2b^2$ , and for  $(\epsilon, 0)$  LDP we identify  $b = \frac{2H\Delta f}{\epsilon}$ , where  $\Delta f$  is the  $l_1$  sensitivity. The  $l_1$  sensitivity in MA-LDP is  $2H^2\sqrt{nd}$ . Therefore,  $b = \frac{4H^3\sqrt{nd}}{\epsilon}$ . Substituting this in Equation (19), we have  $\beta_k^L$  for the Laplace mechanism as

$$\beta_k^L = c_l(nd)^{3/4} H^{3/2} k^{1/4} \log(ndT/\alpha) \sqrt{1/\epsilon}. \tag{21}$$

### C.3 $\beta_k^U$ for the uniform mechanism

The variance of the uniform mechanism is  $a^2/3$ , and for  $(0, \delta)$  LDP we identify  $a = 4H^3\sqrt{\log(\frac{2H}{\delta})}$ . Substituting this in Equation (19), we have  $\beta_k^U$  for the Laplace mechanism as

$$\beta_k^U = c_u(nd)^{3/4} H^{3/2} k^{1/4} \log(ndT/\alpha) (\log(H/\delta))^{1/4}. \tag{22}$$

### C.4 $\beta_k^{BL}$ for the bounded Laplace mechanism

The variance of the bounded Laplace mechanism is  $\frac{2b^2}{1 - \exp(-\frac{B}{b})} - \kappa$ , where  $\kappa = \frac{((B+b)^2 + b^2) \times \exp(-\frac{B}{b})}{1 - \exp(-\frac{B}{b})}$ . Thus, for  $(\epsilon, 0)$  LDP, similar to Laplace with unbounded support, we identify  $b = \frac{4H^3\sqrt{nd}}{\epsilon}$ . Substituting this in Equation (19), we have  $\beta_k^{BL}$  for the bounded Laplace mechanism as

$$\beta_k^{BL} = c_{bl}(nd)^{3/4} \zeta^{1/4} k^{1/4} \log(ndT/\alpha). \tag{23}$$



## D Proof of privacy guarantee of noise adding mechanism

In this Section, we prove the privacy guarantees of the MA-LDP algorithm for all the noise adding mechanisms. To this end, we first find the  $l_1$  and  $l_2$  sensitivity of the information shared to the server by each agent  $i \in N$ . To this end, we first compute the  $l_2$  sensitivity coefficient for the MA-DP algorithm. Let  $\Delta \tilde{u}_{k,h}^i$  and  $\Delta \tilde{\Lambda}_{k,h}^i$  be the noise-free information of agent  $i \in N$  at the  $h$ -th stage of  $k$ -th episode. That is,

$$\begin{aligned}\Delta \tilde{u}_{k,h}^i &= \phi_{V_{k,h+1}^i}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) V_{k,h+1}^i(\mathbf{s}_{k,h+1}) \\ \Delta \tilde{\Lambda}_{k,h}^i &= \phi_{V_{k,h+1}^i}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) \phi_{V_{k,h+1}^i}^\top(\mathbf{s}_{k,h}, \mathbf{a}_{k,h})^\top\end{aligned}$$

For  $\Delta \tilde{u}_{k,h}^i$ , the sensitivity coefficient is upper bounded as

$$\|\Delta \tilde{u}_{k,h}^i - (\Delta \tilde{u}_{k,h}^i)'\|_2 \leq \|\phi_{V_{k,h+1}^i}\| \cdot |V_{k,h+1}^i| + \|\phi_{V_{k,h+1}^i}'\| \cdot |V_{k,h+1}^i| \leq 2H^2 \quad (24)$$

Similarly, the sensitivity of  $\Delta \tilde{\Lambda}_{k,h}^i$  is upper bounded as

$$\begin{aligned}\|\phi_{V^i} \phi_{V^i}^\top - \phi_{V^i}' \phi_{V^i}'^\top\|_F &\leq \|\phi_{V^i} \phi_{V^i}^\top\|_F + \|\phi_{V^i}' \phi_{V^i}'^\top\|_F \\ &= \sqrt{\text{tr}[\phi_{V^i} \phi_{V^i}^\top \phi_{V^i} \phi_{V^i}^\top]} + \sqrt{\text{tr}[\phi_{V^i}' \phi_{V^i}'^\top \phi_{V^i}' \phi_{V^i}'^\top]} \\ &= \phi_{V^i}^\top \phi_{V^i} + \phi_{V^i}'^\top \phi_{V^i}' \\ &\leq 2H^2\end{aligned}$$

Thus,  $l_2$  sensitivity of both the information is  $2H^2$ . Next we find the  $l_1$  sensitivities. Recall, for any matrix  $A \in \mathbb{R}^{l \times l}$ , we have that  $\|A\|_1 \leq \sqrt{l} \|A\|_2$ . Similarly, for any vector  $\mathbf{x} \in \mathbb{R}^l$ , we have  $\|\mathbf{x}\|_1 \leq \sqrt{l} \|\mathbf{x}\|_2$ . Using this property on the information's we have

$$\|\Delta \tilde{u}_{k,h}^i - (\Delta \tilde{u}_{k,h}^i)'\|_2 \leq \sqrt{nd} \|\Delta \tilde{u}_{k,h}^i - (\Delta \tilde{u}_{k,h}^i)'\|_1 \leq 2\sqrt{nd}H^2 \quad (25)$$

Similarly, we have

$$\|\phi_{V^i} \phi_{V^i}^\top - \phi_{V^i}' \phi_{V^i}'^\top\|_F \leq \|\phi_{V^i} \phi_{V^i}^\top\|_F + \|\phi_{V^i}' \phi_{V^i}'^\top\|_F \leq \sqrt{nd} \|\phi_{V^i} \phi_{V^i}^\top\|_1 + \sqrt{nd} \|\phi_{V^i}' \phi_{V^i}'^\top\|_1 \leq 2\sqrt{nd}H^2$$

Thus, the  $l_1$  sensitivity of both the information is  $2\sqrt{nd}H^2$ . Let  $\mathbf{D}_h = (D_h^1, D_h^2, \dots, D_h^n)$  and  $\mathbf{D}'_h = (D_h'^1, D_h'^2, \dots, D_h'^n)$  are the different datasets collected by the server at stage  $h$ . For simplicity of notation, let  $\mathbf{M} = (\mathbf{M}^1, \mathbf{M}^2, \dots, \mathbf{M}^n)$  and let  $\boldsymbol{\alpha} = (\alpha^1, \alpha^2, \dots, \alpha^n)$ . Moreover, let  $(\mathbf{M}, \boldsymbol{\alpha})$  be a possible outcome of the algorithm. Further, let  $\Delta \boldsymbol{\Lambda}_{k,h} = (\Delta \Lambda_{k,h}^1, \Delta \Lambda_{k,h}^2, \dots, \Delta \Lambda_{k,h}^n)$  and  $\Delta \mathbf{u}_{k,h} = (\Delta u_{k,h}^1, \Delta u_{k,h}^2, \dots, \Delta u_{k,h}^n)$  be the information from all the agents. Let  $\mathbf{D}_{1:h-1}$  be the information collected from stage 1 to stage  $h$ , i.e.,  $\mathbf{D}_{1:h-1} = (\mathbf{D}_1, \mathbf{D}_2, \dots, \mathbf{D}_{h-1})$ . Also let  $\mathbf{W}_{k,h} = (W_{k,h}^1, W_{k,h}^2, \dots, W_{k,h}^n)$  and  $\boldsymbol{\xi}_{k,h} = (\xi_{k,h}^1, \xi_{k,h}^2, \dots, \xi_{k,h}^n)$ . Then, we have

