
Offline Reinforcement Learning with Differential Privacy

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The offline reinforcement learning (RL) problem is often motivated by the need to
2 learn data-driven decision policies in financial, legal and healthcare applications.
3 However, the learned policy could retain sensitive information of individuals in the
4 training data (e.g., treatment and outcome of patients), thus susceptible to various
5 privacy risks. We design offline RL algorithms with differential privacy guarantees
6 which provably prevent such risks. These algorithms also enjoy strong instance-
7 dependent learning bounds under both tabular and linear Markov Decision Process
8 (MDP) settings. Our theory and simulation suggest that the privacy guarantee
9 comes at (almost) no drop in utility comparing to the non-private counterpart for a
10 medium-size dataset.

11 1 Introduction

12 Offline Reinforcement Learning (or batch RL) aims to learn a near-optimal policy in an unknown
13 environment¹ through a static dataset gathered from some behavior policy μ . Since offline RL
14 does not require access to the environment, it can be applied to problems where interaction with
15 environment is infeasible, e.g., when collecting new data is costly (trade or finance [Zhang et al.,
16 2020]), risky (autonomous driving [Sallab et al., 2017]) or illegal / unethical (healthcare [Raghu
17 et al., 2017]). In such practical applications, the data used by an RL agent usually contains sensitive
18 information. Take medical history for instance, for each patient, at each time step, the patient reports
19 her health condition (age, disease, etc.), then the doctor decides the treatment (which medicine to use,
20 the dosage of medicine, etc.), finally there is treatment outcome (whether the patient feels good, etc.)
21 and the patient transitions to another health condition. Here, (health condition, treatment, treatment
22 outcome) corresponds to (state, action, reward) and the dataset can be considered as n (number of
23 patients) trajectories sampled from a MDP with horizon H (number of treatment steps). However,
24 learning agents are known to implicitly memorize details of individual training data points verbatim
25 [Carlini et al., 2019], even if they are irrelevant for learning [Brown et al., 2021], which makes offline
26 RL models vulnerable to various privacy attacks.

27 Differential privacy (DP) [Dwork et al., 2006] is a well-established definition of privacy with many
28 desirable properties. A differentially private offline RL algorithm will return a decision policy that
29 is indistinguishable from a policy trained in an alternative universe any individual user is replaced,
30 thereby preventing the aforementioned privacy risks. There is a surge of recent interest in developing
31 RL algorithms with DP guarantees, but they focus mostly on the online setting [Vietri et al., 2020,
32 Garcelon et al., 2021, Liao et al., 2021, Chowdhury and Zhou, 2021, Luyo et al., 2021].

33 Offline RL is arguably more practically relevant than online RL in the applications with sensitive data.
34 For example, in the healthcare domain, online RL requires actively running new exploratory policies

¹The environment is usually characterized by a Markov Decision Process (MDP) in this paper.

35 (clinical trials) with every new patient, which often involves complex ethical / legal clearances,
 36 whereas offline RL uses only historical patient records that are often accessible for research purposes.
 37 Clear communication of the adopted privacy enhancing techniques (e.g., DP) to patients was reported
 38 to further improve data access [Kim et al., 2017].

39 **Our contributions.** In this paper, we present the first provably efficient algorithms for offline RL
 40 with differential privacy. Our contributions are twofold.

- 41 • We design two new pessimism-based algorithms DP-APVI (Algorithm 1) and DP-VAPVI
 42 (Algorithm 2), one for the tabular setting (finite states and actions), the other for the case
 43 with linear function approximation (under linear MDP assumption). Both algorithms enjoy
 44 DP guarantees (pure DP or zCDP) and instance-dependent learning bounds where the cost
 45 of privacy appears as lower order terms.
- 46 • We perform numerical simulations to evaluate and compare the performance of our algorithm
 47 DP-VAPVI (Algorithm 2) with its non-private counterpart VAPVI [Yin et al., 2022] as well
 48 as a popular baseline PEVI [Jin et al., 2021]. The results complement the theoretical findings
 49 by demonstrating the practicality of DP-VAPVI under strong privacy parameters.

50 **Related work.** To our knowledge, differential privacy in offline RL tasks has not been studied before,
 51 except for much simpler cases where the agent only evaluates a single policy [Balle et al., 2016, Xie
 52 et al., 2019]. Balle et al. [2016] privatized first-visit Monte Carlo-Ridge Regression estimator by an
 53 output perturbation mechanism and Xie et al. [2019] used DP-SGD. Neither paper considered offline
 54 learning (or policy optimization), which is our focus.

55 There is a larger body of work on private RL in the online setting, where the goal is to minimize regret
 56 while satisfying either joint differential privacy [Vietri et al., 2020, Chowdhury and Zhou, 2021, Ngo
 57 et al., 2022, Luyo et al., 2021] or local differential privacy [Garcelon et al., 2021, Liao et al., 2021,
 58 Luyo et al., 2021, Chowdhury and Zhou, 2021]. The offline setting introduces new challenges in DP
 59 as we cannot *algorithmically enforce* good “exploration”, but have to work with a static dataset and
 60 privately estimate the uncertainty in addition to the value functions. A private online RL algorithm
 61 can sometimes be adapted for private offline RL too, but those from existing work yield suboptimal
 62 and non-adaptive bounds. We give a more detailed technical comparison in Appendix B.

63 Among non-private offline RL works, we build directly upon non-private offline RL methods known
 64 as Adaptive Pessimistic Value Iteration (APVI, for tabular MDPs) [Yin and Wang, 2021b] and
 65 Variance-Aware Pessimistic Value Iteration (VAPVI, for linear MDPs) [Yin et al., 2022], as they give
 66 the strongest theoretical guarantees to date. We refer readers to Appendix B for a more extensive
 67 review of the offline RL literature. Introducing DP to APVI and VAPVI while retaining the same
 68 sample complexity (modulo lower order terms) require nontrivial modifications to the algorithms.

69 **A remark on technical novelty.** Our algorithms involve substantial technical innovation over
 70 previous works on online DP-RL with joint DP guarantee². Different from previous works, our
 71 DP-APVI (Algorithm 1) operates on Bernstein type pessimism, which requires our algorithm to deal
 72 with conditional variance using private statistics. Besides, our DP-VAPVI (Algorithm 2) replaces the
 73 LSVI technique with variance-aware LSVI (also known as weighted ridge regression, first appears
 74 in [Zhou et al., 2021]). Our DP-VAPVI releases conditional variance privately, and further applies
 75 weighted ridge regression privately. Both approaches ensure tighter instance-dependent bounds on
 76 the suboptimality of the learned policy.

77 2 Problem Setup

78 **Markov Decision Process.** A finite-horizon *Markov Decision Process* (MDP) is denoted by a tuple
 79 $M = (\mathcal{S}, \mathcal{A}, P, r, H, d_1)$ [Sutton and Barto, 2018], where \mathcal{S} is state space and \mathcal{A} is action space. A
 80 non-stationary transition kernel $P_h : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \mapsto [0, 1]$ maps each state action (s_h, a_h) to a probabil-
 81 ity distribution $P_h(\cdot | s_h, a_h)$ and P_h can be different across time. Besides, $r_h : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$ is the ex-
 82 pected immediate reward satisfying $0 \leq r_h \leq 1$, d_1 is the initial state distribution and H is the horizon.
 83 A policy $\pi = (\pi_1, \dots, \pi_H)$ assigns each state $s_h \in \mathcal{S}$ a probability distribution over actions accord-
 84 ing to the map $s_h \mapsto \pi_h(\cdot | s_h)$, $\forall h \in [H]$. A random trajectory $s_1, a_1, r_1, \dots, s_H, a_H, r_H, s_{H+1}$ is
 85 generated according to $s_1 \sim d_1, a_h \sim \pi_h(\cdot | s_h), r_h \sim r_h(s_h, a_h), s_{h+1} \sim P_h(\cdot | s_h, a_h), \forall h \in [H]$.

²Here we only compare our techniques (for offline RL) with the works for online RL under joint DP guarantee, as both settings allow access to the raw data.

86 For tabular MDP, we have $\mathcal{S} \times \mathcal{A}$ is the discrete state-action space and $S := |\mathcal{S}|, A := |\mathcal{A}|$ are finite.
 87 In this work, we assume that r is known³. In addition, we denote the per-step marginal state-action
 88 occupancy $d_h^\pi(s, a)$ as: $d_h^\pi(s, a) := \mathbb{P}[s_h = s | s_1 \sim d_1, \pi] \cdot \pi_h(a | s)$, which is the marginal state-action
 89 probability at time h .

90 **Value function, Bellman (optimality) equations.** The value function $V_h^\pi(\cdot)$ and Q-value func-
 91 tion $Q_h^\pi(\cdot, \cdot)$ for any policy π is defined as: $V_h^\pi(s) = \mathbb{E}_\pi[\sum_{t=h}^H r_t | s_h = s]$, $Q_h^\pi(s, a) =$
 92 $\mathbb{E}_\pi[\sum_{t=h}^H r_t | s_h, a_h = s, a]$, $\forall h, s, a \in [H] \times \mathcal{S} \times \mathcal{A}$. The performance is defined as $v^\pi :=$
 93 $\mathbb{E}_{d_1}[V_1^\pi] = \mathbb{E}_{\pi, d_1}[\sum_{t=1}^H r_t]$. The Bellman (optimality) equations follow $\forall h \in [H]: Q_h^\pi =$
 94 $r_h + P_h V_{h+1}^\pi$, $V_h^\pi = \mathbb{E}_{a \sim \pi_h}[Q_h^\pi]$, $Q_h^* = r_h + P_h V_{h+1}^*$, $V_h^* = \max_a Q_h^*(\cdot, a)$.

95 **Linear MDP [Jin et al., 2020b].** An episodic MDP $(\mathcal{S}, \mathcal{A}, H, P, r)$ is called a linear MDP with
 96 known feature map $\phi: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ if there exist H unknown signed measures $\nu_h \in \mathbb{R}^d$ over \mathcal{S} and
 97 H unknown reward vectors $\theta_h \in \mathbb{R}^d$ such that

$$P_h(s' | s, a) = \langle \phi(s, a), \nu_h(s') \rangle, \quad r_h(s, a) = \langle \phi(s, a), \theta_h \rangle, \quad \forall (h, s, a, s') \in [H] \times \mathcal{S} \times \mathcal{A} \times \mathcal{S}.$$

98 Without loss of generality, we assume $\|\phi(s, a)\|_2 \leq 1$ and $\max(\|\nu_h(\mathcal{S})\|_2, \|\theta_h\|_2) \leq \sqrt{d}$ for all
 99 $h, s, a \in [H] \times \mathcal{S} \times \mathcal{A}$. An important property of linear MDP is that the value functions are linear in
 100 the feature map, which is summarized in Lemma F.14.

101 **Offline setting and the goal.** The offline RL requires the agent to find a policy π in order to maximize
 102 the performance v^π , given only the episodic data $\mathcal{D} = \{(s_h^\tau, a_h^\tau, r_h^\tau, s_{h+1}^\tau)\}_{\tau \in [n]}^{h \in [H]}$ ⁴ rolled out from some
 103 fixed and possibly unknown behavior policy μ , which means we cannot change μ and in particular
 104 we do not assume the functional knowledge of μ . In conclusion, based on the batch data \mathcal{D} and a
 105 targeted accuracy $\epsilon > 0$, the agent seeks to find a policy π_{alg} such that $v^* - v^{\pi_{\text{alg}}} \leq \epsilon$.

106 2.1 Assumptions in offline RL

107 In order to show that our privacy-preserving algorithms can generate near optimal policy, certain
 108 coverage assumptions are needed. In this section, we will list the assumptions we use in this paper.

109 Assumptions for tabular setting.

110 **Assumption 2.1 ([Liu et al., 2019]).** *There exists one optimal policy π^* , such that π^* is fully covered*
 111 *by μ , i.e. $\forall s_h, a_h \in \mathcal{S} \times \mathcal{A}, d_h^{\pi^*}(s_h, a_h) > 0$ only if $d_h^\mu(s_h, a_h) > 0$. Furthermore, we denote the*
 112 *trackable set as $\mathcal{C}_h := \{(s_h, a_h) : d_h^\mu(s_h, a_h) > 0\}$.*

113 Assumption 2.1 is the weakest assumption needed for accurately learning the optimal value v^* by
 114 requiring μ to trace the state-action space of one optimal policy (μ can be agnostic at other locations).
 115 Similar to [Yin and Wang, 2021b], we will use Assumption 2.1 for the tabular part of this paper,
 116 which enables comparison between our sample complexity to the conclusion in [Yin and Wang,
 117 2021b], whose algorithm serves as a non-private baseline.

118 **Assumptions for linear setting.** First, we define the expectation of covariance matrix under the
 119 behavior policy μ for all time step $h \in [H]$ as below:

$$\Sigma_h^p := \mathbb{E}_\mu [\phi(s_h, a_h) \phi(s_h, a_h)^\top]. \quad (1)$$

120 As have been shown in [Wang et al., 2021, Yin et al., 2022], learning a near-optimal policy from
 121 offline data requires coverage assumptions. Here in linear setting, such coverage is characterized by
 122 the minimum eigenvalue of Σ_h^p . Similar to [Yin et al., 2022], we apply the following assumption for
 123 the sake of comparison.

124 **Assumption 2.2 (Feature Coverage, Assumption 2 in [Wang et al., 2021]).** *The data distributions*
 125 *μ satisfy the minimum eigenvalue condition: $\forall h \in [H], \kappa_h := \lambda_{\min}(\Sigma_h^p) > 0$. Furthermore, we*
 126 *denote $\kappa = \min_h \kappa_h$.*

127 2.2 Differential Privacy in offline RL

128 In this work, we aim to design privacy-preserving algorithms for offline RL. We apply differential
 129 privacy as the formal notion of privacy. Below we revisit the definition of differential privacy.

³This is due to the fact that the uncertainty of reward function is dominated by that of transition kernel in RL.

⁴For clarity we use n for tabular MDP and K for linear MDP when referring to the sample complexity.

Definition 2.3 (Differential Privacy [Dwork et al., 2006]). A randomized mechanism M satisfies (ϵ, δ) -differential privacy (ϵ, δ) -DP if for all neighboring datasets U, U' that differ by one data point and for all possible event E in the output range, it holds that

$$\mathbb{P}[M(U) \in E] \leq e^\epsilon \cdot \mathbb{P}[M(U') \in E] + \delta.$$

130 When $\delta = 0$, we say pure DP, while for $\delta > 0$, we say approximate DP.

131 In the problem of offline RL, the dataset consists of several trajectories, therefore one data point in
 132 Definition 2.3 refers to one single trajectory. Hence the definition of Differential Privacy means that
 133 the difference in the distribution of the output policy resulting from replacing one trajectory in the
 134 dataset will be small. In other words, an adversary can not infer much information about any single
 135 trajectory in the dataset from the output policy of the algorithm. For more discussions about our
 136 definition of DP, please refer to Appendix C.1.

137 During the whole paper, we will use zCDP (defined below) as a surrogate for DP, since it enables
 138 cleaner analysis for privacy composition and Gaussian mechanism. The properties of zCDP (e.g.,
 139 composition, conversion formula to DP) are deferred to Appendix F.3.

Definition 2.4 (zCDP [Dwork and Rothblum, 2016, Bun and Steinke, 2016]). A randomized mechanism M satisfies ρ -Zero-Concentrated Differential Privacy (ρ -zCDP), if for all neighboring datasets U, U' and all $\alpha \in (1, \infty)$,

$$D_\alpha(M(U) \| M(U')) \leq \rho\alpha,$$

140 where D_α is the Renyi-divergence [Van Erven and Harremos, 2014].

141 Finally, we go over the definition and privacy guarantee of Gaussian mechanism.

142 **Definition 2.5** (Gaussian Mechanism [Dwork et al., 2014]). Define the ℓ_2 sensitivity of a function
 143 $f : \mathbb{N}^{\mathcal{X}} \mapsto \mathbb{R}^d$ as

$$\Delta_2(f) = \sup_{\text{neighboring } U, U'} \|f(U) - f(U')\|_2.$$

144 The Gaussian mechanism \mathcal{M} with noise level σ is then given by

$$\mathcal{M}(U) = f(U) + \mathcal{N}(0, \sigma^2 I_d).$$

145 **Lemma 2.6** (Privacy guarantee of Gaussian mechanism [Dwork et al., 2014, Bun and Steinke, 2016]).

146 Let $f : \mathbb{N}^{\mathcal{X}} \mapsto \mathbb{R}^d$ be an arbitrary d -dimensional function with ℓ_2 sensitivity Δ_2 . Then for any $\rho > 0$,

147 Gaussian Mechanism with parameter $\sigma^2 = \frac{\Delta_2^2}{2\rho}$ satisfies ρ -zCDP. In addition, for all $0 < \delta, \epsilon < 1$,

148 Gaussian Mechanism with parameter $\sigma = \frac{\Delta_2}{\epsilon} \sqrt{2 \log \frac{1.25}{\delta}}$ satisfies (ϵ, δ) -DP.