$$\begin{aligned}&\mathbb{P}(\forall h \in [H], (\Delta \boldsymbol{\Lambda}_{k,h}, \Delta \mathbf{u}_{k,h}) = (\mathbf{M}, \boldsymbol{\alpha}) \mid \mathbf{D}_{1:h-1}) \\ &\mathbb{P}(\forall h \in [H], ((\Delta \boldsymbol{\Lambda}_{k,h})', (\Delta \mathbf{u}_{k,h})') = (\mathbf{M}, \boldsymbol{\alpha}) \mid \mathbf{D}'_{1:h-1}) \\ &= \prod_{h=1}^H \frac{\mathbb{P}((\mathbf{W}_{k,h}, \boldsymbol{\xi}_{k,h}) = (\mathbf{M} - \Delta \tilde{\boldsymbol{\Lambda}}_{k,h}, \boldsymbol{\alpha} - \Delta \tilde{\mathbf{u}}_{k,h}) \mid \mathbf{D}_{1:h-1})}{\mathbb{P}(((\mathbf{W}_{k,h})', (\boldsymbol{\xi}_{k,h})') = (\mathbf{M} - (\Delta \tilde{\boldsymbol{\Lambda}}_{k,h})', \boldsymbol{\alpha} - (\Delta \tilde{\mathbf{u}}_{k,h})') \mid \mathbf{D}'_{1:h-1})} \\ &= \prod_{h=1}^H \frac{\mathbb{P}((\mathbf{W}_{k,h}, \boldsymbol{\xi}_{k,h}) = (\mathbf{M} - \Delta \tilde{\boldsymbol{\Lambda}}_{k,h}, \boldsymbol{\alpha} - \Delta \tilde{\mathbf{u}}_{k,h}) \mid \mathbf{D}_{h-1})}{\mathbb{P}(((\mathbf{W}_{k,h})', (\boldsymbol{\xi}_{k,h})') = (\mathbf{M} - (\Delta \tilde{\boldsymbol{\Lambda}}_{k,h})', \boldsymbol{\alpha} - (\Delta \tilde{\mathbf{u}}_{k,h})') \mid \mathbf{D}'_{h-1})} \\ &= \prod_{h=1}^H \frac{\mathbb{P}(\mathbf{W}_{k,h} = \mathbf{M} - \Delta \tilde{\boldsymbol{\Lambda}}_{k,h} \mid \mathbf{D}_{h-1}) \times \mathbb{P}(\boldsymbol{\xi}_{k,h} = \boldsymbol{\alpha} - \Delta \tilde{\mathbf{u}}_{k,h} \mid \mathbf{D}_{h-1})}{\mathbb{P}((\mathbf{W}_{k,h})' = \mathbf{M} - (\Delta \tilde{\boldsymbol{\Lambda}}_{k,h})' \mid \mathbf{D}'_{h-1}) \times \mathbb{P}((\boldsymbol{\xi}_{k,h})' = \boldsymbol{\alpha} - (\Delta \tilde{\mathbf{u}}_{k,h})' \mid \mathbf{D}'_{h-1})}\end{aligned}$$

The first and the second equations again uses the Markov property and the last inequality is true because of the independence of the two information one is for the  $\mathbf{W}_{k,h}$  and  $\boldsymbol{\xi}_{k,h}$ .

### D.1 $(\epsilon, \delta)$ privacy for Gaussian mechanism (Theorem 1)

Consider  $\mathbb{P}(\mathbf{W}_{k,h} = \mathbf{M} - \Delta \tilde{\boldsymbol{\Lambda}}_{k,h} \mid \mathbf{D}_{h-1})$ . Since the Gaussian noise in the information is independent across the agents, we have that

$$\mathbb{P}(\mathbf{W}_{k,h} = \mathbf{M} - \Delta \tilde{\boldsymbol{\Lambda}}_{k,h} \mid \mathbf{D}_{h-1}) = \prod_{i \in N} \mathbb{P}(W_{k,h}^i = \mathbf{M}^i - \Delta \tilde{\Lambda}_{k,h}^i \mid \mathbf{D}_{h-1}) \quad (26)$$

Now from Lemma 5, and the sensitivity definition of  $\tilde{\Delta}\tilde{\Lambda}_{k,h}^i$ , for  $\sigma = 4H^3\sqrt{2\log(2.5H/\delta)}/\epsilon$ , then with probability at least  $1 - \delta/2nH$  for each  $W_{k,h}^i$  for each agent  $i \in N$ , we have that

$$\mathbb{P}(W_{k,h}^i = \mathbf{M}^i - \Delta\tilde{\Lambda}_{k,h}^i \mid \mathbf{D}_{h-1}) \leq \exp\left(\frac{\epsilon}{2nH}\right) \times \mathbb{P}((W_{k,h}^i)' = \mathbf{M}^i - (\Delta\tilde{\Lambda}_{k,h}^i)' \mid \mathbf{D}'_{h-1}) \quad (27)$$

Taking the union bound on the agents and applying the composition theorem, we have that with probability at least  $1 - \delta/(2H)$

$$\mathbb{P}(W_{k,h} = \mathbf{M} - \Delta\tilde{\Lambda}_{k,h} \mid \mathbf{D}_{h-1}) \leq \exp\left(\frac{\epsilon}{2H}\right) \times \mathbb{P}((W_{k,h})' = \mathbf{M} - (\Delta\tilde{\Lambda}_{k,h})' \mid \mathbf{D}'_{h-1}) \quad (28)$$

Moreover, for the other term  $\mathbb{P}(\xi_{k,h} = \alpha - \Delta\tilde{\mathbf{u}}_{k,h} \mid \mathbf{D}_{h-1})$  note again that by the independence of the noise distribution, we have

$$\mathbb{P}(\xi_{k,h} = \alpha - \Delta\tilde{\mathbf{u}}_{k,h} \mid \mathbf{D}_{h-1}) = \prod_{i \in N} \mathbb{P}(\xi_{k,h}^i = \alpha^i - \Delta\tilde{u}_{k,h}^i \mid \mathbf{D}_{h-1}) \quad (29)$$

Using the property of the Gaussian distribution, we have

$$\begin{aligned} \frac{\mathbb{P}(\xi_{k,h}^i = \alpha^i - \Delta\tilde{u}_{k,h}^i \mid \mathbf{D}_{h-1})}{\mathbb{P}((\xi_{k,h}^i)' = \alpha^i - (\Delta\tilde{u}_{k,h}^i)' \mid \mathbf{D}'_{h-1})} &= \frac{\exp\left(-\|\alpha^i - \Delta\tilde{u}_{k,h}^i\|^2/2\sigma^2\right)}{\exp\left(-\|\alpha^i - \Delta\tilde{u}_{k,h}^i + (\Delta\tilde{u}_{k,h}^i - (\Delta\tilde{u}_{k,h}^i)')\|^2/2\sigma^2\right)} \\ &= \exp\left(-\frac{\|\alpha^i - \Delta\tilde{u}_{k,h}^i\|^2}{2\sigma^2} + \frac{\|\alpha^i - \Delta\tilde{u}_{k,h}^i + (\Delta\tilde{u}_{k,h}^i - (\Delta\tilde{u}_{k,h}^i)')\|^2}{2\sigma^2}\right) \\ &\leq \exp\left(-\frac{\|\alpha^i - \Delta\tilde{u}_{k,h}^i\|^2}{2\sigma^2} + \frac{\|\alpha^i - \Delta\tilde{u}_{k,h}^i\|^2}{2\sigma^2}\right. \\ &\quad \left. + \frac{\|\Delta\tilde{u}_{k,h}^i - (\Delta\tilde{u}_{k,h}^i)'\|^2}{2\sigma^2}\right) \\ &= \exp\left(\frac{\|\Delta\tilde{u}_{k,h}^i - (\Delta\tilde{u}_{k,h}^i)'\|^2}{2\sigma^2}\right). \end{aligned}$$

Here the first equation is due to the Gaussian mechanism. The first inequality follows from the triangle inequality. Again using the Lemma 5, with probability at least  $1 - \delta/(2H)$  for each agent  $i \in N$ , we have that

$$\mathbb{P}(\xi_{k,h}^i = \alpha^i - \Delta\tilde{u}_{k,h}^i \mid \mathbf{D}_{h-1}) \leq \exp\left(\frac{\epsilon}{2nH}\right) \times \mathbb{P}((\xi_{k,h}^i)' = \alpha^i - (\Delta\tilde{u}_{k,h}^i)' \mid \mathbf{D}'_{h-1}) \quad (30)$$

Therefore, we have with probability  $1 - \delta/(2H)$  that

$$\mathbb{P}(\xi_{k,h} = \alpha - \Delta\tilde{\mathbf{u}}_{k,h} \mid \mathbf{D}_{h-1}) \leq \exp\left(\frac{\epsilon}{2H}\right) \times \mathbb{P}((\xi_{k,h})' = \alpha - (\Delta\tilde{\mathbf{u}}_{k,h})' \mid \mathbf{D}'_{h-1}) \quad (31)$$

Now taking the union bound for  $\mathbf{W}_{k,h}, \xi_{k,h}$  terms and all the stages  $h \in [H]$ , with probability at least  $1 - (2H) \times \delta/(2H) = 1 - \delta$ , we have

$$\log \left[ \frac{\mathbb{P}(\forall h \in [H], (\Delta\mathbf{\Lambda}_{k,h}, \Delta\mathbf{u}_{k,h}) = (\mathbf{M}, \alpha) \mid \mathbf{D}_{1:h-1})}{\mathbb{P}(\forall h \in [H], ((\Delta\mathbf{\Lambda}_{k,h})', (\Delta\mathbf{u}_{k,h})') = (\mathbf{M}, \alpha) \mid \mathbf{D}'_{1:h-1})} \right] \leq \epsilon. \quad (32)$$

Therefore, from Theorem 12, we conclude that Algorithm 1 preserves  $(\epsilon, \delta)$ -LDP property with the Gaussian mechanism.

## D.2 $(\epsilon, 0)$ privacy for Laplace mechanism (Theorem 3)

To prove that algorithm 1 with Laplace noise adding mechanism achieves  $(\epsilon, 0)$ -LDP, we need to show that the above equation is upper bounded by  $e^\epsilon$  with probability 1. Using the independence of the information across the agents, we have that

$$\mathbb{P}(\mathbf{W}_{k,h} = \mathbf{M} - \Delta\tilde{\Lambda}_{k,h} \mid \mathbf{D}_{h-1}) = \prod_{i \in N} \mathbb{P}(W_{k,h}^i = \mathbf{M}^i - \Delta\tilde{\Lambda}_{k,h}^i \mid \mathbf{D}_{h-1}) \quad (33)$$

Now from Theorem 1, and the sensitivity definition of  $\Delta\tilde{\Lambda}_{k,h}^i$ , if we set  $b = \frac{4H^3\sqrt{nd}}{\epsilon}$ , then for each  $W_{k,h}^i$  for each agent  $i \in N$ , we have that

$$\mathbb{P}(W_{k,h}^i = \mathbf{M}^i - \Delta\tilde{\Lambda}_{k,h}^i \mid \mathbf{D}_{h-1}) \leq \exp\left(\frac{\epsilon}{2nH}\right) \times \mathbb{P}((W_{k,h}^i)' = \mathbf{M}^i - (\Delta\tilde{\Lambda}_{k,h}^i)' \mid \mathbf{D}'_{h-1}) \quad (34)$$

Taking the union bound on the agents and applying the composition theorem, we have that with probability 1,

$$\mathbb{P}(W_{k,h} = \mathbf{M} - \Delta\tilde{\Lambda}_{k,h} \mid \mathbf{D}_{h-1}) \leq \exp\left(\frac{\epsilon}{2H}\right) \times \mathbb{P}((W_{k,h})' = \mathbf{M} - (\Delta\tilde{\Lambda}_{k,h})' \mid \mathbf{D}'_{h-1}) \quad (35)$$

Moreover, for the other term  $\mathbb{P}(\xi_{k,h} = \alpha - \Delta\tilde{\mathbf{u}}_{k,h} \mid \mathbf{D}_{h-1})$  note again that by the independence of the noise distribution, we have

$$\mathbb{P}(\xi_{k,h} = \alpha - \Delta\tilde{\mathbf{u}}_{k,h} \mid \mathbf{D}_{h-1}) = \prod_{i \in N} \mathbb{P}(\xi_{k,h}^i = \alpha^i - \Delta\tilde{u}_{k,h}^i \mid \mathbf{D}_{h-1}) \quad (36)$$

Using the property of the Laplace distribution, we have

$$\begin{aligned} \frac{\mathbb{P}(\xi_{k,h}^i = \alpha^i - \Delta\tilde{u}_{k,h}^i \mid \mathbf{D}_{h-1})}{\mathbb{P}((\xi_{k,h}^i)' = \alpha^i - (\Delta\tilde{u}_{k,h}^i)' \mid \mathbf{D}'_{h-1})} &= \prod_{j=1}^{nd} \frac{\exp\left(\frac{-\epsilon|(\alpha^i)_j - (\Delta\tilde{u}_{k,h}^i)_j|}{2nH|\Delta\tilde{u}_{k,h}^i - (\Delta\tilde{u}_{k,h}^i)'|_1}\right)}{\exp\left(\frac{-\epsilon|(\alpha^i)_j - ((\Delta\tilde{u}_{k,h}^i)')_j|}{2nH|\Delta\tilde{u}_{k,h}^i - (\Delta\tilde{u}_{k,h}^i)'|_1}\right)} \\ &= \prod_{j=1}^{nd} \exp\left(\frac{-\epsilon|(\alpha^i)_j - (\Delta\tilde{u}_{k,h}^i)_j| + \epsilon|(\alpha^i)_j - ((\Delta\tilde{u}_{k,h}^i)')_j|}{2nH|\Delta\tilde{u}_{k,h}^i - (\Delta\tilde{u}_{k,h}^i)'|_1}\right) \\ &\leq \prod_{j=1}^{nd} \exp\left(\frac{\epsilon|(\Delta\tilde{u}_{k,h}^i)_j - ((\Delta\tilde{u}_{k,h}^i)')_j|}{2nH|\Delta\tilde{u}_{k,h}^i - (\Delta\tilde{u}_{k,h}^i)'|_1}\right) \\ &= \exp\left(\frac{\epsilon}{2nH}\right) \end{aligned}$$

In above, we use subscript  $j$  to denote the  $j$ -th element of the corresponding vector. Here the first equation is due to the Laplace mechanism. The first inequality follows from the triangle inequality. So with probability 1 for each agent  $i \in N$ , we have that

$$\mathbb{P}(\xi_{k,h}^i = \alpha^i - \Delta\tilde{u}_{k,h}^i \mid \mathbf{D}_{h-1}) \leq \exp\left(\frac{\epsilon}{2nH}\right) \times \mathbb{P}((\xi_{k,h}^i)' = \alpha^i - (\Delta\tilde{u}_{k,h}^i)' \mid \mathbf{D}'_{h-1}) \quad (37)$$

Therefore, we have with probability 1 that

$$\mathbb{P}(\xi_{k,h} = \alpha - \Delta\tilde{\mathbf{u}}_{k,h} \mid \mathbf{D}_{h-1}) \leq \exp\left(\frac{\epsilon}{2H}\right) \times \mathbb{P}((\xi_{k,h})' = \alpha - (\Delta\tilde{\mathbf{u}}_{k,h})' \mid \mathbf{D}'_{h-1}) \quad (38)$$

Now taking the union bound for  $\mathbf{W}_{k,h}, \xi_{k,h}$  terms and all the stages  $h \in [H]$ , with probability 1, we have

$$\log \left[ \frac{\mathbb{P}(\forall h \in [H], (\Delta\mathbf{\Lambda}_{k,h}, \Delta\mathbf{u}_{k,h}) = (\mathbf{M}, \alpha) \mid \mathbf{D}_{1:h-1})}{\mathbb{P}(\forall h \in [H], ((\Delta\mathbf{\Lambda}_{k,h})', (\Delta\mathbf{u}_{k,h})') = (\mathbf{M}, \alpha) \mid \mathbf{D}'_{1:h-1})} \right] \leq \epsilon \quad (39)$$

Therefore, Algorithm 1 preserves  $(\epsilon, 0)$ -LDP property with the Laplace mechanism.