149 We emphasize that the privacy guarantee covers any input data. It does *not* require any distributional
 150 assumptions on the data. The RL-specific assumptions (e.g., linear MDP and coverage assumptions)
 151 are only used for establishing provable utility guarantees.

152 3 Results under tabular MDP: DP-APVI (Algorithm 1)

153 For reinforcement learning, the tabular MDP setting is the most well-studied setting and our first
 154 result applies to this regime. We begin with the construction of private counts.

155 **Private Model-based Components.** Given data $\mathcal{D} = \{(s_h^\tau, a_h^\tau, r_h^\tau, s_{h+1}^\tau)\}_{\tau \in [n]}^{h \in [H]}$, we denote $n_{s_h, a_h} :=$
 156 $\sum_{\tau=1}^n \mathbb{1}[s_h^\tau, a_h^\tau = s_h, a_h]$ be the total counts that visit (s_h, a_h) pair at time h and $n_{s_h, a_h, s_{h+1}} :=$
 157 $\sum_{\tau=1}^n \mathbb{1}[s_h^\tau, a_h^\tau, s_{h+1}^\tau = s_h, a_h, s_{h+1}]$ be the total counts that visit (s_h, a_h, s_{h+1}) pair at time h , then
 158 given the budget ρ for zCDP, we add *independent* Gaussian noises to all the counts:

$$n'_{s_h, a_h} = \{n_{s_h, a_h} + \mathcal{N}(0, \sigma^2)\}^+, \quad n'_{s_h, a_h, s_{h+1}} = \{n_{s_h, a_h, s_{h+1}} + \mathcal{N}(0, \sigma^2)\}^+, \quad \sigma^2 = \frac{2H}{\rho}. \quad (2)$$

159 However, after adding noise, the noisy n' counts may not satisfy $n'_{s_h, a_h} = \sum_{s_{h+1} \in \mathcal{S}} n'_{s_h, a_h, s_{h+1}}$.

160 To address this problem, we choose the private counts of visiting numbers as the solution to the

161 following optimization problem (here $E_\rho = 4\sqrt{\frac{H \log \frac{4HS^2A}{\rho}}{\rho}}$ is chosen as a high probability uniform

162 bound of the noises we add):

$$\begin{aligned} \{\tilde{n}_{s_h, a_h, s'}\}_{s' \in \mathcal{S}} &= \operatorname{argmin}_{\{x_{s'}\}_{s' \in \mathcal{S}}} \max_{s' \in \mathcal{S}} |x_{s'} - n'_{s_h, a_h, s'}| \\ \text{such that } \left| \sum_{s' \in \mathcal{S}} x_{s'} - n'_{s_h, a_h} \right| &\leq \frac{E_\rho}{2} \text{ and } x_{s'} \geq 0, \forall s' \in \mathcal{S}. \end{aligned} \quad (3)$$

$$\tilde{n}_{s_h, a_h} = \sum_{s' \in \mathcal{S}} \tilde{n}_{s_h, a_h, s'}.$$

163 **Remark 3.1** (Some explanations). *The optimization problem above serves as a post-processing step*
 164 *which will not affect the DP guarantee of our algorithm. Briefly speaking, (3) finds a set of noisy*
 165 *counts such that $\tilde{n}_{s_h, a_h} = \sum_{s' \in \mathcal{S}} \tilde{n}_{s_h, a_h, s'}$ and the estimation error for each \tilde{n}_{s_h, a_h} and $\tilde{n}_{s_h, a_h, s'}$*
 166 *is roughly E_ρ .⁵ In contrast, if we directly take the crude approach that $\tilde{n}_{s_h, a_h, s_{h+1}} = n'_{s_h, a_h, s_{h+1}}$*
 167 *and $\tilde{n}_{s_h, a_h} = \sum_{s_{h+1} \in \mathcal{S}} \tilde{n}_{s_h, a_h, s_{h+1}}$, we can only derive $|\tilde{n}_{s_h, a_h} - n_{s_h, a_h}| \leq \tilde{O}(\sqrt{S}E_\rho)$ through*
 168 *concentration on summation of S i.i.d. Gaussian noises. In conclusion, solving the optimization*
 169 *problem (3) enables tight analysis for the lower order term (the additional cost of privacy).*

170 **Remark 3.2** (Computational efficiency). *The optimization problem (3) can be reformulated as:*

$$\min t, \text{ s.t. } |x_{s'} - n'_{s_h, a_h, s'}| \leq t \text{ and } x_{s'} \geq 0 \forall s' \in \mathcal{S}, \left| \sum_{s' \in \mathcal{S}} x_{s'} - n'_{s_h, a_h} \right| \leq \frac{E_\rho}{2}. \quad (4)$$

171 *Note that (4) is a Linear Programming problem with $S + 1$ variables and $2S + 2$ linear constraints*
 172 *(one constraint on absolute value is equivalent to two linear constraints), which can be solved*
 173 *efficiently by the simplex method [Ficken, 2015] or other provably efficient algorithms [Nemhauser*
 174 *and Wolsey, 1988]. Therefore, our Algorithm 1 is computationally friendly.*

175 The private estimation of the transition kernel is defined as:

$$\tilde{P}_h(s' | s_h, a_h) = \frac{\tilde{n}_{s_h, a_h, s'}}{\tilde{n}_{s_h, a_h}}, \quad (5)$$

176 if $\tilde{n}_{s_h, a_h} > E_\rho$ and $\tilde{P}_h(s' | s_h, a_h) = \frac{1}{S}$ otherwise.

177 **Remark 3.3.** *Different from the transition kernel estimate in previous works [Vietri et al., 2020,*
 178 *Chowdhury and Zhou, 2021] that may not be a distribution, we have to ensure that ours is a*
 179 *probability distribution, because our Bernstein type pessimism (line 5 in Algorithm 1) needs to take*
 180 *variance over this transition kernel estimate. The intuition behind the construction of our private*
 181 *transition kernel is that, for those state-action pairs with $\tilde{n}_{s_h, a_h} \leq E_\rho$, we can not distinguish*
 182 *whether the non-zero private count comes from noise or actual visitation. Therefore we only take the*
 183 *empirical estimate of the state-action pairs with sufficiently large \tilde{n}_{s_h, a_h} .*

Algorithm 1 Differentially Private Adaptive Pessimistic Value Iteration (DP-APVI)

- 1: **Input:** Offline dataset $\mathcal{D} = \{(s_h^\tau, a_h^\tau, r_h^\tau, s_{h+1}^\tau)\}_{\tau, h=1}^{n, H}$. Reward function r . Constants $C_1 = \sqrt{2}$, $C_2 = 16$, $C > 1$, failure probability δ , budget for zCDP ρ .
 - 2: **Initialization:** Calculate $\tilde{n}_{s_h, a_h}, \tilde{n}_{s_h, a_h, s_{h+1}}$ as (3), $\tilde{P}_h(s_{h+1} | s_h, a_h)$ as (5). $\tilde{V}_{H+1}(\cdot) \leftarrow 0$. $E_\rho \leftarrow 4\sqrt{\frac{H \log \frac{4HS^2A}{\rho}}{\rho}}$. $\iota \leftarrow \log(HSA/\delta)$.
 - 3: **for** $h = H, H - 1, \dots, 1$ **do**
 - 4: $\tilde{Q}_h(\cdot, \cdot) \leftarrow r_h(\cdot, \cdot) + (\tilde{P}_h \cdot \tilde{V}_{h+1})(\cdot, \cdot)$
 - 5: $\forall s_h, a_h$, let $\Gamma_h(s_h, a_h) \leftarrow C_1 \sqrt{\frac{\operatorname{Var}_{\tilde{P}_{s_h, a_h}}(\tilde{V}_{h+1})^\iota}{\tilde{n}_{s_h, a_h} - E_\rho}} + \frac{C_2 S H E_\rho^\iota}{\tilde{n}_{s_h, a_h}}$ if $\tilde{n}_{s_h, a_h} > E_\rho$, otherwise CH .
 - 6: $\hat{Q}_h^p(\cdot, \cdot) \leftarrow \tilde{Q}_h(\cdot, \cdot) - \Gamma_h(\cdot, \cdot)$.
 - 7: $\bar{Q}_h(\cdot, \cdot) \leftarrow \min\{\hat{Q}_h^p(\cdot, \cdot), H - h + 1\}^+$.
 - 8: $\forall s_h$, let $\hat{\pi}_h(\cdot | s_h) \leftarrow \operatorname{argmax}_{\pi_h} \langle \bar{Q}_h(s_h, \cdot), \pi_h(\cdot | s_h) \rangle$ and $\tilde{V}_h(s_h) \leftarrow \langle \bar{Q}_h(s_h, \cdot), \hat{\pi}_h(\cdot | s_h) \rangle$.
 - 9: **end for**
 - 10: **Output:** $\{\hat{\pi}_h\}$.
-

⁵This conclusion is summarized in Lemma D.3.

184 **Algorithmic design.** Our algorithmic design originates from the idea of pessimism, which holds
 185 conservative view towards the locations with high uncertainty and prefers the locations we have
 186 more confidence about. Based on the Bernstein type pessimism in APVI [Yin and Wang, 2021b], we
 187 design a similar pessimistic algorithm with private counts to ensure differential privacy. If we replace
 188 \tilde{n} and \tilde{P} with n and \hat{P}^6 , then our DP-APVI (Algorithm 1) will degenerate to APVI. Compared to
 189 the pessimism defined in APVI, our pessimistic penalty has an additional term $\tilde{O}\left(\frac{SHE_\rho}{\tilde{n}_{s_h, a_h}}\right)$, which
 190 accounts for the additional pessimism due to our application of private statistics.

191 We state our main theorem about DP-APVI below, the proof sketch is deferred to Appendix D.1 and
 192 detailed proof is deferred to Appendix D due to space limit.

193 **Theorem 3.4.** *DP-APVI (Algorithm 1) satisfies ρ -zCDP. Furthermore, under Assumption 2.1, denote*
 194 $\bar{d}_m := \min_{h \in [H]} \{d_h^\mu(s_h, a_h) : d_h^\mu(s_h, a_h) > 0\}$. *For any $0 < \delta < 1$, there exists constant $c_1 > 0$,*
 195 *such that when $n > c_1 \cdot \max\{H^2, E_\rho\} / \bar{d}_m \cdot \iota$ ($\iota = \log(HSA/\delta)$), with probability $1 - \delta$, the output*
 196 *policy $\hat{\pi}$ of DP-APVI satisfies*

$$0 \leq v^* - v^{\hat{\pi}} \leq 4\sqrt{2} \sum_{h=1}^H \sum_{(s_h, a_h) \in \mathcal{C}_h} d_h^{\pi^*}(s_h, a_h) \sqrt{\frac{\text{Var}_{P_h(\cdot|s_h, a_h)}(V_{h+1}^*(\cdot)) \cdot \iota}{nd_h^\mu(s_h, a_h)}} + \tilde{O}\left(\frac{H^3 + SH^2E_\rho}{n \cdot \bar{d}_m}\right), \quad (6)$$

197 where \tilde{O} hides constants and Polylog terms, $E_\rho = 4\sqrt{\frac{H \log \frac{4HS^2A}{\delta}}{\rho}}$.

198 **Comparison to non-private counterpart APVI [Yin and Wang, 2021b].** According to Theorem
 199 4.1 in [Yin and Wang, 2021b], the sub-optimality bound of APVI is for large enough n , with high
 200 probability, the output $\hat{\pi}$ satisfies:

$$0 \leq v^* - v^{\hat{\pi}} \leq \tilde{O}\left(\sum_{h=1}^H \sum_{(s_h, a_h) \in \mathcal{C}_h} d_h^{\pi^*}(s_h, a_h) \sqrt{\frac{\text{Var}_{P_h(\cdot|s_h, a_h)}(V_{h+1}^*(\cdot))}{nd_h^\mu(s_h, a_h)}}\right) + \tilde{O}\left(\frac{H^3}{n \cdot \bar{d}_m}\right). \quad (7)$$

201 Compared to our Theorem 3.4, the additional sub-optimality bound due to differential privacy is
 202 $\tilde{O}\left(\frac{SH^2E_\rho}{n \cdot \bar{d}_m}\right) = \tilde{O}\left(\frac{SH^{\frac{5}{2}}}{n \cdot \bar{d}_m \sqrt{\rho}}\right) = \tilde{O}\left(\frac{SH^{\frac{5}{2}}}{n \cdot \bar{d}_m \epsilon}\right)$.⁷ In the most popular regime where the privacy budget
 203 ρ or ϵ is a constant, the additional term due to differential privacy appears as a lower order term,
 204 hence becomes negligible as the sample complexity n becomes large.

205 **Comparison to Hoeffding type pessimism.** We can simply revise our algorithm by using Hoeffding
 206 type pessimism, which replaces the pessimism in line 5 with $C_1 H \cdot \sqrt{\frac{\iota}{\tilde{n}_{s_h, a_h} - E_\rho}} + \frac{C_2 SHE_\rho \cdot \iota}{\tilde{n}_{s_h, a_h}}$. Then
 207 with a similar proof schedule, we can arrive at a sub-optimality bound that with high probability,

$$0 \leq v^* - v^{\hat{\pi}} \leq \tilde{O}\left(H \cdot \sum_{h=1}^H \sum_{(s_h, a_h) \in \mathcal{C}_h} d_h^{\pi^*}(s_h, a_h) \sqrt{\frac{1}{nd_h^\mu(s_h, a_h)}}\right) + \tilde{O}\left(\frac{SH^2E_\rho}{n \cdot \bar{d}_m}\right). \quad (8)$$

208 Compared to our Theorem 3.4, our bound is tighter because we express the dominate term by the
 209 system quantities instead of explicit dependence on H (and $\text{Var}_{P_h(\cdot|s_h, a_h)}(V_{h+1}^*(\cdot)) \leq H^2$). In
 210 addition, we highlight that according to Theorem G.1 in [Yin and Wang, 2021b], our main term
 211 nearly matches the non-private minimax lower bound. For more detailed discussions about our main
 212 term and how it subsumes other optimal learning bounds, we refer readers to [Yin and Wang, 2021b].

213 **Apply Laplace Mechanism to achieve pure DP.** To achieve Pure DP instead of ρ -zCDP, we can
 214 simply replace Gaussian Mechanism with Laplace Mechanism (defined as Definition F.19). Given
 215 privacy budget for Pure DP ϵ , since the ℓ_1 sensitivity of $\{n_{s_h, a_h}\} \cup \{n_{s_h, a_h, s_{h+1}}\}$ is $\Delta_1 = 4H$, we
 216 can add *independent* Laplace noises $\text{Lap}(\frac{4H}{\epsilon})$ to each count to achieve ϵ -DP due to Lemma F.20.
 217 Then by using $E_\epsilon = \tilde{O}\left(\frac{H}{\epsilon}\right)$ instead of E_ρ and keeping everything else ((3), (5) and Algorithm 1) the
 218 same, we can reach a similar result to Theorem 3.4 with the same proof schedule. The only difference
 219 is that here the additional learning bound is $\tilde{O}\left(\frac{SH^3}{n \cdot \bar{d}_m \epsilon}\right)$, which still appears as a lower order term.

⁶The non-private empirical estimate, defined as (15) in Appendix D.

⁷Here we apply the second part of Lemma 2.6 to achieve (ϵ, δ) -DP, the notation \tilde{O} also absorbs $\log \frac{1}{\delta}$ (only here δ denotes the privacy budget instead of failure probability).

220 4 Results under linear MDP: DP-VAPVI(Algorithm 2)

221 In large MDPs, to address the computational issues, the technique of function approximation is
 222 widely applied, and linear MDP is a concrete model to study linear function approximations. Our
 223 second result applies to the linear MDP setting. Generally speaking, function approximation reduces
 224 the dimensionality of private releases comparing to the tabular MDPs. We begin with private counts.

225 **Private Model-based Components.** Given the two datasets \mathcal{D} and \mathcal{D}' (both from μ) as in Algorithm
 226 2, we can apply variance-aware pessimistic value iteration to learn a near optimal policy as in
 227 VAPVI [Yin et al., 2022]. To ensure differential privacy, we add *independent* Gaussian noises to the
 228 $5H$ statistics as in DP-VAPVI (Algorithm 2) below. Since there are $5H$ statistics, by the adaptive
 229 composition of zCDP (Lemma F.17), it suffices to keep each count ρ_0 -zCDP, where $\rho_0 = \frac{\rho}{5H}$. In
 230 DP-VAPVI, we use $\phi_1, \phi_2, \phi_3, K_1, K_2$ ⁸ to denote the noises we add. For all ϕ_i , we directly apply
 231 Gaussian Mechanism. For K_i , in addition to the noise matrix $\frac{1}{\sqrt{2}}(Z + Z^\top)$, we also add $\frac{E}{2}I_d$ to
 232 ensure that all K_i are positive definite with high probability (The detailed definition of E, L can be
 233 found in Appendix A).