### D.3 $(0, \delta)$ privacy for uniform mechanism (Theorem 5)

For the uniform distribution  $\epsilon = 0$ , thus setting  $a = 4H^3 \log(2H/\delta)$  satisfies the following with probability at least  $1 - \delta/(2H)$

$$\mathbb{P}(W_{k,h}^i = \mathbf{M}^i - \Delta\tilde{\Lambda}_{k,h}^i \mid \mathbf{D}_{h-1}) \leq \mathbb{P}((W_{k,h}^i)' = \mathbf{M}^i - (\Delta\tilde{\Lambda}_{k,h}^i)' \mid \mathbf{D}'_{h-1}) \quad (40)$$

Moreover, for the other term  $\mathbb{P}(\xi_{k,h} = \alpha - \Delta\tilde{\mathbf{u}}_{k,h} \mid \mathbf{D}_{h-1})$  note again that by the independence of the noise distribution, we have

$$\mathbb{P}(\xi_{k,h} = \alpha - \Delta\tilde{\mathbf{u}}_{k,h} \mid \mathbf{D}_{h-1}) = \prod_{i \in N} \mathbb{P}(\xi_{k,h}^i = \alpha^i - \Delta\tilde{u}_{k,h}^i \mid \mathbf{D}_{h-1}) \quad (41)$$

Using the property of the uniform distribution, we have that

$$\frac{\mathbb{P}(\xi_{k,h}^i = \alpha^i - \Delta\tilde{u}_{k,h}^i \mid \mathbf{D}_{h-1})}{\mathbb{P}((\xi_{k,h}^i)' = \alpha^i - (\Delta\tilde{u}_{k,h}^i)' \mid \mathbf{D}'_{h-1})} = \prod_{j=1}^{nd} \frac{1/2a}{1/2a}$$

$$= 1 = \exp(0)$$

So, for each agent  $i \in N$ , we have with probability at least  $1 - \delta/(2nH)$ , the following

$$\mathbb{P}(\xi_{k,h}^i = \alpha^i - \Delta \tilde{u}_{k,h}^i \mid \mathbf{D}_{h-1}) \leq \exp(0) \times \mathbb{P}((\xi_{k,h}^i)' = \alpha^i - (\Delta \tilde{u}_{k,h}^i)' \mid \mathbf{D}'_{h-1}) \quad (42)$$

Therefore, we have with probability  $1 - \delta/(2H)$  that

$$\mathbb{P}(\xi_{k,h} = \alpha - \Delta \tilde{\mathbf{u}}_{k,h} \mid \mathbf{D}_{h-1}) \leq \exp(0) \times \mathbb{P}((\xi_{k,h})' = \alpha - (\Delta \tilde{\mathbf{u}}_{k,h})' \mid \mathbf{D}'_{h-1}) \quad (43)$$

Now taking the union bound for  $\mathbf{W}_{k,h}, \xi_{k,h}$  terms and all the stages  $h \in [H]$ , with probability  $1 - \delta$ , we have

$$\log \left[ \frac{\mathbb{P}(\forall h \in [H], (\Delta \mathbf{\Lambda}_{k,h}, \Delta \mathbf{u}_{k,h}) = (\mathbf{M}, \alpha) \mid \mathbf{D}_{1:h-1})}{\mathbb{P}(\forall h \in [H], ((\Delta \mathbf{\Lambda}_{k,h})', (\Delta \mathbf{u}_{k,h})') = (\mathbf{M}, \alpha) \mid \mathbf{D}'_{1:h-1})} \right] = 0 \quad (44)$$

Therefore, Algorithm 1 preserves  $(0, \delta)$ -LDP property with the uniform mechanism.

#### D.4 $(\epsilon, 0)$ privacy for bounded Laplace mechanism (Theorem 7)

The proof of  $(\epsilon, 0)$  LDP for the bounded Laplace is exactly same as the Laplace as given in Appendix D.2, so we avoid writing it.

## E Proof of Regret Bounds

In this Section, we give the proof the upper bound on the regret of MA-LDP with different noise adding mechanisms. We first bound the regret for any noise mechanism, and then use the noise mechanism to give the explicit bound for each noise mechanism. Consider the following difference

$$\begin{aligned} V_h^{\star,i}(\mathbf{s}_{k,h}) - V_h^{\pi_{k,i}}(\mathbf{s}_{k,h}) &\stackrel{(i)}{\leq} V_{k,h}^i(\mathbf{s}_{k,h}) - V_h^{\pi_{k,i}}(\mathbf{s}_{k,h}) \\ &= \max_{\mathbf{a}} Q_{k,h}^i(\mathbf{s}_{k,h}, \mathbf{a}) - \max_{\mathbf{a}} Q_h^{\pi_{k,i}}(\mathbf{s}_{k,h}, \mathbf{a}) \\ &\stackrel{(ii)}{\leq} Q_{k,h}^i(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) - Q_h^{\pi_{k,i}}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) \\ &\stackrel{(iii)}{\leq} \bar{r}_h(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}; \mathbf{w}_{k,h}^i) + \left\langle \hat{\boldsymbol{\theta}}_{k,h}^i, \phi_{V_{k,h+1}^i}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) \right\rangle \\ &\quad + \beta_k \|\Sigma_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h})\|_2 \\ &\quad - \bar{r}_h(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}; \mathbf{w}_{k,h}^i) - \mathbb{P}_h V_{h+1}^{\pi_{k,i}}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) \\ &= \left\langle \hat{\boldsymbol{\theta}}_{k,h}^i, \phi_{V_{k,h+1}^i}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) \right\rangle + \beta_k \|\Sigma_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h})\|_2 \\ &\quad - \mathbb{P}_h V_{h+1}^{\pi_{k,i}}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) + \mathbb{P}_h V_{k,h+1}^i(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) - \mathbb{P}_h V_{k,h+1}^i(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) \\ &= \left\langle \hat{\boldsymbol{\theta}}_{k,h}^i - \boldsymbol{\theta}_h^{\star}, \phi_{V_{k,h+1}^i}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) \right\rangle + \beta_k \|\Sigma_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h})\|_2 \\ &\quad - \mathbb{P}_h V_{h+1}^{\pi_{k,i}}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) + \mathbb{P}_h V_{k,h+1}^i(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) \\ &\stackrel{(iv)}{\leq} \|\Sigma_{k,h}^{i-1/2} \hat{\boldsymbol{\theta}}_{k,h}^i - \boldsymbol{\theta}_h^{\star}\|_2 \|\Sigma_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h})\|_2 - \mathbb{P}_h V_{h+1}^{\pi_{k,i}}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) \\ &\quad + \beta_k \|\Sigma_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h})\|_2 + \mathbb{P}_h V_{k,h+1}^i(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) \\ &\stackrel{(v)}{\leq} 2\beta_k \|\Sigma_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h})\|_2 - \mathbb{P}_h V_{h+1}^{\pi_{k,i}}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) + \mathbb{P}_h V_{k,h+1}^i(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) \quad (46) \end{aligned}$$

(i) follows from the previous Lemma 2, in (ii) we replace the max over all the actions by  $\mathbf{a}_{k,h}$ . The inequality (iii) uses the update of state-action value function from line 10 of the MA-DP algorithm. In (iv) we use the Cauchy-Schwartz inequality. Finally (v) follows from the Lemma 3. Apart from the above, we also have the following

$$V_{h+1}^i(\mathbf{s}_{k,h}) - V_{h+1}^{\pi_{k,i}}(\mathbf{s}_{k,h}) \leq V_{h+1}^i(\mathbf{s}_{k,h}) \leq H \quad (47)$$

Combining the Equations (46) and (47) we have the following:

$$V_{h+1}^i(\mathbf{s}_{k,h}) - V_{h+1}^{\pi_{k,i}}(\mathbf{s}_{k,h})$$

$$\begin{aligned}
 &\leq \min\{H, 2\beta_k \|\sum_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h})\|_2 - \mathbb{P}_h V_{h+1}^{\pi_{k,i}}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) + \mathbb{P}_h V_{k,h+1}^i(\mathbf{s}_{k,h}, \mathbf{a}_{k,h})\} \\
 &\leq \min\{H, 2\beta_k \|\sum_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h})\|_2\} - \mathbb{P}_h V_{h+1}^{\pi_{k,i}}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) + \mathbb{P}_h V_{k,h+1}^i(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}),
 \end{aligned}$$

where the second inequality holds because  $V_{k,h+1}^i \geq V_{h+1}^{*,i} \geq V_{h+1}^{\pi_{k,i}}$ . Adding  $V_{h+1}^{\pi_{k,i}}(\mathbf{s}_{k,h+1}) - V_{k,h+1}^i(\mathbf{s}_{k,h+1})$  to both sides in the above equation we have the following:

$$\begin{aligned}
 &V_{h+1}^i(\mathbf{s}_{k,h}) - V_{h+1}^{\pi_{k,i}}(\mathbf{s}_{k,h}) + [V_{h+1}^{\pi_{k,i}}(\mathbf{s}_{k,h+1}) - V_{k,h+1}^i(\mathbf{s}_{k,h+1})] \\
 &\leq \min\{H, 2\beta_k \|\sum_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h})\|_2\} - \mathbb{P}_h V_{h+1}^{\pi_{k,i}}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) + \mathbb{P}_h V_{k,h+1}^i(\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) \\
 &\quad + [V_{h+1}^{\pi_{k,i}}(\mathbf{s}_{k,h+1}) - V_{k,h+1}^i(\mathbf{s}_{k,h+1})]
 \end{aligned} \tag{48}$$

Summing these inequalities for  $k = 1, 2, \dots, K$  and stages  $h = h', \dots, H$ , we have

$$\begin{aligned}
 \sum_{k=1}^K [V_{k,h'}^i(\mathbf{s}_{k,h'}) - V_{h'}^{\pi_{k,i}}(\mathbf{s}_{k,h'})] &\leq 2 \sum_{k=1}^K \sum_{h=h'}^H \beta_k \min\{1, \|\sum_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h})\|_2\} \\
 &\quad + \sum_{k=1}^K \sum_{h=h'}^H [[\mathbb{P}_h(V_{k,h+1}^i - V_{h+1}^{\pi_{k,i}})](\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) - [V_{k,h+1}^i - V_{h+1}^{\pi_{k,i}}](\mathbf{s}_{k,h+1})]
 \end{aligned} \tag{49}$$

Define the following event  $\mathcal{E}_4$  as

$$\mathcal{E}_4 = \left\{ \forall h' \in [H], \sum_{k=1}^K \sum_{h=h'}^H [[\mathbb{P}_h(V_{k,h+1}^i - V_{h+1}^{\pi_{k,i}})](\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) - [V_{k,h+1}^i - V_{h+1}^{\pi_{k,i}}](\mathbf{s}_{k,h+1})] \leq 4H \sqrt{2T \log(2H/\alpha)} \right\} \tag{50}$$

Since,  $[[\mathbb{P}_h(V_{k,h+1}^i - V_{h+1}^{\pi_{k,i}})](\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) - [V_{k,h+1}^i - V_{h+1}^{\pi_{k,i}}](\mathbf{s}_{k,h+1})]$  forms the martingale difference sequence and it is less than  $4H$ , i.e.,

$$[[\mathbb{P}_h(V_{k,h+1}^i - V_{h+1}^{\pi_{k,i}})](\mathbf{s}_{k,h}, \mathbf{a}_{k,h}) - [V_{k,h+1}^i - V_{h+1}^{\pi_{k,i}}](\mathbf{s}_{k,h+1})] \leq 4H \tag{51}$$

Applying the Azuma-Hoeffdings inequality, we have that  $\mathcal{E}_4$  holds with probability at least  $1 - \alpha/2$ . That is  $\mathbb{P}(\mathcal{E}_4) \geq 1 - \alpha/2$ . Recall,  $\Sigma \succeq \lambda I$ , and choosing  $h' = 1$  we have

$$\begin{aligned}
 \sum_{k=1}^K \sum_{h=1}^H \beta_k \min\{1, \|\sum_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h})\|_2\} &\leq \beta_K \sum_{k=1}^K \sum_{h=1}^H \min\{1, \|\sum_{k,h}^{i-1/2} \phi_{V_{k,h+1}^i}(\mathbf{s}_{k,h}, \mathbf{a}_{k,h})\|_2\} \\
 &\leq H \beta_K \sqrt{2ndK \log(1 + K/\lambda)}
 \end{aligned} \tag{52}$$

The above inequalities hold because of the Cauchy-Schwartz and the Theorem 7. Finally, on the events  $\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3, \mathcal{E}_4$ , we conclude with probability at least  $1 - \alpha$ , we have the exact expression of the regret as follows:

$$R_K \leq H \beta_K \sqrt{2ndK \log(1 + K/\lambda)} + 4H \sqrt{2T \log(2H/\alpha)}. \tag{53}$$

### E.1 Regret bound for Gaussian mechanism (Theorem 2)

To complete the regret bound for MA-LDP algorithm with the Gaussian mechanism, we substitute  $\beta_K^G$  given in Equation (20) in the regret expression given in Equation (53). Recall,

$$\beta_K^G = c_g (nd)^{3/4} H^{3/2} K^{1/4} \log(ndT/\alpha) (\log(H/\delta))^{1/4} \sqrt{1/\epsilon}. \tag{54}$$

Thus, the regret of MA-LDP algorithm with Gaussian noise adding mechanism is given by

$$\begin{aligned}
 R_K^G &\leq c_g H (nd)^{3/4} H^{3/2} K^{1/4} \log(ndT/\alpha) (\log(H/\delta))^{1/4} \sqrt{1/\epsilon} \\
 &\quad \cdot \sqrt{2ndK \log(1 + K/\lambda)} + 4H \sqrt{2T \log(2H/\alpha)} \\
 &\leq \tilde{\mathcal{O}}(n^{5/4} d^{5/4} H^{7/4} T^{3/4} \log(ndT/\alpha) (\log(H/\delta))^{1/4} \sqrt{1/\epsilon})
 \end{aligned} \tag{55}$$

This ends the proof of the regret upper bound for the Gaussian mechanism.

### E.2 Regret bound for Laplace mechanism (Theorem 4)

To complete the regret bound for MA-LDP algorithm with the Laplace mechanism, we substitute  $\beta_K^L$  given in Equation (21) in the regret expression given in Equation (53). Recall,

$$\beta_K^L = c_l(nd)^{3/4} H^{3/2} K^{1/4} \log(ndT/\alpha) \sqrt{1/\epsilon}. \quad (56)$$

Thus, the regret of MA-LDP algorithm with Gaussian noise adding mechanism is given by

$$\begin{aligned} R_K^L &\leq c_l H (nd)^{3/4} H^{3/2} K^{1/4} \log(ndT/\alpha) \sqrt{1/\epsilon} \\ &\quad \cdot \sqrt{2ndK \log(1 + K/\lambda) + 4H \sqrt{2T \log(2H/\alpha)}} \\ &\leq \tilde{O}(n^{5/4} d^{5/4} H^{7/4} T^{3/4} \log(ndT/\alpha) \sqrt{1/\epsilon}) \end{aligned} \quad (57)$$

This ends the proof of the regret upper bound for the Laplace mechanism.

### E.3 Regret bound for uniform mechanism (Theorem 6)

To complete the regret bound for MA-LDP algorithm with the uniform mechanism, we substitute  $\beta_K^U$  given in Equation (22) in the regret expression given in Equation (53). Recall,

$$\beta_K^U = c_u(nd)^{3/4} H^{3/2} k^{1/4} \log(ndT/\alpha) (\log(H/\delta))^{1/4}. \quad (58)$$

Thus, the regret of MA-LDP algorithm with uniform noise adding mechanism is given by

$$\begin{aligned} R_K^L &\leq c_u H (nd)^{3/4} H^{3/2} k^{1/4} \log(ndT/\alpha) (\log(H/\delta))^{1/4} \\ &\quad \cdot \sqrt{2ndK \log(1 + K/\lambda) + 4H \sqrt{2T \log(2H/\alpha)}} \\ &\leq \tilde{O}(n^{5/4} d^{5/4} H^{7/4} T^{3/4} \log(ndT/\alpha) (\log(1/\delta))^{1/4}) \end{aligned} \quad (59)$$

This ends the proof of the regret upper bound for the uniform mechanism.

### E.4 Regret bound for bounded Laplace mechanism (Theorem 8)

To complete the regret bound for MA-LDP algorithm with the Bounded Laplace mechanism, we substitute  $\beta_K^{BL}$  given in Equation (23) in the regret expression given in Equation (53). Recall,

$$\beta_K^{BL} = c_{bl}(nd)^{3/4} \zeta^{1/4} K^{1/4} \log(ndT/\alpha). \quad (60)$$

and  $\zeta = \frac{2b^2}{1 - \exp(-\frac{B}{b})} - \kappa$ , where  $\kappa = \frac{((B+b)^2 + b^2) \times \exp(-\frac{B}{b})}{1 - \exp(-\frac{B}{b})}$ . Thus, the regret of MA-LDP algorithm with Bounded Laplace noise adding mechanism is given by

$$\begin{aligned} R_K^{BL} &\leq c_{bl} H (nd)^{3/4} \zeta^{1/4} K^{1/4} \log(ndT/\alpha) \\ &\quad \cdot \sqrt{2ndK \log(1 + K/\lambda) + 4H \sqrt{2T \log(2H/\alpha)}} \\ &\leq \tilde{O}(n^{5/4} d^{5/4} H^{1/4} T^{3/4} \zeta^{1/4} \log(ndT/\alpha)) \end{aligned} \quad (61)$$

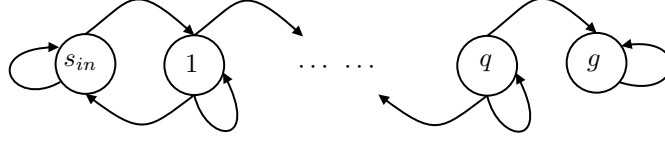
This ends the proof of the regret upper bound for the Bounded Laplace mechanism.