Algorithm 2 Differentially Private Variance-Aware Pessimistic Value Iteration (DP-VAPVI)

- 1: **Input:** Dataset $\mathcal{D} = \{(s_h^\tau, a_h^\tau, r_h^\tau, s_{h+1}^\tau)\}_{\tau, h=1}^{K, H}$ $\mathcal{D}' = \{(\bar{s}_h^\tau, \bar{a}_h^\tau, \bar{r}_h^\tau, \bar{s}_{h+1}^\tau)\}_{\tau, h=1}^{K, H}$. Budget for zCDP ρ . Failure probability δ . Universal constant C .
 - 2: **Initialization:** Set $\rho_0 \leftarrow \frac{\rho}{5H}$, $\tilde{V}_{H+1}(\cdot) \leftarrow 0$. Sample $\phi_1 \sim \mathcal{N}\left(0, \frac{2H^4}{\rho_0} I_d\right)$, $\phi_2, \phi_3 \sim \mathcal{N}\left(0, \frac{2H^2}{\rho_0} I_d\right)$, $K_1, K_2 \leftarrow \frac{E}{2}I_d + \frac{1}{\sqrt{2}}(Z + Z^\top)$, where $Z_{i,j} \sim \mathcal{N}\left(0, \frac{1}{4\rho_0}\right)$ (i.i.d.), $E = \tilde{O}\left(\sqrt{\frac{Hd}{\rho}}\right)$. Set $D \leftarrow \tilde{O}\left(\frac{H^2L}{\kappa} + \frac{H^4E\sqrt{d}}{\kappa^{3/2}} + H^3\sqrt{d}\right)$.
 - 3: **for** $h = H, H-1, \dots, 1$ **do**
 - 4: Set $\tilde{\Sigma}_h \leftarrow \sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \phi(\bar{s}_h^\tau, \bar{a}_h^\tau)^\top + \lambda I + K_1$
 - 5: Set $\tilde{\beta}_h \leftarrow \tilde{\Sigma}_h^{-1} [\sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \cdot \tilde{V}_{h+1}(\bar{s}_{h+1}^\tau)^2 + \phi_1]$
 - 6: Set $\tilde{\theta}_h \leftarrow \tilde{\Sigma}_h^{-1} [\sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \cdot \tilde{V}_{h+1}(\bar{s}_{h+1}^\tau) + \phi_2]$
 - 7: Set $[\widetilde{\text{Var}}_h \tilde{V}_{h+1}](\cdot, \cdot) \leftarrow \langle \phi(\cdot, \cdot), \tilde{\beta}_h \rangle_{[0, (H-h+1)^2]} - [\langle \phi(\cdot, \cdot), \tilde{\theta}_h \rangle_{[0, H-h+1]}]^2$
 - 8: Set $\tilde{\sigma}_h(\cdot, \cdot)^2 \leftarrow \max\{1, \widetilde{\text{Var}}_h \tilde{V}_{h+1}(\cdot, \cdot)\}$
 - 9: Set $\tilde{\Lambda}_h \leftarrow \sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top / \tilde{\sigma}_h^2(s_h^\tau, a_h^\tau) + \lambda I + K_2$
 - 10: Set $\tilde{w}_h \leftarrow \tilde{\Lambda}_h^{-1} \left(\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \cdot (r_h^\tau + \tilde{V}_{h+1}(s_{h+1}^\tau)) / \tilde{\sigma}_h^2(s_h^\tau, a_h^\tau) + \phi_3 \right)$
 - 11: Set $\Gamma_h(\cdot, \cdot) \leftarrow C\sqrt{d} \cdot \left(\phi(\cdot, \cdot)^\top \tilde{\Lambda}_h^{-1} \phi(\cdot, \cdot) \right)^{1/2} + \frac{D}{K}$
 - 12: Set $\tilde{Q}_h(\cdot, \cdot) \leftarrow \phi(\cdot, \cdot)^\top \tilde{w}_h - \Gamma_h(\cdot, \cdot)$
 - 13: Set $\hat{Q}_h(\cdot, \cdot) \leftarrow \min\{\tilde{Q}_h(\cdot, \cdot), H-h+1\}^+$
 - 14: Set $\hat{\pi}_h(\cdot | \cdot) \leftarrow \operatorname{argmax}_{\pi_h} \langle \hat{Q}_h(\cdot, \cdot), \pi_h(\cdot | \cdot) \rangle_{\mathcal{A}}$, $\tilde{V}_h(\cdot) \leftarrow \max_{\pi_h} \langle \hat{Q}_h(\cdot, \cdot), \pi_h(\cdot | \cdot) \rangle_{\mathcal{A}}$
 - 15: **end for**
 - 16: **Output:** $\{\hat{\pi}_h\}_{h=1}^H$.
-

234 Below we will show the algorithmic design of DP-VAPVI (Algorithm 2). For the offline dataset,
 235 we divide it into two independent parts with equal length: $\mathcal{D} = \{(s_h^\tau, a_h^\tau, r_h^\tau, s_{h+1}^\tau)\}_{\tau \in [K]}^{h \in [H]}$ and
 236 $\mathcal{D}' = \{(\bar{s}_h^\tau, \bar{a}_h^\tau, \bar{r}_h^\tau, \bar{s}_{h+1}^\tau)\}_{\tau \in [K]}^{h \in [H]}$. One for estimating variance and the other for calculating Q -values.

237 **Estimating conditional variance.** The first part (line 4 to line 8) aims to estimate the condi-
 238 tional variance of \tilde{V}_{h+1} via the definition of variance: $[\text{Var}_h \tilde{V}_{h+1}](s, a) = [P_h(\tilde{V}_{h+1})^2](s, a) -$
 239 $([P_h \tilde{V}_{h+1}](s, a))^2$. For the first term, by the definition of linear MDP, it holds that
 240 $[P_h \tilde{V}_{h+1}^2](s, a) = \phi(s, a)^\top \int_{\mathcal{S}} \tilde{V}_{h+1}^2(s') d\nu_h(s') = \langle \phi, \int_{\mathcal{S}} \tilde{V}_{h+1}^2(s') d\nu_h(s') \rangle$. We can estimate
 241 $\beta_h = \int_{\mathcal{S}} \tilde{V}_{h+1}^2(s') d\nu_h(s')$ by applying ridge regression. Below is the output of ridge regression
 242 with raw statistics without noise:

⁸We need to add noise to each of the $5H$ counts, therefore for ϕ_1 , we actually sample H i.i.d samples $\phi_{1,h}$, $h = 1, \dots, H$ from the distribution of ϕ_1 . Then we add $\phi_{1,h}$ to $\sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \cdot \tilde{V}_{h+1}(\bar{s}_{h+1}^\tau)^2, \forall h \in [H]$. For simplicity, we use ϕ_1 to represent all the $\phi_{1,h}$. The procedure applied to the other $4H$ statistics are similar.

$$\operatorname{argmin}_{\beta \in \mathbb{R}^d} \sum_{k=1}^K \left[\left\langle \phi(\bar{s}_h^k, \bar{a}_h^k), \beta \right\rangle - \tilde{V}_{h+1}(\bar{s}_{h+1}^k) \right]^2 + \lambda \|\beta\|_2^2 = \bar{\Sigma}_h^{-1} \sum_{k=1}^K \phi(\bar{s}_h^k, \bar{a}_h^k) \tilde{V}_{h+1}^2(\bar{s}_{h+1}^k),$$

243 where definition of $\bar{\Sigma}_h$ can be found in Appendix A. Instead of using the raw statistics, we replace
 244 them with private ones with Gaussian noises as in line 5. The second term is estimated similarly in
 245 line 6. The final estimator is defined as in line 8: $\tilde{\sigma}_h(\cdot, \cdot)^2 = \max\{1, \widetilde{\text{Var}}_h \tilde{V}_{h+1}(\cdot, \cdot)\}$ ⁹.

246 **Variance-weighted LSVI.** Instead of directly applying LSVI [Jin et al., 2021], we can solve the
 247 variance-weighted LSVI (line 10). The result of variance-weighted LSVI with non-private statistics
 248 is shown below:

$$\operatorname{argmin}_{w \in \mathbb{R}^d} \lambda \|w\|_2^2 + \sum_{k=1}^K \frac{\left[\langle \phi(s_h^k, a_h^k), w \rangle - r_h^k - \tilde{V}_{h+1}(s_{h+1}^k) \right]^2}{\tilde{\sigma}_h^2(s_h^k, a_h^k)} = \hat{\Lambda}_h^{-1} \sum_{k=1}^K \frac{\phi(s_h^k, a_h^k) \cdot \left[r_h^k + \tilde{V}_{h+1}(s_{h+1}^k) \right]}{\tilde{\sigma}_h^2(s_h^k, a_h^k)},$$

249 where definition of $\hat{\Lambda}_h$ can be found in Appendix A. For the sake of differential privacy, we use
 250 private statistics instead and derive the \tilde{w}_h as in line 10.

251 **Our private pessimism.** Notice that if we remove all the Gaussian noises we add, our DP-VAPVI
 252 (Algorithm 2) will degenerate to VAPVI [Yin et al., 2022]. We design a similar pessimistic penalty
 253 using private statistics (line 11), with additional $\frac{D}{K}$ accounting for the extra pessimism due to DP.

254 **Main theorem.** We state our main theorem about DP-VAPVI below, the proof sketch is deferred to
 255 Appendix E.1 and detailed proof is deferred to Appendix E due to space limit. Note that quantities
 256 \mathcal{M}_i, L, E can be found in Appendix A and briefly, $L = \tilde{O}(\sqrt{H^3 d / \rho})$, $E = \tilde{O}(\sqrt{H d / \rho})$. For the
 257 sample complexity lower bound, within the practical regime where the privacy budget is not very
 258 small, $\max\{\mathcal{M}_i\}$ is dominated by $\max\{\tilde{O}(H^{12} d^3 / \kappa^5), \tilde{O}(H^{14} d / \kappa^5)\}$, which also appears in the
 259 sample complexity lower bound of VAPVI [Yin et al., 2022]. The $\sigma_V^2(s, a)$ in Theorem 4.1 is defined
 260 as $\max\{1, \text{Var}_{P_h}(V)(s, a)\}$ for any V .

261 **Theorem 4.1.** DP-VAPVI (Algorithm 2) satisfies ρ -zCDP. Furthermore, let K be the number of
 262 episodes. Under the condition that $K > \max\{\mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_3, \mathcal{M}_4\}$ and $\sqrt{d} > \xi$, where $\xi :=$
 263 $\sup_{V \in [0, H], s' \sim P_h(s, a), h \in [H]} \left| \frac{r_h + V(s') - (\mathcal{T}_h V)(s, a)}{\sigma_V(s, a)} \right|$, for any $0 < \lambda < \kappa$, with probability $1 - \delta$, for
 264 all policy π simultaneously, the output $\hat{\pi}$ of DP-VAPVI satisfies

$$v^\pi - v^{\hat{\pi}} \leq \tilde{O} \left(\sqrt{d} \cdot \sum_{h=1}^H \mathbb{E}_\pi \left[\sqrt{\phi(\cdot, \cdot)^\top \Lambda_h^{-1} \phi(\cdot, \cdot)} \right] \right) + \frac{DH}{K}, \quad (9)$$

265 where $\Lambda_h = \sum_{k=1}^K \frac{\phi(s_h^k, a_h^k) \cdot \phi(s_h^k, a_h^k)^\top}{\sigma_{\tilde{V}_{h+1}(s_h^k, a_h^k)}^2} + \lambda I_d$, $D = \tilde{O} \left(\frac{H^2 L}{\kappa} + \frac{H^4 E \sqrt{d}}{\kappa^{3/2}} + H^3 \sqrt{d} \right)$ and \tilde{O} hides
 266 constants and Polylog terms.

267 In particular, define $\Lambda_h^* = \sum_{k=1}^K \frac{\phi(s_h^k, a_h^k) \cdot \phi(s_h^k, a_h^k)^\top}{\sigma_{V_{h+1}^*(s_h^k, a_h^k)}^2} + \lambda I_d$, we have with probability $1 - \delta$,

$$v^* - v^{\hat{\pi}} \leq \tilde{O} \left(\sqrt{d} \cdot \sum_{h=1}^H \mathbb{E}_{\pi^*} \left[\sqrt{\phi(\cdot, \cdot)^\top \Lambda_h^{*-1} \phi(\cdot, \cdot)} \right] \right) + \frac{DH}{K}. \quad (10)$$

268 **Comparison to non-private counterpart VAPVI [Yin et al., 2022].** Plugging in the definition
 269 of L, E (Appendix A), under the meaningful case that the privacy budget is not very large, DH is
 270 dominated by $\tilde{O} \left(\frac{H^{\frac{13}{2}} d / \kappa^{\frac{3}{2}}}{\sqrt{\rho}} \right)$. According to Theorem 3.2 in [Yin et al., 2022], the sub-optimality
 271 bound of VAPVI is for sufficiently large K , with high probability, the output $\hat{\pi}$ satisfies:

$$v^* - v^{\hat{\pi}} \leq \tilde{O} \left(\sqrt{d} \cdot \sum_{h=1}^H \mathbb{E}_{\pi^*} \left[\sqrt{\phi(\cdot, \cdot)^\top \Lambda_h^{*-1} \phi(\cdot, \cdot)} \right] \right) + \frac{2H^4 \sqrt{d}}{K}. \quad (11)$$

⁹The $\max\{1, \cdot\}$ operator here is for technical reason only: we want a lower bound for each variance estimate.

272 Compared to our Theorem 4.1, the additional sub-optimality bound due to differential privacy is
 273 $\tilde{O}\left(\frac{H^{\frac{11}{2}} d/\kappa^{\frac{3}{2}}}{\sqrt{\rho} \cdot K}\right) = \tilde{O}\left(\frac{H^{\frac{11}{2}} d/\kappa^{\frac{3}{2}}}{\epsilon \cdot K}\right)$.¹⁰ In the most popular regime where the privacy budget ρ or ϵ is
 274 a constant, the additional term due to differential privacy also appears as a lower order term.

275 **Instance-dependent sub-optimality bound.** Similar to DP-APVI (Algorithm 1), our DP-VAPVI
 276 (Algorithm 2) also enjoys instance-dependent sub-optimality bound. First, the main term in (10)
 277 improves PEVI [Jin et al., 2021] over $O(\sqrt{d})$ on feature dependence. Also, our main term admits no
 278 explicit dependence on H , thus improves the sub-optimality bound of PEVI on horizon dependence.
 279 For more detailed discussions about our main term, we refer readers to [Yin et al., 2022].

280 5 Simulations

281 In this section, we carry out simulations to evaluate the performance of our DP-VAPVI (Algorithm 2),
 282 and compare it with its non-private counterpart VAPVI [Yin et al., 2022] and another pessimism-based
 283 algorithm PEVI [Jin et al., 2021] which does not have privacy guarantee.

284 **Experimental setting.** We evaluate DP-VAPVI (Algorithm 2) on a synthetic linear MDP example that
 285 originates from the linear MDP in [Min et al., 2021, Yin et al., 2022] but with some modifications.¹¹
 286 For details of the linear MDP setting, please refer to Appendix G. The two MDP instances we use
 287 both have horizon $H = 20$. We compare different algorithms in figure 1(a), while in figure 1(b), we
 288 compare our DP-VAPVI with different privacy budgets. When doing empirical evaluation, we do not
 289 split the data for DP-VAPVI or VAPVI and for DP-VAPVI, we run the simulation for 5 times and
 290 take the average performance.

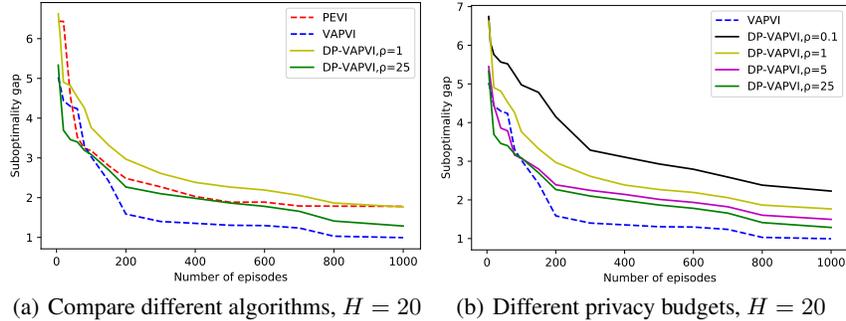


Figure 1: Comparison between performance of PEVI, VAPVI and DP-VAPVI (with different privacy budgets) under the linear MDP example described above. In each figure, y-axis represents sub-optimality gap $v^* - v^{\hat{\pi}}$ while x-axis denotes the number of episodes K . The horizons are fixed to be $H = 20$. The number of episodes takes value from 5 to 1000.