## F More details of experiments

Here we give more details of the experiments and the other results stated in the main paper. The following figure shows the MDP we consider in the experiments.

Recall, the feature we use in the experiments are as follows: Let  $S(s)$  is the set all feasible states from state  $s$ .

$$\phi(s'|s, \mathbf{a}) = \begin{cases} (\phi(s^1 | s^1, \mathbf{a}^1), \dots, \phi(s^n | s^n, \mathbf{a}^n)), & \text{if } s \neq \mathbf{g}, s' \in S(s) \\ \mathbf{0}_{nd}, & \text{if } s \neq \mathbf{g}, s' \notin S(s) \\ \mathbf{0}_{nd}, & \text{if } s = \mathbf{g}, s' \neq \mathbf{g} \\ (\mathbf{0}_{nd-1}, \alpha(s)), & \text{if } s = \mathbf{g}, s' = \mathbf{g}, \end{cases} \quad (62)$$


 Figure 2: Network with  $q + 2$  nodes.

where we identify  $\alpha(s)$  as  $\alpha(s) = \frac{|S(s)|}{n} \left\{ \frac{x_0}{2} + x_{q+1} + \sum_{j=1}^q \frac{x_j}{3} \right\}$ . Here  $x_0, x_1, \dots, x_q, x_{q+1}$  are the number of agents at the nodes  $s_{in}, 1, \dots, q, g$  respectively in the state  $s$ . The local features  $\phi(s^i | s^i, \mathbf{a}^i)$  are defined as

$$\phi(s^i | s^i, \mathbf{a}^i) = \begin{cases} \left( -\mathbf{a}^i, \frac{1-\delta}{n} \right)^\top, & \text{if } s^i = s^{i'} = s_{in} \\ \left( \mathbf{a}^i, \frac{\delta}{n} \right)^\top, & \text{if } s^i = s_{in}, s^{i'} = g \\ \left( -\mathbf{a}^i, \frac{1-\delta_{j,j}}{n} \right)^\top, & \text{if } s^i = s^{i'} = j \in \{1, 2, \dots, q\}, \\ \left( \mathbf{a}^i, \frac{\delta_{j,j+1}}{n} \right)^\top, & \text{if } s^i = j, s^{i'} = j+1, \forall j \in \{1, 2, \dots, q\}, \\ \left( \mathbf{0}_{d-1}, \frac{\delta_{j,j}-\delta_{j,j+1}}{n} \right)^\top, & \text{if } s^i = j, s^{i'} = j-1, \forall j \in \{1, 2, \dots, q\}, \\ \mathbf{0}_d^\top, & \text{if } s^i = g, s^{i'} = s_{in} \\ \left( \mathbf{0}_{d-1}, \frac{1}{n} \right)^\top, & \text{if } s^i = g, s^{i'} = g. \end{cases}$$

Here  $\delta_{j,j} \geq \delta_{j,j+1}$ , and  $\mathbf{0}_d^\top = (0, 0, \dots, 0)^\top$  of  $d$  dimension. Moreover, the transition probability parameters for any state  $s$  are taken as  $\theta(s) = \left( \theta^1, \frac{1}{\alpha(s)}, \theta^2, \frac{1}{\alpha(s)}, \dots, \theta^n, \frac{1}{\alpha(s)} \right)$  where  $\theta^i \in \left\{ -\frac{\Delta}{n(d-1)}, \frac{\Delta}{n(d-1)} \right\}^{d-1}$ , and  $\Delta < \delta$ .

### F.1 Proof of Lemma 4

*Proof.* We consider two cases. In case 1,  $s \neq g$  and case 2,  $s = g$ .

**Case 01:** ( $s \neq g$ ). Without loss of generality we consider the following state  $s = \underbrace{(s_{init}, \dots, s_{init})}_{x_0 \text{ times}}, \underbrace{1, 1, \dots, 1}_{x_1 \text{ times}}, \underbrace{2, 2, \dots, 2}_{x_2 \text{ times}}, \dots, \underbrace{q-1, q-1, \dots, q-1}_{x_{q-1} \text{ times}}, \underbrace{q, q, \dots, q}_{x_q \text{ times}}, \underbrace{g, g, \dots, g}_{(n-x_0-x_1-\dots-x_q) \text{ times}}$ , i.e.,  $x_0$  agents are at  $s_{init}$ ,  $x_1$  agents are at node 1,  $x_2$  agents at node 2, and so on, finally remaining  $n - x_0 - x_1 - \dots - x_q$  agents are at  $g$ . Consider an agent  $i$ , who is at  $s_{init}$  node. Let  $|S(s)|$  denotes the number of next states feasible from state  $s$ . A simple calculation shows that  $|S(s)| = 2^{x_0} \times 3^{x_1} \times 3^{x_2} \dots \times \dots \times 3^{x_q} \times 1^{n-x_0-x_1-\dots-x_q}$ . Out of these possible next states, there are exactly  $\frac{|S(s)|}{2}$  states in which agent  $i$  will remain at  $s_{init}$ , and in  $\frac{|S(s)|}{2}$  states the agent  $i$  moves to node 1. The probability that the next node of agent  $i$  is  $s_{init}$  given that the current node of agent  $i$  is  $s_{init}$  is given by  $-\langle \mathbf{a}^i, \theta^i \rangle + \frac{1-\delta}{n} \times \frac{1}{\alpha(s)}$ . And the probability that the next node of agent  $i$  is 1 given that the current node of agent  $i$  is  $s_{init}$  is  $\langle \mathbf{a}^i, \theta^i \rangle + \frac{\delta}{n} \times \frac{1}{\alpha(s)}$ . These probabilities are obtained using the features defined in Equation (62). Since, this is true for all the agents  $1, 2, \dots, x_0$  which are at  $s_{init}$ . So, the contribution to the probability term from these  $x_0$  agents who are at  $s_{init}$  is

$$\begin{aligned} & \sum_{i=1}^{x_0} \left\{ \left( -\langle \mathbf{a}^i, \theta^i \rangle + \frac{1-\delta}{n} \times \frac{1}{\alpha(s)} \right) \times \frac{|S(s)|}{2} \right\} + \sum_{i=1}^{x_0} \left\{ \left( \langle \mathbf{a}^i, \theta^i \rangle + \frac{\delta}{n} \times \frac{1}{\alpha(s)} \right) \times \frac{|S(s)|}{2} \right\} \\ & = \frac{|S(s)|}{2} \times \frac{1}{n} \times \frac{x_0}{\alpha(s)}. \end{aligned} \quad (63)$$

Next consider an agent  $i$  who is at node  $j \in \{1, 2, \dots, q\}$ . Out of next possible states, the number of next possible states where agent  $i$  will remain at node  $j$  is  $\frac{|S(s)|}{3}$ , move to node  $j+1$  is  $\frac{|S(s)|}{3}$  states and moves to node  $j-1$  is  $\frac{|S(s)|}{3}$ . The probability of staying at node  $j$  is  $-\langle \mathbf{a}^i, \theta^i \rangle + \frac{1-\delta_{j,j}}{n} \times \frac{1}{\alpha(s)}$ ; moving to node  $j+1$  is  $\langle \mathbf{a}^i, \theta^i \rangle + \frac{\delta_{j,j+1}}{n} \times \frac{1}{\alpha(s)}$ ; and probability of going to node  $j-1$  is  $\frac{\delta_{j,j}-\delta_{j,j+1}}{n} \times \frac{1}{\alpha(s)}$ . This is true for all the agents who are at node  $j$  in the state  $s$ .