291 **Results and discussions.** From Figure 1, we can observe that DP-VAPVI (Algorithm 2) performs
 292 slightly worse than its non-private version VAPVI [Yin et al., 2022]. This is due to the fact that
 293 we add Gaussian noise to each count. However, as the size of dataset goes larger, the performance
 294 of DP-VAPVI will converge to that of VAPVI, which supports our theoretical conclusion that the
 295 cost of privacy only appears as lower order terms. For DP-VAPVI with larger privacy budget, the
 296 scale of noise will be smaller, thus the performance will be closer to VAPVI, as shown in figure
 297 1(b). Furthermore, in most cases, DP-VAPVI still outperforms PEVI, which does not have privacy
 298 guarantee. This arises from our privatization of variance-aware LSVI instead of LSVI.

299 6 Conclusion

300 In this work, we take the first steps towards the well-motivated task of designing private offline RL
 301 algorithms. We propose algorithms for both tabular MDPs and linear MDPs, and show that they
 302 enjoy instance-dependent sub-optimality bounds while guaranteeing differential privacy (either zCDP
 303 or pure DP). Our results highlight that the cost of privacy only appears as lower order terms, thus
 304 become negligible as the number of samples goes large.

¹⁰Here we apply the second part of Lemma 2.6 to achieve (ϵ, δ) -DP, the notation \tilde{O} also absorbs $\log \frac{1}{\delta}$ (only here δ denotes the privacy budget instead of failure probability).

¹¹We keep the state space $\mathcal{S} = \{1, 2\}$, action space $\mathcal{A} = \{1, \dots, 100\}$ and feature map of state-action pairs while we choose stochastic transition (instead of the original deterministic transition) and more complex reward.

305 References

- 306 Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic
307 bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320, 2011.
- 308 Naman Agarwal and Karan Singh. The price of differential privacy for online learning. In *International
309 Conference on Machine Learning*, pages 32–40. PMLR, 2017.
- 310 Alex Ayoub, Zeyu Jia, Csaba Szepesvari, Mengdi Wang, and Lin Yang. Model-based reinforcement
311 learning with value-targeted regression. In *International Conference on Machine Learning*, pages
312 463–474. PMLR, 2020.
- 313 Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforce-
314 ment learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume
315 70*, pages 263–272. JMLR. org, 2017.
- 316 Borja Balle, Maziar Gomrokchi, and Doina Precup. Differentially private policy evaluation. In
317 *International Conference on Machine Learning*, pages 2130–2138. PMLR, 2016.
- 318 Debabrota Basu, Christos Dimitrakakis, and Aristide Tossou. Differential privacy for multi-armed
319 bandits: What is it and what is its cost? *arXiv preprint arXiv:1905.12298*, 2019.
- 320 Gavin Brown, Mark Bun, Vitaly Feldman, Adam Smith, and Kunal Talwar. When is memorization of
321 irrelevant training data necessary for high-accuracy learning? In *ACM SIGACT Symposium on
322 Theory of Computing*, pages 123–132, 2021.
- 323 Mark Bun and Thomas Steinke. Concentrated differential privacy: Simplifications, extensions, and
324 lower bounds. In *Theory of Cryptography Conference*, pages 635–658. Springer, 2016.
- 325 Qi Cai, Zhuoran Yang, Chi Jin, and Zhaoran Wang. Provably efficient exploration in policy optimiza-
326 tion. In *International Conference on Machine Learning*, pages 1283–1294. PMLR, 2020.
- 327 Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. The secret sharer: Evalu-
328 ating and testing unintended memorization in neural networks. In *USENIX Security Symposium
329 (USENIX Security 19)*, pages 267–284, 2019.
- 330 Xiaoyu Chen, Kai Zheng, Zixin Zhou, Yunchang Yang, Wei Chen, and Liwei Wang. (locally) differ-
331 entially private combinatorial semi-bandits. In *International Conference on Machine Learning*,
332 pages 1757–1767. PMLR, 2020.
- 333 Herman Chernoff et al. A measure of asymptotic efficiency for tests of a hypothesis based on the
334 sum of observations. *The Annals of Mathematical Statistics*, 23(4):493–507, 1952.
- 335 Sayak Ray Chowdhury and Xingyu Zhou. Differentially private regret minimization in episodic
336 markov decision processes. *arXiv preprint arXiv:2112.10599*, 2021.
- 337 Sayak Ray Chowdhury, Xingyu Zhou, and Ness Shroff. Adaptive control of differentially private
338 linear quadratic systems. In *2021 IEEE International Symposium on Information Theory (ISIT)*,
339 pages 485–490. IEEE, 2021.
- 340 Chris Cundy and Stefano Ermon. Privacy-constrained policies via mutual information regularized
341 policy gradients. *arXiv preprint arXiv:2012.15019*, 2020.
- 342 Christoph Dann, Tor Lattimore, and Emma Brunskill. Unifying pac and regret: Uniform pac bounds
343 for episodic reinforcement learning. In *Advances in Neural Information Processing Systems*, pages
344 5713–5723, 2017.
- 345 Cynthia Dwork and Guy N Rothblum. Concentrated differential privacy. *arXiv preprint
346 arXiv:1603.01887*, 2016.
- 347 Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in
348 private data analysis. In *Theory of cryptography conference*, pages 265–284. Springer, 2006.
- 349 Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Found. Trends
350 Theor. Comput. Sci.*, 9(3-4):211–407, 2014.

- 351 Frederick Arthur Ficken. *The simplex method of linear programming*. Courier Dover Publications,
352 2015.
- 353 Pratik Gajane, Tanguy Urvoy, and Emilie Kaufmann. Corrupt bandits for preserving local privacy. In
354 *Algorithmic Learning Theory*, pages 387–412. PMLR, 2018.
- 355 Evrard Garcelon, Vianney Perchet, Ciara Pike-Burke, and Matteo Pirota. Local differential privacy
356 for regret minimization in reinforcement learning. *Advances in Neural Information Processing*
357 *Systems*, 34, 2021.
- 358 Abhradeep Guha Thakurta and Adam Smith. (nearly) optimal algorithms for private online learning
359 in full-information and bandit settings. *Advances in Neural Information Processing Systems*, 26,
360 2013.
- 361 Bingshan Hu, Zhiming Huang, and Nishant A Mehta. Optimal algorithms for private online learning
362 in a stochastic environment. *arXiv preprint arXiv:2102.07929*, 2021.
- 363 Chi Jin, Akshay Krishnamurthy, Max Simchowitz, and Tiancheng Yu. Reward-free exploration for
364 reinforcement learning. In *International Conference on Machine Learning*, pages 4870–4879.
365 PMLR, 2020a.
- 366 Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement
367 learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143.
368 PMLR, 2020b.
- 369 Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl? In
370 *International Conference on Machine Learning*, pages 5084–5096. PMLR, 2021.
- 371 Hyeoneui Kim, Elizabeth Bell, Jihoon Kim, Amy Sitapati, Joe Ramsdell, Claudiu Farcas, Dexter
372 Friedman, Stephanie Feudjio Feupe, and Lucila Ohno-Machado. iconcur: informed consent
373 for clinical data and bio-sample use for research. *Journal of the American Medical Informatics*
374 *Association*, 24(2):380–387, 2017.
- 375 Jonathan Lebensold, William Hamilton, Borja Balle, and Doina Precup. Actor critic with differentially
376 private critic. *arXiv preprint arXiv:1910.05876*, 2019.
- 377 Chonghua Liao, Jiafan He, and Quanquan Gu. Locally differentially private reinforcement learning
378 for linear mixture markov decision processes. *arXiv preprint arXiv:2110.10133*, 2021.
- 379 Yao Liu, Adith Swaminathan, Alekh Agarwal, and Emma Brunskill. Off-policy policy gradient with
380 state distribution correction. In *Uncertainty in Artificial Intelligence*, 2019.
- 381 Paul Luyo, Evrard Garcelon, Alessandro Lazaric, and Matteo Pirota. Differentially private explo-
382 ration in reinforcement learning with linear representation. *arXiv preprint arXiv:2112.01585*,
383 2021.
- 384 Andreas Maurer and Massimiliano Pontil. Empirical bernstein bounds and sample variance penaliza-
385 tion. *Conference on Learning Theory*, 2009.
- 386 Yifei Min, Tianhao Wang, Dongruo Zhou, and Quanquan Gu. Variance-aware off-policy evaluation
387 with linear function approximation. *Advances in neural information processing systems*, 34, 2021.
- 388 George Nemhauser and Laurence Wolsey. Polynomial-time algorithms for linear programming.
389 *Integer and Combinatorial Optimization*, pages 146–181, 1988.
- 390 Dung Daniel Ngo, Giuseppe Vietri, and Zhiwei Steven Wu. Improved regret for differentially private
391 exploration in linear mdp. *arXiv preprint arXiv:2202.01292*, 2022.
- 392 Hajime Ono and Tsubasa Takahashi. Locally private distributed reinforcement learning. *arXiv*
393 *preprint arXiv:2001.11718*, 2020.
- 394 Dan Qiao and Yu-Xiang Wang. Near-optimal differentially private reinforcement learning. *arXiv*
395 *preprint arXiv:2212.04680*, 2022.

- 396 Dan Qiao, Ming Yin, Ming Min, and Yu-Xiang Wang. Sample-efficient reinforcement learning
397 with $\log\log(T)$ switching cost. In *Proceedings of the 39th International Conference on Machine*
398 *Learning*, pages 18031–18061. PMLR, 2022.
- 399 Aniruddh Raghu, Matthieu Komorowski, Leo Anthony Celi, Peter Szolovits, and Marzyeh Ghas-
400 semi. Continuous state-space models for optimal sepsis treatment: a deep reinforcement learning
401 approach. In *Machine Learning for Healthcare Conference*, pages 147–163, 2017.
- 402 Rachel Redberg and Yu-Xiang Wang. Privately publishable per-instance privacy. *Advances in Neural*
403 *Information Processing Systems*, 34, 2021.
- 404 Ahmad EL Sallab, Mohammed Abdou, Etienne Perot, and Senthil Yogamani. Deep reinforcement
405 learning framework for autonomous driving. *Electronic Imaging*, 2017(19):70–76, 2017.
- 406 Roshan Shariff and Or Sheffet. Differentially private contextual linear bandits. *Advances in Neural*
407 *Information Processing Systems*, 31, 2018.
- 408 Karthik Sridharan. A gentle introduction to concentration inequalities. *Dept. Comput. Sci., Cornell*
409 *Univ., Tech. Rep*, 2002.
- 410 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 411 Aristide Charles Yedia Tossou and Christos Dimitrakakis. Achieving privacy in the adversarial
412 multi-armed bandit. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- 413 Tim Van Erven and Peter Harremo. Rényi divergence and kullback-leibler divergence. *IEEE*
414 *Transactions on Information Theory*, 60(7):3797–3820, 2014.
- 415 Giuseppe Vietri, Borja Balle, Akshay Krishnamurthy, and Steven Wu. Private reinforcement learning
416 with pac and regret guarantees. In *International Conference on Machine Learning*, pages 9754–
417 9764. PMLR, 2020.
- 418 Baoxiang Wang and Nidhi Hegde. Privacy-preserving q-learning with functional noise in continuous
419 spaces. *Advances in Neural Information Processing Systems*, 32, 2019.
- 420 Ruosong Wang, Dean P Foster, and Sham M Kakade. What are the statistical limits of offline rl with
421 linear function approximation? *International Conference on Learning Representations*, 2021.
- 422 Tengyang Xie, Philip S Thomas, and Gerome Miklau. Privacy preserving off-policy evaluation. *arXiv*
423 *preprint arXiv:1902.00174*, 2019.
- 424 Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-consistent
425 pessimism for offline reinforcement learning. *Advances in neural information processing systems*,
426 2021a.
- 427 Tengyang Xie, Nan Jiang, Huan Wang, Caiming Xiong, and Yu Bai. Policy finetuning: Bridg-
428 ing sample-efficient offline and online reinforcement learning. *Advances in neural information*
429 *processing systems*, 2021b.
- 430 Ming Yin and Yu-Xiang Wang. Asymptotically efficient off-policy evaluation for tabular rein-
431 forcement learning. In *International Conference on Artificial Intelligence and Statistics*, pages
432 3948–3958. PMLR, 2020.
- 433 Ming Yin and Yu-Xiang Wang. Optimal uniform ope and model-based offline reinforcement learning
434 in time-homogeneous, reward-free and task-agnostic settings. *Advances in neural information*
435 *processing systems*, 2021a.
- 436 Ming Yin and Yu-Xiang Wang. Towards instance-optimal offline reinforcement learning with
437 pessimism. *Advances in neural information processing systems*, 34, 2021b.
- 438 Ming Yin, Yu Bai, and Yu-Xiang Wang. Near-optimal provable uniform convergence in offline policy
439 evaluation for reinforcement learning. In *International Conference on Artificial Intelligence and*
440 *Statistics*, pages 1567–1575. PMLR, 2021.

- 441 Ming Yin, Yaqi Duan, Mengdi Wang, and Yu-Xiang Wang. Near-optimal offline reinforcement
442 learning with linear representation: Leveraging variance information with pessimism. *arXiv*
443 *preprint arXiv:2203.05804*, 2022.
- 444 Andrea Zanette. Exponential lower bounds for batch reinforcement learning: Batch rl can be
445 exponentially harder than online rl. In *International Conference on Machine Learning*, pages
446 12287–12297. PMLR, 2021.
- 447 Andrea Zanette, Martin J Wainwright, and Emma Brunskill. Provable benefits of actor-critic methods
448 for offline reinforcement learning. *Advances in neural information processing systems*, 2021.
- 449 Zihao Zhang, Stefan Zohren, and Stephen Roberts. Deep reinforcement learning for trading. *The*
450 *Journal of Financial Data Science*, 2(2):25–40, 2020.
- 451 Kai Zheng, Tianle Cai, Weiran Huang, Zhenguo Li, and Liwei Wang. Locally differentially private
452 (contextual) bandits learning. *Advances in Neural Information Processing Systems*, 33:12300–
453 12310, 2020.
- 454 Dongruo Zhou, Quanquan Gu, and Csaba Szepesvari. Nearly minimax optimal reinforcement
455 learning for linear mixture markov decision processes. In *Conference on Learning Theory*, pages
456 4532–4576. PMLR, 2021.
- 457 Xingyu Zhou. Differentially private reinforcement learning with linear function approximation. *arXiv*
458 *preprint arXiv:2201.07052*, 2022.

459 **A Notation List**

460 **A.1 Notations for tabular MDP**

E_ρ	$4\sqrt{\frac{H \log \frac{4HS^2A}{\delta}}{\rho}}$
n	The original counts of visitation
n'	The noisy counts, as defined in (2)
\tilde{n}	Final choice of private counts, as defined in (3)
461 \tilde{P}	Private estimate of transition kernel, as defined in (5)
\hat{P}	Non-private estimate of transition kernel, as defined in (15)
ι	$\log \frac{HSA}{\delta}$
ρ	Budget for zCDP
δ	Failure probability

462 **A.2 Notations for linear MDP**

L	$2H\sqrt{\frac{5Hd \log(\frac{10Hd}{\delta})}{\rho}}$
E	$\sqrt{\frac{10Hd}{\rho}} \left(2 + \left(\frac{\log(5c_1H/\delta)}{c_2d} \right)^{\frac{2}{3}} \right)$
D	$\tilde{O} \left(\frac{H^2L}{\kappa} + \frac{H^4E\sqrt{d}}{\kappa^{3/2}} + H^3\sqrt{d} \right)$
$\hat{\Lambda}_h$	$\sum_{k=1}^K \phi(s_h^k, a_h^k) \phi(s_h^k, a_h^k)^\top / \tilde{\sigma}_h^2(s_h^k, a_h^k) + \lambda I_d$
$\tilde{\Lambda}_h$	$\sum_{k=1}^K \phi(s_h^k, a_h^k) \phi(s_h^k, a_h^k)^\top / \tilde{\sigma}_h^2(s_h^k, a_h^k) + \lambda I_d + K_2$
$\tilde{\Lambda}_h^p$	$\mathbb{E}_{\mu, h} [\tilde{\sigma}_h^{-2}(s, a) \phi(s, a) \phi(s, a)^\top]$
Λ_h	$\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top / \sigma_{V_{h+1}}^2(s_h^\tau, a_h^\tau) + \lambda I$
Λ_h^p	$\mathbb{E}_{\mu, h} [\sigma_{V_{h+1}}^{-2}(s, a) \phi(s, a) \phi(s, a)^\top]$
Λ_h^*	$\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top / \sigma_{V_{h+1}^*}^2(s_h^\tau, a_h^\tau) + \lambda I$
$\bar{\Sigma}_h$	$\sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \phi(\bar{s}_h^\tau, \bar{a}_h^\tau)^\top + \lambda I_d$
$\tilde{\Sigma}_h$	$\sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \phi(\bar{s}_h^\tau, \bar{a}_h^\tau)^\top + \lambda I_d + K_1$
463 Σ_h^p	$\mathbb{E}_{\mu, h} [\phi(s, a) \phi(s, a)^\top]$
κ	$\min_h \lambda_{\min}(\Sigma_h^p)$
$\sigma_V^2(s, a)$	$\max\{1, \text{Var}_{P_h}(V)(s, a)\}$ for any V
$\sigma_h^{*2}(s, a)$	$\max\{1, \text{Var}_{P_h} V_{h+1}^*(s, a)\}$
$\tilde{\sigma}_h^2(s, a)$	$\max\{1, \widetilde{\text{Var}}_h \tilde{V}_{h+1}(s, a)\}$
\mathcal{M}_1	$\max\{2\lambda, 128 \log(2dH/\delta), \frac{128H^4 \log(2dH/\delta)}{\kappa^2}, \frac{\sqrt{2}L}{\sqrt{d\kappa}}\}$
\mathcal{M}_2	$\max\{\tilde{O}(H^{12}d^3/\kappa^5), \tilde{O}(H^{14}d/\kappa^5)\}$
\mathcal{M}_3	$\max\left\{ \frac{512H^4 \log(\frac{2dH}{\delta})}{\kappa^2}, \frac{4\lambda H^2}{\kappa} \right\}$
\mathcal{M}_4	$\max\left\{ \frac{H^2L^2}{d\kappa}, \frac{H^6E^2}{\kappa^2}, H^4\kappa \right\}$
ρ	Budget for zCDP
δ	Failure probability (<i>not</i> the δ of (ϵ, δ) -DP)
ξ	$\sup_{V \in [0, H], s' \sim P_h(s, a), h \in [H]} \left \frac{r_h + V(s') - (\mathcal{T}_h V)(s, a)}{\sigma_V(s, a)} \right $

464 B Extended related work

465 **Online reinforcement learning under JDP or LDP.** For online RL, some recent works analyze
466 this setting under *Joint Differential Privacy (JDP)*, which requires the RL agent to minimize regret
467 while handling user’s raw data privately. Under tabular MDP, Vietri et al. [2020] design PUCB by
468 revising UBEV [Dann et al., 2017]. Private-UCB-VI [Chowdhury and Zhou, 2021] results from
469 UCBVI (with bonus-1) [Azar et al., 2017]. However, both works privatize Hoeffding type bonus,
470 which lead to sub-optimal regret bound. Under linear MDP, Private LSVI-UCB [Ngo et al., 2022]
471 and Privacy-Preserving LSVI-UCB [Luyo et al., 2021] are private versions of LSVI-UCB [Jin et al.,
472 2020b], while LinOpt-VI-Reg [Zhou, 2022] and Privacy-Preserving UCRL-VTR [Luyo et al., 2021]
473 generalize UCRL-VTR [Ayoub et al., 2020]. However, these works are usually based on the LSVI
474 technique [Jin et al., 2020b] (unweighted ridge regression), which does not ensure optimal regret
475 bound.

476 In addition to JDP, another common privacy guarantee for online RL is *Local Differential Privacy*
477 (*LDP*), LDP is a stronger definition of DP since it requires that the user’s data is protected before the
478 RL agent has access to it. Under LDP, Garcelon et al. [2021] reach a regret lower bound and design
479 LDP-OB1 which has matching regret upper bound. The result is generalized by Liao et al. [2021] to
480 linear mixture setting. Later, Luyo et al. [2021] provide an unified framework for analyzing JDP and
481 LDP under linear setting.

482 **Some other differentially private learning algorithms.** There are some other works about dif-
483 ferentially private online learning [Guha Thakurta and Smith, 2013, Agarwal and Singh, 2017, Hu
484 et al., 2021] and various settings of bandit [Shariff and Sheffet, 2018, Gajane et al., 2018, Basu et al.,
485 2019, Zheng et al., 2020, Chen et al., 2020, Tossou and Dimitrakakis, 2017]. For the reinforcement
486 learning setting, Wang and Hegde [2019] propose privacy-preserving Q-learning to protect the reward
487 information. Ono and Takahashi [2020] study the problem of distributed reinforcement learning
488 under LDP. Lebensold et al. [2019] present an actor critic algorithm with differentially private critic.
489 Cundy and Ermon [2020] tackle DP-RL under the policy gradient framework. Chowdhury et al.
490 [2021] consider the adaptive control of differentially private linear quadratic (LQ) systems.

491 **Offline reinforcement learning under tabular MDP.** Under tabular MDP, there are several works
492 achieving optimal sub-optimality/sample complexity bounds under different coverage assumptions.
493 For the problem of off-policy evaluation (OPE), Yin and Wang [2020] uses Tabular-MIS estimator to
494 achieve asymptotic efficiency. In addition, the idea of uniform OPE is used to achieve the optimal
495 sample complexity $O(H^3/d_m\epsilon^2)$ [Yin et al., 2021] for non-stationary MDP and the optimal sample
496 complexity $O(H^2/d_m\epsilon^2)$ [Yin and Wang, 2021a] for stationary MDP, where d_m is the lower bound
497 for state-action occupancy. Such uniform convergence idea also supports some works regarding
498 online exploration [Jin et al., 2020a, Qiao et al., 2022]. For offline RL with single concentrability
499 assumption, Xie et al. [2021b] arrive at the optimal sample complexity $O(H^3SC^*/\epsilon^2)$. Recently,
500 Yin and Wang [2021b] propose APVI which can lead to instance-dependent sub-optimality bound,
501 which subsumes previous optimal results under several assumptions.

502 **Offline reinforcement learning under linear MDP.** Recently, many works focus on offline RL
503 under linear representation. Jin et al. [2021] present PEVI which applies the idea of pessimistic value
504 iteration (the idea originates from [Jin et al., 2020b]), and PEVI is provably efficient for offline RL
505 under linear MDP. Yin et al. [2022] improve the sub-optimality bound in [Jin et al., 2021] by replacing
506 LSVI by variance-weighted LSVI. Xie et al. [2021a] consider Bellman consistent pessimism for
507 general function approximation, and their result improves the sample complexity in [Jin et al., 2021]
508 by order $O(d)$ (shown in Theorem 3.2). However, there is no improvement on horizon dependence.
509 Zanette et al. [2021] propose a new offline actor-critic algorithm that naturally incorporates the
510 pessimism principle. Besides, Wang et al. [2021], Zanette [2021] study the statistical hardness of
511 offline RL with linear representations by presenting exponential lower bounds.

512 C More discussions

513 C.1 Discussions about definition of DP

514 For a concrete motivating example, please refer to the first paragraph of Introduction. We remark that
515 our definition of DP is consistent with Joint DP and Local DP defined under the online RL setting

516 where JDP/LDP also cast each user as one trajectory and provide user-wise privacy protection. For
 517 detailed definitions and more discussions about JDP/LDP, please refer to [Qiao and Wang \[2022\]](#).

518 C.2 Tightness of our results

519 We believe our bounds for offline RL with DP is tight. To the best of our knowledge, APVI and
 520 VAPVI provide the tightest bound under tabular MDP and linear MDP, respectively. The suboptimality
 521 bounds of our algorithms match these two in the main term, with some lower order additional terms.
 522 The leading terms are known to match multiple information-theoretical lower bounds for offline
 523 RL simultaneously (this was illustrated in [Yin and Wang \[2021b\]](#), [Yin et al. \[2022\]](#)), for this reason
 524 our bound cannot be improved in general. For the lower order terms, the dependence on sample
 525 complexity n and privacy budget ϵ : $\tilde{O}(\frac{1}{n\epsilon})$ is optimal since policy learning is a special case of ERM
 526 problems and such dependence is optimal in DP-ERM. In addition, we believe the dependence on
 527 other parameters (H, S, A, d) in the lower order term is tight due to our special tricks as (3) and
 528 Lemma E.6.

529 D Proof of Theorem 3.4

530 D.1 Proof sketch

531 Since the whole proof for privacy guarantee is not very complex, we present it in Section D.2 below
 532 and only sketch the proof for suboptimality bound.

533 First of all, we bound the scale of noises we add to show that the \tilde{n} derived from (3) are close to real
 534 visitation numbers. Therefore, denoting the non-private empirical transition kernel by \hat{P} (detailed
 535 definition in (15)), we can show that $\|\tilde{P} - \hat{P}\|_1$ and $|\sqrt{\text{Var}_{\tilde{P}}(V)} - \sqrt{\text{Var}_{\hat{P}}(V)}|$ are small.

536 Next, resulting from the conditional independence of \tilde{V}_{h+1} and \tilde{P}_h , we apply Empirical Bernstein's
 537 inequality to get $|(\tilde{P}_h - P_h)\tilde{V}_{h+1}| \lesssim \sqrt{\text{Var}_{\tilde{P}_h(\cdot|s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) / \tilde{n}_{s_h, a_h}} + SHE_\rho / \tilde{n}_{s_h, a_h}$. Together with our
 538 definition of private pessimism and the key lemma: extended value difference (Lemma F.7 and F.8),
 539 we can bound the suboptimality of our output policy $\hat{\pi}$ by:

$$v^* - v^{\hat{\pi}} \lesssim \sum_{h=1}^H \sum_{(s_h, a_h) \in \mathcal{C}_h} d_h^{\pi^*}(s_h, a_h) \sqrt{\frac{\text{Var}_{\tilde{P}_h(\cdot|s_h, a_h)}(\tilde{V}_{h+1}(\cdot))}{\tilde{n}_{s_h, a_h}}} + SHE_\rho / \tilde{n}_{s_h, a_h}. \quad (12)$$

540 Finally, we further bound the above suboptimality via replacing private statistics by non-private ones.
 541 Specifically, we replace \tilde{n} by n , \tilde{P} by P and \tilde{V} by V^* . Due to (12), we have $\|\tilde{V} - V^*\|_\infty \lesssim \sqrt{\frac{1}{n\bar{d}_m}}$.

542 Together with the upper bounds of $\|\tilde{P} - \hat{P}\|_1$ and $|\sqrt{\text{Var}_{\tilde{P}}(V)} - \sqrt{\text{Var}_{\hat{P}}(V)}|$, we have

$$\begin{aligned} & \sqrt{\frac{\text{Var}_{\tilde{P}_h(\cdot|s_h, a_h)}(\tilde{V}_{h+1}(\cdot))}{\tilde{n}_{s_h, a_h}}} \lesssim \sqrt{\frac{\text{Var}_{\tilde{P}_h(\cdot|s_h, a_h)}(V_{h+1}^*(\cdot))}{\tilde{n}_{s_h, a_h}}} + \frac{1}{n\bar{d}_m} \\ & \lesssim \sqrt{\frac{\text{Var}_{\hat{P}_h(\cdot|s_h, a_h)}(V_{h+1}^*(\cdot))}{\tilde{n}_{s_h, a_h}}} + \frac{1}{n\bar{d}_m} \lesssim \sqrt{\frac{\text{Var}_{P_h(\cdot|s_h, a_h)}(V_{h+1}^*(\cdot))}{\tilde{n}_{s_h, a_h}}} + \frac{1}{n\bar{d}_m} \\ & \lesssim \sqrt{\frac{\text{Var}_{P_h(\cdot|s_h, a_h)}(V_{h+1}^*(\cdot))}{nd_h^{\pi^*}(s_h, a_h)}} + \frac{1}{n\bar{d}_m}. \end{aligned} \quad (13)$$

543 The final bound using non-private statistics results from (12) and (13).

544 D.2 Proof of the privacy guarantee

545 The privacy guarantee of DP-APVI (Algorithm 1) is summarized by Lemma D.1 below.

546 **Lemma D.1** (Privacy analysis of DP-APVI (Algorithm 1)). *DP-APVI (Algorithm 1) satisfies ρ -zCDP.*

547 *Proof of Lemma D.1.* The ℓ_2 sensitivity of $\{n_{s_h, a_h}\}$ is $\sqrt{2H}$. According to Lemma 2.6, the Gaussian
548 Mechanism used on $\{n_{s_h, a_h}\}$ with $\sigma^2 = \frac{2H}{\rho}$ satisfies $\frac{\rho}{2}$ -zCDP. Similarly, the Gaussian Mechanism
549 used on $\{n_{s_h, a_h, s_{h+1}}\}$ with $\sigma^2 = \frac{2H}{\rho}$ also satisfies $\frac{\rho}{2}$ -zCDP. Combining these two results, due to the
550 composition of zCDP (Lemma F.16), the construction of $\{n'\}$ satisfies ρ -zCDP. Finally, DP-APVI
551 satisfies ρ -zCDP because the output $\hat{\pi}$ is post processing of $\{n'\}$. \square

552 D.3 Proof of the sub-optimality bound

553 D.3.1 Utility analysis

554 First of all, the following Lemma D.2 gives a high probability bound for $|n' - n|$.

555 **Lemma D.2.** Let $E_\rho = 2\sqrt{2}\sigma\sqrt{\log \frac{4HS^2A}{\delta}} = 4\sqrt{\frac{H \log \frac{4HS^2A}{\delta}}{\rho}}$, then with probability $1 - \delta$, for all
556 s_h, a_h, s_{h+1} , it holds that

$$|n'_{s_h, a_h} - n_{s_h, a_h}| \leq \frac{E_\rho}{2}, \quad |n'_{s_h, a_h, s_{h+1}} - n_{s_h, a_h, s_{h+1}}| \leq \frac{E_\rho}{2}. \quad (14)$$

557 *Proof of Lemma D.2.* The inequalities directly result from the concentration inequality of Gaussian
558 distribution and a union bound. \square

559 According to the utility analysis above, we have the following Lemma D.3 giving a high probability
560 bound for $|\tilde{n} - n|$.

Lemma D.3. Under the high probability event in Lemma D.2, for all s_h, a_h, s_{h+1} , it holds that

$$|\tilde{n}_{s_h, a_h} - n_{s_h, a_h}| \leq E_\rho, \quad |\tilde{n}_{s_h, a_h, s_{h+1}} - n_{s_h, a_h, s_{h+1}}| \leq E_\rho.$$

Proof of Lemma D.3. When the event in Lemma D.2 holds, the original counts $\{n_{s_h, a_h, s'}\}_{s' \in \mathcal{S}}$ is a
feasible solution to the optimization problem, which means that

$$\max_{s'} |\tilde{n}_{s_h, a_h, s'} - n'_{s_h, a_h, s'}| \leq \max_{s'} |n_{s_h, a_h, s'} - n'_{s_h, a_h, s'}| \leq \frac{E_\rho}{2}.$$

Due to the second part of (14), it holds that for any s_h, a_h, s_{h+1} ,

$$|\tilde{n}_{s_h, a_h, s_{h+1}} - n_{s_h, a_h, s_{h+1}}| \leq |\tilde{n}_{s_h, a_h, s_{h+1}} - n'_{s_h, a_h, s_{h+1}}| + |n'_{s_h, a_h, s_{h+1}} - n_{s_h, a_h, s_{h+1}}| \leq E_\rho.$$

For the second part, because of the constraints in the optimization problem, it holds that

$$|\tilde{n}_{s_h, a_h} - n'_{s_h, a_h}| \leq \frac{E_\rho}{2}.$$

Due to the first part of (14), it holds that for any s_h, a_h ,

$$|\tilde{n}_{s_h, a_h} - n_{s_h, a_h}| \leq |\tilde{n}_{s_h, a_h} - n'_{s_h, a_h}| + |n'_{s_h, a_h} - n_{s_h, a_h}| \leq E_\rho.$$

561 \square

562 Let the non-private empirical estimate be:

$$\hat{P}_h(s'|s_h, a_h) = \frac{n_{s_h, a_h, s'}}{n_{s_h, a_h}}, \quad (15)$$

563 if $n_{s_h, a_h} > 0$ and $\hat{P}_h(s'|s_h, a_h) = \frac{1}{S}$ otherwise. We will show that the private transition kernel \tilde{P} is
564 close to \hat{P} by the Lemma D.4 and Lemma D.5 below.