Therefore, the contribution in the overall probability from this agent is

$$\begin{aligned} & \sum_{i=1}^{x_j} \left\{ \left( -\langle \mathbf{a}^i, \boldsymbol{\theta}^i \rangle + \frac{1 - \delta_{j,j}}{n} \times \frac{1}{\alpha(\mathbf{s})} \right) \times \frac{|S(\mathbf{s})|}{3} \right\} + \sum_{i=1}^{x_j} \left\{ \left( \langle \mathbf{a}^i, \boldsymbol{\theta}^i \rangle + \frac{\delta_{j,j+1}}{n} \times \frac{1}{\alpha(\mathbf{s})} \right) \times \frac{|S(\mathbf{s})|}{3} \right\} \\ & + \sum_{i=1}^{x_j} \left( \frac{\delta_{j,j} - \delta_{j,j+1}}{n} \times \frac{1}{\alpha(\mathbf{s})} \right) \times \frac{|S(\mathbf{s})|}{3} = \frac{|S(\mathbf{s})|}{3} \times \frac{1}{n} \times \frac{x_j}{\alpha(\mathbf{s})}. \end{aligned} \quad (64)$$

The above expression is valid for any node  $j \in \{1, 2, \dots, q\}$ . Finally consider the agent who is at node  $g$ , the number of next states in which the agent stays at node  $g$  is  $|S(\mathbf{s})|$ . Let  $x_{q+1} = n - x_0 - x_1 - \dots - x_q$ . The probability of this is  $\frac{1}{n} \frac{x_{q+1}}{\alpha(\mathbf{s})}$ , so Therefore, the contribution in the probability from the agent who is at node  $g$  is

$$\sum_{i=1}^{x_{q+1}} |S(\mathbf{s})| \times \frac{1}{n} \times \frac{1}{\alpha(\mathbf{s})} = |S(\mathbf{s})| \times \frac{1}{n} \times \frac{x_{q+1}}{\alpha(\mathbf{s})} \quad (65)$$

Adding Equations (63), (64) and (65), we have

$$\begin{aligned} \sum_{s' \neq g} \langle \phi(s'|s, \mathbf{a}), \boldsymbol{\theta}(s) \rangle &= \left( \frac{|S(\mathbf{s})|}{2} \times \frac{1}{n} \times \frac{x_0}{\alpha(\mathbf{s})} \right) + \sum_{j=1}^q \left( \frac{|S(\mathbf{s})|}{3} \times \frac{1}{n} \times \frac{x_j}{\alpha(\mathbf{s})} \right) + \left( |S(\mathbf{s})| \times \frac{1}{n} \times \frac{x_{q+1}}{\alpha(\mathbf{s})} \right) \\ &= \frac{|S(\mathbf{s})|}{n\alpha(\mathbf{s})} \left\{ \frac{x_0}{2} + x_{q+1} + \sum_{j=1}^q \frac{x_j}{3} \right\} \end{aligned}$$

Since, we set  $\alpha(\mathbf{s}) = \frac{|S(\mathbf{s})|}{n} \left\{ \frac{x_0}{2} + x_{q+1} + \sum_{j=1}^q \frac{x_j}{3} \right\}$ , we have that the above summation as 1.

**Case 02:** ( $s = g$ ). For this case, the probability is

$$\begin{aligned} \sum_{s'} \langle \phi(s'|s = \mathbf{g}, \mathbf{a}), \boldsymbol{\theta}(s) \rangle &= \sum_{s' \neq \mathbf{g}} \langle \phi(s'|s = \mathbf{g}, \mathbf{a}), \boldsymbol{\theta}(s) \rangle + \langle \phi(s' = \mathbf{g}|s = \mathbf{g}, \mathbf{a}), \boldsymbol{\theta}(s) \rangle \\ &= \langle \mathbf{0}, \boldsymbol{\theta}(s) \rangle + \langle (\mathbf{0}_{nd-1}, \alpha(\mathbf{s})), \boldsymbol{\theta}(s) \rangle = 1 \end{aligned}$$

Therefore, in both cases, we have

$$\sum_{s'} \langle \phi(s'|s = \mathbf{g}, \mathbf{a}), \boldsymbol{\theta}(s) \rangle = 1, \quad \forall s, \mathbf{a}.$$

The other two statements of the Lemma follow by feature design and model parameter space.  $\square$

## G Some useful results

### G.1 Equivalence of the optimization problems

To start with we show the equivalence of the optimization problems we obtain from the least square minimizer of the global reward function. Recall the optimization problem is

$$\min_{\mathbf{w}} \mathbb{E}_{s, \mathbf{a}} [\bar{r}(s, \mathbf{a}) - \bar{r}(s, \mathbf{a}; \mathbf{w})]^2. \quad (\text{OP 1})$$

We prove the following key Proposition which enables the decentralized working of our algorithm.

**Proposition 3** (Zhang et al. (2018), Trivedi and Hemachandra (2022)). *The optimization problem in Eq. (OP 1) is equivalently characterized as (both have the same stationary points)*

$$\min_{\mathbf{w}} \sum_{i=1}^n \mathbb{E}_{s, \mathbf{a}} [r^i(s, \mathbf{a}) - \bar{r}(s, \mathbf{a}; \mathbf{w})]^2. \quad (\text{OP 2})$$

*Proof.* Taking the first order derivative of the objective function in optimization problem (OP 1) w.r.t.  $\mathbf{w}$ , we have:

$$\begin{aligned} & -2 \times \mathbb{E}_{s, \mathbf{a}} [\bar{r}(s, \mathbf{a}) - \bar{r}(s, \mathbf{a}; \mathbf{w})] \times \nabla_{\mathbf{w}} \bar{r}(s, \mathbf{a}; \mathbf{w}) \\ &= -2 \times \mathbb{E}_{s, \mathbf{a}} \left[ \frac{1}{n} \sum_{i \in N} r^i(s, \mathbf{a}) - \bar{r}(s, \mathbf{a}; \mathbf{w}) \right] \times \nabla_{\mathbf{w}} \bar{r}(s, \mathbf{a}; \mathbf{w}), \end{aligned}$$



$$\begin{aligned}
 &= -\frac{2}{n} \times \mathbb{E}_{s,\mathbf{a}} \left[ \sum_{i \in N} r^i(s, \mathbf{a}) - n \cdot \bar{r}(s, \mathbf{a}; \mathbf{w}) \right] \times \nabla_{\mathbf{w}} \bar{r}(s, \mathbf{a}; \mathbf{w}), \\
 &= -\frac{2}{n} \times \mathbb{E}_{s,\mathbf{a}} \left[ \sum_{i \in N} (r^i(s, \mathbf{a}) - \bar{r}(s, \mathbf{a}; \mathbf{w})) \right] \times \nabla_{\mathbf{w}} \bar{r}(s, \mathbf{a}; \mathbf{w}).
 \end{aligned}$$

Ignoring the factor  $\frac{1}{n}$  in the above equation, we exactly have the first order derivative of the objective function in [OP 2](#). Thus, both optimization problems have the same stationary points. Hence, [OP 1](#) is an *equivalent characterization* of [OP 2](#).  $\square$

## G.2 Lemma for Gaussian mechanism

For the Gaussian mechanism, we can only hope for the  $(\epsilon, \delta)$  LDP. In this subsection, we show that our MA-LDP algorithm with Gaussian mechanism indeed preserves the differential privacy. To this end, we have following Theorem (Theorem A.2 in [Liao et al. \(2021\)](#))

**Theorem 12** (Theorem A.2 [Liao et al. \(2021\)](#)). *If the privacy loss  $c$  satisfy  $\mathbb{P}_{o \sim \mathcal{M}(d)}[c(o; \mathcal{M}, \mathbf{aux}, d, d') > \epsilon] \leq \delta$  for all auxiliary input  $\mathbf{aux}$  and neighboring data sets  $d, d'$ , then the mechanism  $\mathcal{M}$  satisfies  $(\epsilon, \delta)$ -LDP property.*

The basic idea involved in the above Theorem is that if we set the privacy parameter  $\delta$  then for any auxiliary input  $\mathbf{aux}$  the probability that any outcome  $o$  obtained using the privacy preserving mechanism incurs at least  $\epsilon$  privacy loss then the mechanism  $\mathcal{M}$  satisfies the  $(\epsilon, \delta)$ -LDP. The proof of this Theorem is available in [Liao et al. \(2021\)](#); [Abadi et al. \(2016\)](#) and is based on the construction of an event containing all possible outcomes for which the absolute privacy loss is at least  $\epsilon$ .