565 **Lemma D.4.** Under the high probability event of Lemma D.3, for s_h, a_h , if $\tilde{n}_{s_h, a_h} \geq 3E_\rho$, it holds
566 that

$$\left\| \tilde{P}_h(\cdot|s_h, a_h) - \hat{P}_h(\cdot|s_h, a_h) \right\|_1 \leq \frac{5SE_\rho}{\tilde{n}_{s_h, a_h}}. \quad (16)$$

567 *Proof of Lemma D.4.* If $\tilde{n}_{s_h, a_h} \geq 3E_\rho$ and the conclusion in Lemma D.3 hold, we have

$$\begin{aligned}
& \left\| \tilde{P}_h(\cdot | s_h, a_h) - \hat{P}_h(\cdot | s_h, a_h) \right\|_1 \leq \sum_{s' \in \mathcal{S}} \left| \tilde{P}_h(s' | s_h, a_h) - \hat{P}_h(s' | s_h, a_h) \right| \\
& \leq \sum_{s' \in \mathcal{S}} \left(\frac{\tilde{n}_{s_h, a_h, s'} + E_\rho}{\tilde{n}_{s_h, a_h} - E_\rho} - \frac{\tilde{n}_{s_h, a_h, s'}}{\tilde{n}_{s_h, a_h}} \right) \\
& \leq \sum_{s' \in \mathcal{S}} \left[\left(\frac{1}{\tilde{n}_{s_h, a_h}} + \frac{2E_\rho}{\tilde{n}_{s_h, a_h}^2} \right) (\tilde{n}_{s_h, a_h, s'} + E_\rho) - \frac{\tilde{n}_{s_h, a_h, s'}}{\tilde{n}_{s_h, a_h}} \right] \\
& \leq \frac{SE_\rho}{\tilde{n}_{s_h, a_h}} + \frac{2E_\rho}{\tilde{n}_{s_h, a_h}} + \frac{2SE_\rho^2}{\tilde{n}_{s_h, a_h}^2} \\
& \leq \frac{5SE_\rho}{\tilde{n}_{s_h, a_h}}.
\end{aligned} \tag{17}$$

568 The second inequality is because $\frac{\tilde{n}_{s_h, a_h, s'} - E_\rho}{\tilde{n}_{s_h, a_h} + E_\rho} \leq \frac{n_{s_h, a_h, s'}}{n_{s_h, a_h}} \leq \frac{\tilde{n}_{s_h, a_h, s'} + E_\rho}{\tilde{n}_{s_h, a_h} - E_\rho}$ and $\frac{\tilde{n}_{s_h, a_h, s'} + E_\rho}{\tilde{n}_{s_h, a_h} - E_\rho} -$
569 $\frac{\tilde{n}_{s_h, a_h, s'}}{\tilde{n}_{s_h, a_h}} \geq \frac{\tilde{n}_{s_h, a_h, s'}}{\tilde{n}_{s_h, a_h}} - \frac{\tilde{n}_{s_h, a_h, s'} - E_\rho}{\tilde{n}_{s_h, a_h} + E_\rho}$. The third inequality is because of Lemma F.6. The last
570 inequality is because $\tilde{n}_{s_h, a_h} \geq 3E_\rho$. \square

571 **Lemma D.5.** Let $V \in \mathbb{R}^S$ be any function with $\|V\|_\infty \leq H$, under the high probability event of
572 Lemma D.3, for s_h, a_h , if $\tilde{n}_{s_h, a_h} \geq 3E_\rho$, it holds that

$$\left| \sqrt{\text{Var}_{\tilde{P}_h(\cdot | s_h, a_h)}(V)} - \sqrt{\text{Var}_{\hat{P}_h(\cdot | s_h, a_h)}(V)} \right| \leq 4H \sqrt{\frac{SE_\rho}{\tilde{n}_{s_h, a_h}}}. \tag{18}$$

Proof of Lemma D.5. For s_h, a_h such that $\tilde{n}_{s_h, a_h} \geq 3E_\rho$, we use $\tilde{P}(\cdot)$ and $\hat{P}(\cdot)$ instead of $\tilde{P}_h(\cdot | s_h, a_h)$ and $\hat{P}_h(\cdot | s_h, a_h)$ for simplicity. Because of Lemma D.4, we have

$$\left\| \tilde{P}(\cdot) - \hat{P}(\cdot) \right\|_1 \leq \frac{5SE_\rho}{\tilde{n}_{s_h, a_h}}.$$

573 Therefore, it holds that

$$\begin{aligned}
& \left| \sqrt{\text{Var}_{\tilde{P}(\cdot)}(V)} - \sqrt{\text{Var}_{\hat{P}(\cdot)}(V)} \right| \leq \sqrt{|\text{Var}_{\tilde{P}(\cdot)}(V) - \text{Var}_{\hat{P}(\cdot)}(V)|} \\
& \leq \sqrt{\sum_{s' \in \mathcal{S}} \left| \hat{P}(s') - \tilde{P}(s') \right| V(s')^2 + \left| \sum_{s' \in \mathcal{S}} \left[\hat{P}(s') + \tilde{P}(s') \right] V(s') \right| \cdot \sum_{s' \in \mathcal{S}} \left| \hat{P}(s') - \tilde{P}(s') \right| V(s')} \\
& \leq \sqrt{H^2 \left\| \tilde{P}(\cdot) - \hat{P}(\cdot) \right\|_1 + 2H^2 \left\| \tilde{P}(\cdot) - \hat{P}(\cdot) \right\|_1} \\
& \leq 4H \sqrt{\frac{SE_\rho}{\tilde{n}_{s_h, a_h}}}.
\end{aligned} \tag{19}$$

574 The second inequality is due to the definition of variance. \square

575 D.3.2 Validity of our pessimistic penalty

576 Now we are ready to present the key lemma (Lemma D.6) below to justify our use of Γ as the
577 pessimistic penalty.

578 **Lemma D.6.** Under the high probability event of Lemma D.3, with probability $1 - \delta$, for any s_h, a_h ,
579 if $\tilde{n}_{s_h, a_h} \geq 3E_\rho$ (which implies $n_{s_h, a_h} > 0$), it holds that

$$\left| (\tilde{P}_h - P_h) \cdot \tilde{V}_{h+1}(s_h, a_h) \right| \leq \sqrt{\frac{2\text{Var}_{\tilde{P}_h(\cdot | s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) \cdot \iota}{\tilde{n}_{s_h, a_h} - E_\rho} + \frac{16SHE_\rho \cdot \iota}{\tilde{n}_{s_h, a_h}}}, \tag{20}$$

580 where \tilde{V} is the private version of estimated V function, which appears in Algorithm 1 and $\iota =$
581 $\log(HSA/\delta)$.

Proof of Lemma D.6.

$$\begin{aligned}
& \left| (\tilde{P}_h - P_h) \cdot \tilde{V}_{h+1}(s_h, a_h) \right| \leq \left| (\tilde{P}_h - \hat{P}_h) \cdot \tilde{V}_{h+1}(s_h, a_h) \right| + \left| (\hat{P}_h - P_h) \cdot \tilde{V}_{h+1}(s_h, a_h) \right| \\
& \leq H \left\| \tilde{P}_h(\cdot | s_h, a_h) - \hat{P}_h(\cdot | s_h, a_h) \right\|_1 + \left| (\hat{P}_h - P_h) \cdot \tilde{V}_{h+1}(s_h, a_h) \right| \\
& \leq \frac{5SHE_\rho}{\tilde{n}_{s_h, a_h}} + \left| (\hat{P}_h - P_h) \cdot \tilde{V}_{h+1}(s_h, a_h) \right|,
\end{aligned} \tag{21}$$

582 where the third inequality is due to Lemma D.4.

583 Next, recall $\hat{\pi}_{h+1}$ in Algorithm 1 is computed backwardly therefore only depends on sample tuple
584 from time $h + 1$ to H . As a result, $\tilde{V}_{h+1} = \langle \tilde{Q}_{h+1}, \hat{\pi}_{h+1} \rangle$ also only depends on the sample tuple
585 from time $h + 1$ to H and some Gaussian noise that is independent to the offline dataset. On the
586 other side, by the definition, \hat{P}_h only depends on the sample tuples from time h to $h + 1$. Therefore
587 \tilde{V}_{h+1} and \hat{P}_h are *Conditionally* independent (This trick is also used in [Yin et al., 2021] and [Yin and
588 Wang, 2021b]), by Empirical Bernstein's inequality (Lemma F.4) and a union bound, with probability
589 $1 - \delta$, for all s_h, a_h such that $\tilde{n}_{s_h, a_h} \geq 3E_\rho$,

$$\left| (\hat{P}_h - P_h) \cdot \tilde{V}_{h+1}(s_h, a_h) \right| \leq \sqrt{\frac{2\text{Var}_{\hat{P}_h(\cdot | s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) \cdot \iota}{n_{s_h, a_h}}} + \frac{7H \cdot \iota}{3n_{s_h, a_h}}. \tag{22}$$

590 Therefore, we have

$$\begin{aligned}
& \left| (\tilde{P}_h - P_h) \cdot \tilde{V}_{h+1}(s_h, a_h) \right| \leq \sqrt{\frac{2\text{Var}_{\hat{P}_h(\cdot | s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) \cdot \iota}{n_{s_h, a_h}}} + \frac{7H \cdot \iota}{3n_{s_h, a_h}} + \frac{5SHE_\rho}{\tilde{n}_{s_h, a_h}} \\
& \leq \sqrt{\frac{2\text{Var}_{\hat{P}_h(\cdot | s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) \cdot \iota}{n_{s_h, a_h}}} + \frac{9SHE_\rho \cdot \iota}{\tilde{n}_{s_h, a_h}} \\
& \leq \frac{9SHE_\rho \cdot \iota}{\tilde{n}_{s_h, a_h}} + \sqrt{\frac{2\text{Var}_{\hat{P}_h(\cdot | s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) \cdot \iota}{n_{s_h, a_h}}} + 4\sqrt{2H} \sqrt{\frac{SE_\rho \cdot \iota}{\tilde{n}_{s_h, a_h} \cdot n_{s_h, a_h}}} \\
& \leq \sqrt{\frac{2\text{Var}_{\hat{P}_h(\cdot | s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) \cdot \iota}{n_{s_h, a_h}}} + \frac{16SHE_\rho \cdot \iota}{\tilde{n}_{s_h, a_h}} \\
& \leq \sqrt{\frac{2\text{Var}_{\hat{P}_h(\cdot | s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) \cdot \iota}{\tilde{n}_{s_h, a_h} - E_\rho}} + \frac{16SHE_\rho \cdot \iota}{\tilde{n}_{s_h, a_h}}.
\end{aligned} \tag{23}$$

591 The second and forth inequality is because when $\tilde{n}_{s_h, a_h} \geq 3E_\rho$, $n_{s_h, a_h} \geq \frac{2\tilde{n}_{s_h, a_h}}{3}$. Specifically,
592 these two inequalities are also because usually we only care about the case when $SE_\rho \geq 1$, which is
593 equivalent to ρ being not very large. The third inequality is due to Lemma D.5. The last inequality is
594 due to Lemma D.3. \square

595 Note that the previous Lemmas rely on the condition that \tilde{n} is not very small ($\tilde{n}_{s_h, a_h} \geq 3E_\rho$). Below
596 we state the Multiplicative Chernoff bound (Lemma D.7 and Remark D.8) to show that under our
597 condition in Theorem 3.4, for $(s_h, a_h) \in \mathcal{C}_h$, \tilde{n}_{s_h, a_h} will be larger than $3E_\rho$ with high probability.

598 **Lemma D.7** (Lemma B.1 in [Yin and Wang, 2021b]). *For any $0 < \delta < 1$, there exists an absolute*
599 *constant c_1 such that when total episode $n > c_1 \cdot 1/d_m \cdot \log(HSA/\delta)$, then with probability $1 - \delta$,*
600 *$\forall h \in [H]$*

$$n_{s_h, a_h} \geq n \cdot d_h^\mu(s_h, a_h)/2, \quad \forall (s_h, a_h) \in \mathcal{C}_h.$$

601 Furthermore, we denote

$$\mathcal{E} := \{n_{s_h, a_h} \geq n \cdot d_h^\mu(s_h, a_h)/2, \forall (s_h, a_h) \in \mathcal{C}_h, h \in [H]\}. \tag{24}$$

602 then equivalently $P(\mathcal{E}) > 1 - \delta$.

603 In addition, we denote

$$\mathcal{E}' := \{n_{s_h, a_h} \leq \frac{3}{2}n \cdot d_h^\mu(s_h, a_h), \forall (s_h, a_h) \in \mathcal{C}_h, h \in [H]\}. \quad (25)$$

604 then similarly $P(\mathcal{E}') > 1 - \delta$.

605 **Remark D.8.** According to Lemma D.7, for any failure probability δ , there exists some constant
 606 $c_1 > 0$ such that when $n \geq \frac{c_1 E_\rho \cdot \iota}{d_m}$, with probability $1 - \delta$, for all $(s_h, a_h) \in \mathcal{C}_h$, $n_{s_h, a_h} \geq 4E_\rho$.
 607 Therefore, under the condition of Theorem 3.4 and the high probability events in Lemma D.3
 608 and Lemma D.7, it holds that for all $(s_h, a_h) \in \mathcal{C}_h$, $\tilde{n}_{s_h, a_h} \geq 3E_\rho$ while for all $(s_h, a_h) \notin \mathcal{C}_h$,
 609 $\tilde{n}_{s_h, a_h} \leq E_\rho$.

610 **Lemma D.9.** Define $(\mathcal{T}_h V)(\cdot, \cdot) := r_h(\cdot, \cdot) + (P_h V)(\cdot, \cdot)$ for any $V \in \mathbb{R}^S$. Note $\hat{\pi}$, \bar{Q}_h , \tilde{V}_h are
 611 defined in Algorithm 1 and denote $\xi_h(s, a) = (\mathcal{T}_h \tilde{V}_{h+1})(s, a) - \bar{Q}_h(s, a)$. Then it holds that

$$V_1^{\pi^*}(s) - V_1^{\hat{\pi}}(s) \leq \sum_{h=1}^H \mathbb{E}_{\pi^*} [\xi_h(s_h, a_h) \mid s_1 = s] - \sum_{h=1}^H \mathbb{E}_{\hat{\pi}} [\xi_h(s_h, a_h) \mid s_1 = s]. \quad (26)$$

612 Furthermore, (26) holds for all $V_h^{\pi^*}(s) - V_h^{\hat{\pi}}(s)$.

613 *Proof of Lemma D.9.* Lemma D.9 is a direct corollary of Lemma F.8 with $\pi = \pi^*$, $\hat{Q}_h = \bar{Q}_h$,
 614 $\hat{V}_h = \tilde{V}_h$ and $\hat{\pi} = \hat{\pi}$ in Algorithm 1, we can obtain this result since by the definition of $\hat{\pi}$ in
 615 Algorithm 1, $\langle \bar{Q}_h(s_h, \cdot), \pi_h(\cdot | s_h) - \hat{\pi}_h(\cdot | s_h) \rangle \leq 0$. The proof for $V_h^{\pi^*}(s) - V_h^{\hat{\pi}}(s)$ is identical. \square

616 Next we prove the asymmetric bound for ξ_h , which is the key to the proof.

617 **Lemma D.10** (Private version of Lemma D.6 in [Yin and Wang, 2021b]). Denote $\xi_h(s, a) =$
 618 $(\mathcal{T}_h \tilde{V}_{h+1})(s, a) - \bar{Q}_h(s, a)$, where \tilde{V}_{h+1} and \bar{Q}_h are the quantities in Algorithm 1 and $\mathcal{T}_h(V) :=$
 619 $r_h + P_h \cdot V$ for any $V \in \mathbb{R}^S$. Then under the high probability events in Lemma D.3 and Lemma D.6,
 620 for any h, s_h, a_h such that $\tilde{n}_{s_h, a_h} > 3E_\rho$, we have

$$\begin{aligned} 0 &\leq \xi_h(s_h, a_h) = (\mathcal{T}_h \tilde{V}_{h+1})(s_h, a_h) - \bar{Q}_h(s_h, a_h) \\ &\leq 2\sqrt{\frac{2\text{Var}_{\tilde{P}_h(\cdot | s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) \cdot \iota}{\tilde{n}_{s_h, a_h} - E_\rho} + \frac{32SHE_\rho \cdot \iota}{\tilde{n}_{s_h, a_h}}}, \end{aligned}$$

621 where $\iota = \log(HSA/\delta)$.

622 *Proof of Lemma D.10. The first inequality:* We first prove $\xi_h(s_h, a_h) \geq 0$ for all (s_h, a_h) , such
 623 that $\tilde{n}_{s_h, a_h} \geq 3E_\rho$.

624 Indeed, if $\hat{Q}_h^p(s_h, a_h) < 0$, then $\bar{Q}_h(s_h, a_h) = 0$. In this case, $\xi_h(s_h, a_h) = (\mathcal{T}_h \tilde{V}_{h+1})(s_h, a_h) \geq 0$
 625 (note $\tilde{V}_h \geq 0$ by the definition). If $\hat{Q}_h^p(s_h, a_h) \geq 0$, then by definition $\bar{Q}_h(s_h, a_h) =$
 626 $\min\{\hat{Q}_h^p(s_h, a_h), H - h + 1\}^+ \leq \hat{Q}_h^p(s_h, a_h)$ and this implies

$$\begin{aligned} \xi_h(s_h, a_h) &\geq (\mathcal{T}_h \tilde{V}_{h+1})(s_h, a_h) - \hat{Q}_h^p(s_h, a_h) \\ &= (P_h - \tilde{P}_h) \cdot \tilde{V}_{h+1}(s_h, a_h) + \Gamma_h(s_h, a_h) \\ &\geq -\sqrt{\frac{2\text{Var}_{\tilde{P}_h(\cdot | s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) \cdot \iota}{\tilde{n}_{s_h, a_h} - E_\rho} - \frac{16SHE_\rho \cdot \iota}{\tilde{n}_{s_h, a_h}}} + \Gamma_h(s_h, a_h) = 0, \end{aligned}$$

627 where the second inequality uses Lemma D.6, and the last equation uses Line 5 of Algorithm 1.

628 **The second inequality:** Then we prove $\xi_h(s_h, a_h) \leq 2\sqrt{\frac{2\text{Var}_{\tilde{P}_h(\cdot | s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) \cdot \iota}{\tilde{n}_{s_h, a_h} - E_\rho} + \frac{32SHE_\rho \cdot \iota}{\tilde{n}_{s_h, a_h}}}$ for
 629 all (s_h, a_h) such that $\tilde{n}_{s_h, a_h} \geq 3E_\rho$.

630 First, since by construction $\tilde{V}_h \leq H - h + 1$ for all $h \in [H]$, this implies

$$\hat{Q}_h^p = \tilde{Q}_h - \Gamma_h \leq \tilde{Q}_h = r_h + (\tilde{P}_h \cdot \tilde{V}_{h+1}) \leq 1 + (H - h) = H - h + 1$$

631 which is because $r_h \leq 1$ and \tilde{P}_h is a probability distribution. Therefore, we have the equivalent
 632 definition

$$\bar{Q}_h := \min\{\widehat{Q}_h^p, H - h + 1\}^+ = \max\{\widehat{Q}_h^p, 0\} \geq \widehat{Q}_h^p.$$

633 Then it holds that

$$\begin{aligned} \xi_h(s_h, a_h) &= (\mathcal{T}_h \tilde{V}_{h+1})(s_h, a_h) - \bar{Q}_h(s_h, a_h) \leq (\mathcal{T}_h \tilde{V}_{h+1})(s_h, a_h) - \widehat{Q}_h^p(s_h, a_h) \\ &= (\mathcal{T}_h \tilde{V}_{h+1})(s_h, a_h) - \tilde{Q}_h(s_h, a_h) + \Gamma_h(s_h, a_h) \\ &= (P_h - \tilde{P}_h) \cdot \tilde{V}_{h+1}(s_h, a_h) + \Gamma_h(s_h, a_h) \\ &\leq \sqrt{\frac{2\text{Var}_{\tilde{P}_h(\cdot|s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) \cdot \iota}{\tilde{n}_{s_h, a_h} - E_\rho} + \frac{16SHE_\rho \cdot \iota}{\tilde{n}_{s_h, a_h}}} + \Gamma_h(s_h, a_h) \\ &= 2\sqrt{\frac{2\text{Var}_{\tilde{P}_h(\cdot|s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) \cdot \iota}{\tilde{n}_{s_h, a_h} - E_\rho} + \frac{32SHE_\rho \cdot \iota}{\tilde{n}_{s_h, a_h}}}. \end{aligned}$$

634 The proof is complete by combining the two parts. \square

635 D.3.3 Reduction to augmented absorbing MDP

636 Before we prove the theorem, we need to construct an augmented absorbing MDP to bridge \tilde{V} and
 637 V^* . According to Line 5 in Algorithm 1, the locations with $\tilde{n}_{s_h, a_h} \leq E_\rho$ is heavily penalized with
 638 penalty of order $\tilde{O}(H)$. Therefore we can prove that under the high probability event in Remark D.8,
 639 $d_h^{\hat{\pi}}(s_h, a_h) > 0$ only if $d_h^\mu(s_h, a_h) > 0$ by induction, where $\hat{\pi}$ is the output of Algorithm 1. The
 640 conclusion holds for $h = 1$. Assume it holds for some $h > 1$ that $d_h^{\hat{\pi}}(s_h, a_h) > 0$ only if $d_h^\mu(s_h, a_h) >$
 641 0 , then for any $s_{h+1} \in \mathcal{S}$ such that $d_{h+1}^{\hat{\pi}}(s_{h+1}) > 0$, it holds that $d_{h+1}^\mu(s_{h+1}) > 0$, which leads to
 642 the conclusion that $d_{h+1}^{\hat{\pi}}(s_{h+1}, a_{h+1}) > 0$ only if $d_{h+1}^\mu(s_{h+1}, a_{h+1}) > 0$. To summarize, we have

$$d_h^{\pi_0}(s_h, a_h) > 0 \text{ only if } d_h^\mu(s_h, a_h) > 0, \pi_0 \in \{\pi^*, \hat{\pi}\}. \quad (27)$$

643 Let us define M^\dagger by adding one absorbing state s_h^\dagger for all $h \in \{2, \dots, H\}$, therefore the augmented
 644 state space $\mathcal{S}^\dagger = \mathcal{S} \cup \{s_h^\dagger\}$ and the transition and reward is defined as follows: (recall $\mathcal{C}_h :=$
 645 $\{(s_h, a_h) : d_h^\mu(s_h, a_h) > 0\}$)

$$P_h^\dagger(\cdot | s_h, a_h) = \begin{cases} P_h(\cdot | s_h, a_h) & s_h, a_h \in \mathcal{C}_h, \\ \delta_{s_{h+1}^\dagger} & s_h = s_h^\dagger \text{ or } s_h, a_h \notin \mathcal{C}_h, \end{cases} \quad r_h^\dagger(s_h, a_h) = \begin{cases} r_h(s_h, a_h) & s_h, a_h \in \mathcal{C}_h \\ 0 & s_h = s_h^\dagger \text{ or } s_h, a_h \notin \mathcal{C}_h \end{cases}$$

646 and we further define for any π ,

$$V_h^{\dagger\pi}(s) = \mathbb{E}_\pi^\dagger \left[\sum_{t=h}^H r_t^\dagger \middle| s_h = s \right], v^{\dagger\pi} = \mathbb{E}_\pi^\dagger \left[\sum_{t=1}^H r_t^\dagger \right] \quad \forall h \in [H], \quad (28)$$

647 where \mathbb{E}^\dagger means taking expectation under the absorbing MDP M^\dagger .

648 Note that because π^* and $\hat{\pi}$ are fully covered by μ (27), it holds that

$$v^{\dagger\pi^*} = v^{\pi^*}, \quad v^{\dagger\hat{\pi}} = v^{\hat{\pi}}. \quad (29)$$

649 Define $(\mathcal{T}_h^\dagger V)(\cdot, \cdot) := r_h^\dagger(\cdot, \cdot) + (P_h^\dagger V)(\cdot, \cdot)$ for any $V \in \mathbb{R}^{\mathcal{S}^\dagger}$. Note $\hat{\pi}, \bar{Q}_h, \tilde{V}_h$ are defined
 650 in Algorithm 1 (we extend the definition by letting $\tilde{V}_h(s_h^\dagger) = 0$ and $\bar{Q}_h(s_h^\dagger, \cdot) = 0$) and denote
 651 $\xi_h^\dagger(s, a) = (\mathcal{T}_h^\dagger \tilde{V}_{h+1})(s, a) - \bar{Q}_h(s, a)$. Using identical proof to Lemma D.9, we have

$$V_1^{\dagger\pi^*}(s) - V_1^{\dagger\hat{\pi}}(s) \leq \sum_{h=1}^H \mathbb{E}_{\pi^*}^\dagger \left[\xi_h^\dagger(s_h, a_h) \middle| s_1 = s \right] - \sum_{h=1}^H \mathbb{E}_{\hat{\pi}}^\dagger \left[\xi_h^\dagger(s_h, a_h) \middle| s_1 = s \right], \quad (30)$$

652 where $V_1^{\dagger\pi}$ is defined in (28). Furthermore, (30) holds for all $V_h^{\dagger\pi^*}(s) - V_h^{\dagger\hat{\pi}}(s)$.

653 **D.3.4 Finalize our result with non-private statistics**

654 For those $(s_h, a_h) \in \mathcal{C}_h$, $\xi_h^\dagger(s_h, a_h) = r_h(s_h, a_h) + P_h \tilde{V}_{h+1}(s_h, a_h) - \bar{Q}_h(s_h, a_h) = \xi_h(s_h, a_h)$.

655 For those $(s_h, a_h) \notin \mathcal{C}_h$ or $s_h = s_h^\dagger$, we have $\xi_h^\dagger(s_h, a_h) = 0$.

656 Therefore, by (30) and Lemma D.10, under the high probability events in Lemma D.3, Lemma D.6
657 and Lemma D.7, we have for all $t \in [H]$, $s \in \mathcal{S}$ (\mathcal{S} does not include the absorbing state s_t^\dagger),

$$\begin{aligned}
V_t^{\dagger\pi^*}(s) - V_t^{\dagger\hat{\pi}}(s) &\leq \sum_{h=t}^H \mathbb{E}_{\pi^*}^\dagger \left[\xi_h^\dagger(s_h, a_h) \mid s_t = s \right] - \sum_{h=t}^H \mathbb{E}_{\hat{\pi}}^\dagger \left[\xi_h^\dagger(s_h, a_h) \mid s_t = s \right] \\
&\leq \sum_{h=t}^H \mathbb{E}_{\pi^*}^\dagger \left[\xi_h^\dagger(s_h, a_h) \mid s_t = s \right] - 0 \\
&\leq \sum_{h=t}^H \mathbb{E}_{\pi^*}^\dagger \left[2 \sqrt{\frac{2\text{Var}_{\tilde{P}_h(\cdot|s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) \cdot \iota}{\tilde{n}_{s_h, a_h} - E_\rho} + \frac{32SHE_\rho \cdot \iota}{\tilde{n}_{s_h, a_h}}} \mid s_t = s \right] \cdot \mathbb{1}((s_h, a_h) \in \mathcal{C}_h) \\
&\leq \sum_{h=t}^H \mathbb{E}_{\pi^*}^\dagger \left[2 \sqrt{\frac{2\text{Var}_{\tilde{P}_h(\cdot|s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) \cdot \iota}{n_{s_h, a_h} - 2E_\rho} + \frac{32SHE_\rho \cdot \iota}{n_{s_h, a_h} - E_\rho}} \mid s_t = s \right] \cdot \mathbb{1}((s_h, a_h) \in \mathcal{C}_h) \quad (31) \\
&\leq \sum_{h=t}^H \mathbb{E}_{\pi^*}^\dagger \left[4 \sqrt{\frac{\text{Var}_{\tilde{P}_h(\cdot|s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) \cdot \iota}{n_{s_h, a_h}} + \frac{128SHE_\rho \cdot \iota}{3n_{s_h, a_h}}} \mid s_t = s \right] \cdot \mathbb{1}((s_h, a_h) \in \mathcal{C}_h) \\
&\leq \sum_{h=t}^H \mathbb{E}_{\pi^*}^\dagger \left[4 \sqrt{\frac{2\text{Var}_{\tilde{P}_h(\cdot|s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) \cdot \iota}{nd_h^\mu(s_h, a_h)} + \frac{256SHE_\rho \cdot \iota}{3nd_h^\mu(s_h, a_h)}} \mid s_t = s \right] \cdot \mathbb{1}((s_h, a_h) \in \mathcal{C}_h)
\end{aligned}$$

658 The second and third inequality are because of Lemma D.10, Remark D.8 and the fact that either
659 $\xi^\dagger = 0$ or $\xi^\dagger = \xi$ while $(s_h, a_h) \in \mathcal{C}_h$. The fourth inequality is due to Lemma D.3. The fifth inequality
660 is because of Remark D.8. The last inequality is by Lemma D.7.

661 Below we present a crude bound of $\left| V_t^{\dagger\pi^*}(s) - \tilde{V}_t(s) \right|$, which can be further used to bound the main
662 term in the main result.

663 **Lemma D.11** (Self-bounding, private version of Lemma D.7 in [Yin and Wang, 2021b]). *Under the*
664 *high probability events in Lemma D.3, Lemma D.6 and Lemma D.7, it holds that for all $t \in [H]$ and*
665 *$s \in \mathcal{S}$,*

$$\left| V_t^{\dagger\pi^*}(s) - \tilde{V}_t(s) \right| \leq \frac{4\sqrt{2\iota}H^2}{\sqrt{n \cdot \bar{d}_m}} + \frac{256SH^2E_\rho \cdot \iota}{3n \cdot \bar{d}_m}.$$

666 where \bar{d}_m is defined in Theorem 3.4.

667 *Proof of Lemma D.11.* According to (31), since $\text{Var}_{\tilde{P}_h(\cdot|s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) \leq H^2$, we have for all
668 $t \in [H]$,

$$\left| V_t^{\dagger\pi^*}(s) - V_t^{\dagger\hat{\pi}}(s) \right| \leq \frac{4\sqrt{2\iota}H^2}{\sqrt{n \cdot \bar{d}_m}} + \frac{256SH^2E_\rho \cdot \iota}{3n \cdot \bar{d}_m} \quad (32)$$

669 Next, apply Lemma F.7 by setting $\pi = \hat{\pi}$, $\pi' = \pi^*$, $\hat{Q} = \bar{Q}$, $\hat{V} = \tilde{V}$ under M^\dagger , then we have

$$\begin{aligned}
V_t^{\dagger\pi^*}(s) - \tilde{V}_t(s) &= \sum_{h=t}^H \mathbb{E}_{\pi^*}^\dagger \left[\xi_h^\dagger(s_h, a_h) \mid s_t = s \right] + \sum_{h=t}^H \mathbb{E}_{\pi^*}^\dagger \left[(\bar{Q}_h(s_h, \cdot), \pi_h^*(\cdot | s_h) - \hat{\pi}_h(\cdot | s_h)) \mid s_t = s \right] \\
&\leq \sum_{h=t}^H \mathbb{E}_{\pi^*}^\dagger \left[\xi_h^\dagger(s_h, a_h) \mid s_t = s \right] \\
&\leq \frac{4\sqrt{2\iota}H^2}{\sqrt{n \cdot \bar{d}_m}} + \frac{256SH^2E_\rho \cdot \iota}{3n \cdot \bar{d}_m}.
\end{aligned} \tag{33}$$

670 Also, apply Lemma F.7 by setting $\pi = \pi' = \hat{\pi}$, $\hat{Q} = \bar{Q}$, $\hat{V} = \tilde{V}$ under M^\dagger , then we have

$$\tilde{V}_t(s) - V_t^{\dagger\hat{\pi}}(s) = - \sum_{h=t}^H \mathbb{E}_{\hat{\pi}}^\dagger \left[\xi_h^\dagger(s_h, a_h) \mid s_t = s \right] \leq 0. \tag{34}$$

671 The proof is complete by combing (32), (33) and (34). \square

672 Now we are ready to bound $\sqrt{\text{Var}_{\tilde{P}_h(\cdot | s_h, a_h)}(\tilde{V}_{h+1}(\cdot))}$ by $\sqrt{\text{Var}_{P_h(\cdot | s_h, a_h)}(V_{h+1}^{\dagger\pi^*}(\cdot))}$. Under the
673 high probability events in Lemma D.3, Lemma D.6 and Lemma D.7, with probability $1 - \delta$, it holds
674 that for all $(s_h, a_h) \in \mathcal{C}_h$,

$$\begin{aligned}
\sqrt{\text{Var}_{\tilde{P}_h(\cdot | s_h, a_h)}(\tilde{V}_{h+1}(\cdot))} &\leq \sqrt{\text{Var}_{\tilde{P}_h(\cdot | s_h, a_h)}(V_{h+1}^{\dagger\pi^*}(\cdot))} + \left\| \tilde{V}_{h+1} - V_{h+1}^{\dagger\pi^*} \right\|_{\infty, s \in \mathcal{S}} \\
&\leq \sqrt{\text{Var}_{\tilde{P}_h(\cdot | s_h, a_h)}(V_{h+1}^{\dagger\pi^*}(\cdot))} + \frac{4\sqrt{2\iota}H^2}{\sqrt{n \cdot \bar{d}_m}} + \frac{256SH^2E_\rho \cdot \iota}{3n \cdot \bar{d}_m} \\
&\leq \sqrt{\text{Var}_{\hat{P}_h(\cdot | s_h, a_h)}(V_{h+1}^{\dagger\pi^*}(\cdot))} + \frac{4\sqrt{2\iota}H^2}{\sqrt{n \cdot \bar{d}_m}} + \frac{256SH^2E_\rho \cdot \iota}{3n \cdot \bar{d}_m} + 4H\sqrt{\frac{SE_\rho}{\tilde{n}_{s_h, a_h}}} \\
&\leq \sqrt{\text{Var}_{\hat{P}_h(\cdot | s_h, a_h)}(V_{h+1}^{\dagger\pi^*}(\cdot))} + \frac{4\sqrt{2\iota}H^2}{\sqrt{n \cdot \bar{d}_m}} + \frac{256SH^2E_\rho \cdot \iota}{3n \cdot \bar{d}_m} + 8H\sqrt{\frac{SE_\rho}{n \cdot \bar{d}_m}} \\
&\leq \sqrt{\text{Var}_{P_h(\cdot | s_h, a_h)}(V_{h+1}^{\dagger\pi^*}(\cdot))} + \frac{4\sqrt{2\iota}H^2}{\sqrt{n \cdot \bar{d}_m}} + \frac{256SH^2E_\rho \cdot \iota}{3n \cdot \bar{d}_m} + 8H\sqrt{\frac{SE_\rho}{n \cdot \bar{d}_m}} + 3H\sqrt{\frac{\iota}{n \cdot \bar{d}_m}} \\
&\leq \sqrt{\text{Var}_{P_h(\cdot | s_h, a_h)}(V_{h+1}^{\dagger\pi^*}(\cdot))} + \frac{9\sqrt{\iota}H^2}{\sqrt{n \cdot \bar{d}_m}} + \frac{256SH^2E_\rho \cdot \iota}{3n \cdot \bar{d}_m} + 8H\sqrt{\frac{SE_\rho}{n \cdot \bar{d}_m}}.
\end{aligned} \tag{35}$$

675 The second inequality is because of Lemma D.11. The third inequality is due to Lemma D.5. The
676 fourth inequality comes from Lemma D.3 and Remark D.8. The fifth inequality holds with probability
677 $1 - \delta$ because of Lemma F.5 and a union bound.