**Lemma 5** (Gaussian Mechanism [Dwork et al. \(2006\)](#); [Liao et al. \(2021\)](#)). *Let  $f : \mathcal{N}^{\mathcal{X}} \rightarrow \mathbb{R}^d$  be an arbitrary  $d$ -dimensional function (a query), and define the  $l_2$  sensitivity as  $\Delta_2 f = \max_{adj(x,y)} \|f(x) - f(y)\|_2$ , where  $adj(x, y)$  indicates that  $x, y$  are different at one entry only. For any  $0 \leq \epsilon \leq 1$  and  $c^2 > 2 \log(1.25/\delta)$ , the Gaussian mechanism with parameter  $\sigma \geq c\Delta_2 f/\epsilon$  is  $(\epsilon, \delta)$ -LDP.*

The following Analogous Lemma for the Laplace Mechanism is also available in Theorem 3.6 of [Dwork and Roth \(2014\)](#).

**Lemma 6** (Laplace Mechanism; Theorem 3.6 [Dwork and Roth \(2014\)](#)). *Let  $f : \mathcal{N}^{\mathcal{X}} \rightarrow \mathbb{R}^d$  be an arbitrary  $d$ -dimensional function (a query), and define the  $l_1$  sensitivity as  $\Delta f = \max_{\|x-y\|_1=1} \|f(x) - f(y)\|_1$ . For any  $0 \leq \epsilon \leq 1$  the Laplace mechanism with parameter  $b = \frac{\Delta f}{\epsilon}$  preserves  $(\epsilon, 0)$  differential privacy.*

## G.3 Other important results

**Lemma 7** (Lemma 11, [Abbasi-Yadkori et al. \(2011\)](#)). *Let  $\{\phi_t\}_{t \geq 0}$  be the bounded sequence in  $\mathbb{R}^d$  satisfying  $\sup_{t \geq 0} \|\phi_t\| \leq 1$ . Let  $\Lambda_0 \in \mathbb{R}^{d \times d}$  be a positive definite matrix. For any  $t \geq 0$ , we define  $\Lambda_t = \Lambda_0 + \sum_{j=0}^t \phi_j^\top \phi_j$ . Then, if the smallest eigenvalue of  $\Lambda_0$  satisfies  $\lambda_{\min}(\Lambda) \geq 1$ , we have*

$$\log \left[ \frac{\det(\Lambda_t)}{\det(\Lambda_0)} \right] \leq \sum_{j=1}^t \phi_j^\top \Lambda_{j-1}^{-1} \phi_j \leq 2 \log \left[ \frac{\det(\Lambda_t)}{\det(\Lambda_0)} \right] \quad (66)$$

**Theorem 13** (Theorem 2 of [Zhou et al. \(2021\)](#)). *Let  $\{\mathcal{G}_t\}_{t=1}^\infty$  be the filtration. Let  $\{x_t, \eta_t\}_{t \geq 1}$  be a stochastic process so that  $x_t \in \mathbb{R}^d$  is  $\mathcal{G}_t$ -measurable and  $\eta_t$  be  $\mathcal{G}_{t+1}$  measurable. Fix,  $R, L, \sigma, \lambda, \mu^* \in \mathbb{R}^d$ . For  $t \geq 1$  let  $y_t = \langle \mu^*, x_t \rangle + \eta_t$  and suppose that  $\eta_t, x_t$  also satisfy*

$$\|\eta_t\| \leq R, \mathbb{E}[\eta_t^2 | \mathcal{G}_t] \leq \sigma^2, \|x_t\|_2 \leq L \quad (67)$$

Then, for any  $0 < \delta < 1$ , with probability at least  $1 - \delta$  we have

$$\forall t \geq 0, \left\| \sum_{i=1}^t x_i \eta_i \right\|_{Z_t^{-1}} \leq \beta_t, \|\mu_t - \mu^*\|_{Z_t} \leq \beta_t + \sqrt{\lambda} \|\mu^*\|_2, \quad (68)$$

where for  $t \geq 1$ ,  $\mu_t = Z_t^{-1} b_t$ ,  $Z_t = \lambda I + \sum_{i=1}^t x_i x_i^\top$ ,  $b_t = \sum_{i=1}^t y_i x_i$  and

$$\beta_t = 8\sigma \sqrt{d \log(1 + tL^2/(d\lambda)) \log(4t^2/\delta)} + 4R \log(4t^2/\delta) \quad (69)$$

**Lemma 8** (Kushner-Clark Lemma [Kushner and Yin \(2003\)](#); [Metivier and Priouret \(1984\)](#)). *Let  $\mathcal{X} \subseteq \mathbb{R}^p$  be a compact set and let  $h : \mathcal{X} \rightarrow \mathbb{R}^p$  be a continuous function. Consider the following recursion in  $p$ -dimensions*

$$x_{t+1} = \Gamma\{x_t + \gamma_t[h(x_t) + \zeta_t + \beta_t]\}. \quad (70)$$

Let  $\hat{\Gamma}(\cdot)$  be transformed projection operator defined for any  $x \in \mathcal{X} \subseteq \mathbb{R}^p$  as

$$\hat{\Gamma}(h(x)) = \lim_{0 < \eta \rightarrow 0} \left\{ \frac{\Gamma(x + \eta h(x)) - x}{\eta} \right\},$$

then the ODE associated with Equation (70) is  $\dot{x} = \hat{\Gamma}(h(x))$ .

**Assumption 4.** *Kushner-Clark lemma requires the following assumptions*

1. *Stepsize  $\{\gamma_t\}_{t \geq 0}$  satisfy  $\sum_t \gamma_t = \infty$ , and  $\gamma_t \rightarrow 0$  as  $t \rightarrow \infty$ .*
2. *The sequence  $\{\beta_t\}_{t \geq 0}$  is a bounded random sequence with  $\beta_t \rightarrow 0$  almost surely as  $t \rightarrow \infty$ .*
3. *For any  $\epsilon > 0$ , the sequence  $\{\zeta_t\}_{t \geq 0}$  satisfy*

$$\lim_t \mathbb{P} \left( \sup_{p \geq t} \left\| \sum_{\tau=t}^p \gamma_\tau \zeta_\tau \right\| \geq \epsilon \right) = 0.$$

*Kushner-Clark lemma is as follows: suppose that ODE  $\dot{x} = \hat{\Gamma}(h(x))$  has a compact set  $\mathcal{K}^*$  as its asymptotically stable equilibria, then under Assumption 4,  $x_t$  in Equation (70) converges almost surely to  $\mathcal{K}^*$  as  $t \rightarrow \infty$ .*

#### G.4 Generating a bounded Laplace distribution, $\mathcal{BL}$ ([Ross, 2022](#))

The cdf of the bounded Laplace random variable,  $\mathcal{BL}$  can be obtained as follows:

$$F_{\mathcal{BL}}(x) = \int_{-B}^x \frac{1}{2b(1 - \exp(\frac{-B}{b}))} \exp\left(\frac{-|x|}{b}\right) \quad (71)$$

$$= \int_{-B}^0 \frac{1}{2b(1 - \exp(\frac{-B}{b}))} \exp\left(\frac{x}{b}\right) + \int_0^x \frac{1}{2b(1 - \exp(\frac{-B}{b}))} \exp\left(\frac{-x}{b}\right) \quad (72)$$

$$= \frac{1}{2b(1 - \exp(\frac{-B}{b}))} [b - b \exp(\frac{-B}{b})] - \frac{1}{2b(1 - \exp(\frac{-B}{b}))} [b \exp(\frac{-x}{b}) - b] \quad (73)$$

$$= \frac{2 - \exp(\frac{-B}{b}) - \exp(\frac{-x}{b})}{2(1 - \exp(\frac{-B}{b}))} \quad (74)$$

To simulate this, let  $u = F_{\mathcal{BL}}(x)$ , where  $u \sim \mathcal{U}(0, 1)$ , then

$$u = \frac{2 - \exp(\frac{-B}{b}) - \exp(\frac{-x}{b})}{2(1 - \exp(\frac{-B}{b}))} \quad (75)$$

$$\exp(\frac{-x}{b}) = 2 - \exp(\frac{-B}{b}) - 2u \left(1 - \exp(\frac{-B}{b})\right) \quad (76)$$

$$\frac{-x}{b} = \log \left( 2 - \exp(\frac{-B}{b}) - 2u \left(1 - \exp(\frac{-B}{b})\right) \right) \quad (77)$$

$$x = -b \log \left( 2 - \exp(\frac{-B}{b}) - 2u \left(1 - \exp(\frac{-B}{b})\right) \right) \quad (78)$$

Note, this  $x$  will follow the Laplace distribution with bounded support.