678 Finally, by plugging (35) into (31) and averaging over s_1 , we finally have with probability $1 - 4\delta$,

$$\begin{aligned}
v^{\pi^*} - v^{\hat{\pi}} &= v^{\dagger\pi^*} - v^{\dagger\hat{\pi}} \leq \sum_{h=1}^H \mathbb{E}_{\pi^*}^\dagger \left[4\sqrt{\frac{2\text{Var}_{\tilde{P}_h(\cdot | s_h, a_h)}(\tilde{V}_{h+1}(\cdot)) \cdot \iota}{nd_h^\mu(s_h, a_h)} + \frac{256SHE_\rho \cdot \iota}{3nd_h^\mu(s_h, a_h)}} \right] \\
&\leq 4\sqrt{2} \sum_{h=1}^H \mathbb{E}_{\pi^*}^\dagger \left[\sqrt{\frac{\text{Var}_{P_h(\cdot | s_h, a_h)}(V_{h+1}^{\dagger\pi^*}(\cdot)) \cdot \iota}{nd_h^\mu(s_h, a_h)}} \right] + \tilde{O}\left(\frac{H^3 + SH^2E_\rho}{n \cdot \bar{d}_m}\right) \\
&= 4\sqrt{2} \sum_{h=1}^H \sum_{(s_h, a_h) \in \mathcal{C}_h} d_h^{\pi^*}(s_h, a_h) \sqrt{\frac{\text{Var}_{P_h(\cdot | s_h, a_h)}(V_{h+1}^{\dagger\pi^*}(\cdot)) \cdot \iota}{nd_h^\mu(s_h, a_h)}} + \tilde{O}\left(\frac{H^3 + SH^2E_\rho}{n \cdot \bar{d}_m}\right) \\
&= 4\sqrt{2} \sum_{h=1}^H \sum_{(s_h, a_h) \in \mathcal{C}_h} d_h^{\pi^*}(s_h, a_h) \sqrt{\frac{\text{Var}_{P_h(\cdot | s_h, a_h)}(V_{h+1}^{\dagger\pi^*}(\cdot)) \cdot \iota}{nd_h^\mu(s_h, a_h)}} + \tilde{O}\left(\frac{H^3 + SH^2E_\rho}{n \cdot \bar{d}_m}\right),
\end{aligned} \tag{36}$$

679 where \tilde{O} absorbs constants and Polylog terms. The first equation is due to (29). The first inequality is
680 because of (31). The second inequality comes from (35) and our assumption that $n \cdot \tilde{d}_m \geq c_1 H^2$.
681 The second equation uses the fact that $d_h^{\pi^*}(s_h, a_h) = d_h^{\dagger \pi^*}(s_h, a_h)$, for all (s_h, a_h) . The last
682 equation is because for any (s_h, a_h, s_{h+1}) such that $d_h^{\pi^*}(s_h, a_h) > 0$ and $P_h(s_{h+1}|s_h, a_h) > 0$,
683 $V_{h+1}^{\dagger \pi^*}(s_{h+1}) = V_{h+1}^*(s_{h+1})$.

684 D.4 Put everything together

685 Combining Lemma D.1 and (36), the proof of Theorem 3.4 is complete.

686 E Proof of Theorem 4.1

687 E.1 Proof sketch

688 Since the whole proof for privacy guarantee is not very complex, we present it in Section E.2 below
689 and only sketch the proof for suboptimality bound.

690 First of all, by extended value difference (Lemma F.7 and F.8), we can convert bounding the subopti-
691 mality gap of $v^* - v^{\hat{\pi}}$ to bounding $\sum_{h=1}^H \mathbb{E}_\pi [\Gamma_h(s_h, a_h)]$, given that $|(\mathcal{T}_h \tilde{V}_{h+1} - \tilde{\mathcal{T}}_h \tilde{V}_{h+1})(s, a)| \leq$
692 $\Gamma_h(s, a)$ for all s, a, h . To bound $(\mathcal{T}_h \tilde{V}_{h+1} - \tilde{\mathcal{T}}_h \tilde{V}_{h+1})(s, a)$, according to our analysis about the
693 upper bound of the noises we add, we can decompose $(\mathcal{T}_h \tilde{V}_{h+1} - \tilde{\mathcal{T}}_h \tilde{V}_{h+1})(s, a)$ to lower order terms
694 ($\tilde{O}(\frac{1}{K})$) and the following key quantity:

$$\phi(s, a)^\top \hat{\Lambda}_h^{-1} \left[\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \cdot \left(r_h^\tau + \tilde{V}_{h+1}(s_{h+1}^\tau) - \left(\mathcal{T}_h \tilde{V}_{h+1} \right)(s_h^\tau, a_h^\tau) \right) / \tilde{\sigma}_h^2(s_h^\tau, a_h^\tau) \right]. \quad (37)$$

695 For the term above, we prove an upper bound of $\left\| \sigma_{\tilde{V}_{h+1}}^2 - \tilde{\sigma}_h^2 \right\|_\infty$, so we can convert $\tilde{\sigma}_h^2$ to $\sigma_{\tilde{V}_{h+1}}^2$.

696 Next, since $\text{Var} \left[r_h^\tau + \tilde{V}_{h+1}(s_{h+1}^\tau) - \left(\mathcal{T}_h \tilde{V}_{h+1} \right)(s_h^\tau, a_h^\tau) \mid s_h^\tau, a_h^\tau \right] \approx \sigma_{\tilde{V}_{h+1}}^2$, we can apply Bern-
697 stein's inequality for self-normalized martingale (Lemma F.10) as in Yin et al. [2022] for deriving
698 tighter bound.

699 Finally, we replace the private statistics by non-private ones. More specifically, we convert $\sigma_{\tilde{V}_{h+1}}^2$ to
700 σ_h^{*2} (Λ_h^{-1} to Λ_h^{*-1}) by combining the crude upper bound of $\left\| \tilde{\mathcal{V}} - V^* \right\|_\infty$ and matrix concentrations.

701 E.2 Proof of the privacy guarantee

702 The privacy guarantee of DP-VAPVI (Algorithm 2) is summarized by Lemma E.1 below.

703 **Lemma E.1** (Privacy analysis of DP-VAPVI (Algorithm 2)). *DP-VAPVI (Algorithm 2) satisfies*
704 *ρ -zCDP.*

Proof of Lemma E.1. For $\sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \cdot \tilde{V}_{h+1}(\bar{s}_{h+1}^\tau)^2$, the ℓ_2 sensitivity is $2H^2$. For
 $\sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \cdot \tilde{V}_{h+1}(\bar{s}_{h+1}^\tau)$ and $\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \cdot \left(r_h^\tau + \tilde{V}_{h+1}(s_{h+1}^\tau) \right) / \tilde{\sigma}_h^2(s_h^\tau, a_h^\tau)$, the ℓ_2 sensi-
sitivity is $2H$. Therefore according to Lemma 2.6, the use of Gaussian Mechanism (the additional
noises ϕ_1, ϕ_2, ϕ_3) ensures ρ_0 -zCDP for each counter. For $\sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \phi(\bar{s}_h^\tau, \bar{a}_h^\tau)^\top + \lambda I$ and
 $\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top / \tilde{\sigma}_h^2(s_h^\tau, a_h^\tau) + \lambda I$, according to Appendix D in [Redberg and Wang,
2021], the per-instance ℓ_2 sensitivity is

$$\|\Delta_x\|_2 = \frac{1}{\sqrt{2}} \sup_{\phi: \|\phi\|_2 \leq 1} \|\phi \phi^\top\|_F = \frac{1}{\sqrt{2}} \sup_{\phi: \|\phi\|_2 \leq 1} \sqrt{\sum_{i,j} \phi_i^2 \phi_j^2} = \frac{1}{\sqrt{2}}.$$

705 Therefore the use of Gaussian Mechanism (the additional noises K_1, K_2) also ensures ρ_0 -zCDP for
 706 each counter.¹² Combining these results, according to Lemma F.17, the whole algorithm satisfies
 707 $5H\rho_0 = \rho$ -zCDP. \square

708 E.3 Proof of the sub-optimality bound

709 E.3.1 Utility analysis and some preparation

710 We begin with the following high probability bound of the noises we add.

711 **Lemma E.2** (Utility analysis). *Let $L = 2H\sqrt{\frac{d}{\rho_0} \log(\frac{10Hd}{\delta})} = 2H\sqrt{\frac{5Hd \log(\frac{10Hd}{\delta})}{\rho}}$ and*
 712 $E = \sqrt{\frac{2d}{\rho_0}} \left(2 + \left(\frac{\log(5c_1H/\delta)}{c_2d} \right)^{\frac{2}{3}} \right) = \sqrt{\frac{10Hd}{\rho}} \left(2 + \left(\frac{\log(5c_1H/\delta)}{c_2d} \right)^{\frac{2}{3}} \right)$ *for some universal constants*
 713 c_1, c_2 . *Then with probability $1 - \delta$, the following inequalities hold simultaneously:*

$$\begin{aligned} & \text{For all } h \in [H], \|\phi_1\|_2 \leq HL, \|\phi_2\|_2 \leq L, \|\phi_3\|_2 \leq L. \\ & \text{For all } h \in [H], K_1, K_2 \text{ are symmetric and positive definite and } \|K_i\|_2 \leq E, i \in \{1, 2\}. \end{aligned} \quad (38)$$

714 *Proof of Lemma E.2.* The second line of (38) results from Lemma 19 in [Redberg and Wang, 2021]
 715 and Weyl's Inequality. The first line of (38) directly results from the concentration inequality for
 716 Gaussian distribution and a union bound. \square

717 Define the Bellman update error $\zeta_h(s, a) := (\mathcal{T}_h \tilde{V}_{h+1})(s, a) - \hat{Q}_h(s, a)$ and recall
 718 $\hat{\pi}_h(s) = \operatorname{argmax}_{\pi_h} \langle \hat{Q}_h(s, \cdot), \pi_h(\cdot | s) \rangle_{\mathcal{A}}$, then because of Lemma F.8,

$$V_1^\pi(s) - V_1^{\hat{\pi}}(s) \leq \sum_{h=1}^H \mathbb{E}_\pi [\zeta_h(s_h, a_h) | s_1 = s] - \sum_{h=1}^H \mathbb{E}_{\hat{\pi}} [\zeta_h(s_h, a_h) | s_1 = s]. \quad (39)$$

719 Define $\tilde{\mathcal{T}}_h \tilde{V}_{h+1}(\cdot, \cdot) = \phi(\cdot, \cdot)^\top \tilde{w}_h$. Then similar to Lemma D.10, we have the following lemma
 720 showing that in order to bound the sub-optimality, it is sufficient to bound the pessimistic penalty.

721 **Lemma E.3** (Lemma C.1 in [Yin et al., 2022]). *Suppose with probability $1 - \delta$, it holds for all*
 722 $s, a, h \in \mathcal{S} \times \mathcal{A} \times [H]$ *that $|(\mathcal{T}_h \tilde{V}_{h+1} - \tilde{\mathcal{T}}_h \tilde{V}_{h+1})(s, a)| \leq \Gamma_h(s, a)$, then it implies $\forall s, a, h \in$*
 723 $\mathcal{S} \times \mathcal{A} \times [H]$, $0 \leq \zeta_h(s, a) \leq 2\Gamma_h(s, a)$. *Furthermore, with probability $1 - \delta$, it holds for any policy*
 724 π *simultaneously,*

$$V_1^\pi(s) - V_1^{\hat{\pi}}(s) \leq \sum_{h=1}^H 2 \cdot \mathbb{E}_\pi [\Gamma_h(s_h, a_h) | s_1 = s].$$

725 *Proof of Lemma E.3.* We first show given $|(\mathcal{T}_h \tilde{V}_{h+1} - \tilde{\mathcal{T}}_h \tilde{V}_{h+1})(s, a)| \leq \Gamma_h(s, a)$, then $0 \leq$
 726 $\zeta_h(s, a) \leq 2\Gamma_h(s, a), \forall s, a, h \in \mathcal{S} \times \mathcal{A} \times [H]$.

727 **Step 1:** The first step is to show $0 \leq \zeta_h(s, a), \forall s, a, h \in \mathcal{S} \times \mathcal{A} \times [H]$.

728 Indeed, if $\bar{Q}_h(s, a) \leq 0$, then by definition $\hat{Q}_h(s, a) = 0$ and therefore $\zeta_h(s, a) := (\mathcal{T}_h \tilde{V}_{h+1})(s, a) -$
 729 $\hat{Q}_h(s, a) = (\mathcal{T}_h \tilde{V}_{h+1})(s, a) \geq 0$. If $\bar{Q}_h(s, a) > 0$, then $\hat{Q}_h(s, a) \leq \bar{Q}_h(s, a)$ and

$$\begin{aligned} \zeta_h(s, a) & := (\mathcal{T}_h \tilde{V}_{h+1})(s, a) - \hat{Q}_h(s, a) \geq (\mathcal{T}_h \tilde{V}_{h+1})(s, a) - \bar{Q}_h(s, a) \\ & = (\mathcal{T}_h \tilde{V}_{h+1})(s, a) - (\tilde{\mathcal{T}}_h \tilde{V}_{h+1})(s, a) + \Gamma_h(s, a) \geq 0. \end{aligned}$$

730 **Step 2:** The second step is to show $\zeta_h(s, a) \leq 2\Gamma_h(s, a), \forall s, a, h \in \mathcal{S} \times \mathcal{A} \times [H]$.

Under the assumption that $|(\mathcal{T}_h \tilde{V}_{h+1} - \tilde{\mathcal{T}}_h \tilde{V}_{h+1})(s, a)| \leq \Gamma_h(s, a)$, we have

$$\bar{Q}_h(s, a) = (\tilde{\mathcal{T}}_h \tilde{V}_{h+1})(s, a) - \Gamma_h(s, a) \leq (\mathcal{T}_h \tilde{V}_{h+1})(s, a) \leq H - h + 1,$$

¹²For more detailed explanation, we refer the readers to Appendix D of [Redberg and Wang, 2021].

731 which implies that $\widehat{Q}_h(s, a) = \max(\bar{Q}_h(s, a), 0)$. Therefore, it holds that

$$\begin{aligned}\zeta_h(s, a) &:= (\mathcal{T}_h \widetilde{V}_{h+1})(s, a) - \widehat{Q}_h(s, a) \leq (\mathcal{T}_h \widetilde{V}_{h+1})(s, a) - \bar{Q}_h(s, a) \\ &= (\mathcal{T}_h \widetilde{V}_{h+1})(s, a) - (\widetilde{\mathcal{T}}_h \widetilde{V}_{h+1})(s, a) + \Gamma_h(s, a) \leq 2 \cdot \Gamma_h(s, a).\end{aligned}$$

732 For the last statement, denote $\mathfrak{F} := \{0 \leq \zeta_h(s, a) \leq 2\Gamma_h(s, a), \forall s, a, h \in \mathcal{S} \times \mathcal{A} \times [H]\}$. Note
733 conditional on \mathfrak{F} , then by (39), $V_1^\pi(s) - V_1^{\widehat{\pi}}(s) \leq \sum_{h=1}^H 2 \cdot \mathbb{E}_\pi[\Gamma_h(s_h, a_h) \mid s_1 = s]$ holds for any
734 policy π almost surely. Therefore,

$$\begin{aligned}& \mathbb{P} \left[\forall \pi, V_1^\pi(s) - V_1^{\widehat{\pi}}(s) \leq \sum_{h=1}^H 2 \cdot \mathbb{E}_\pi[\Gamma_h(s_h, a_h) \mid s_1 = s]. \right] \\ &= \mathbb{P} \left[\forall \pi, V_1^\pi(s) - V_1^{\widehat{\pi}}(s) \leq \sum_{h=1}^H 2 \cdot \mathbb{E}_\pi[\Gamma_h(s_h, a_h) \mid s_1 = s] \mid \mathfrak{F} \right] \cdot \mathbb{P}[\mathfrak{F}] \\ &+ \mathbb{P} \left[\forall \pi, V_1^\pi(s) - V_1^{\widehat{\pi}}(s) \leq \sum_{h=1}^H 2 \cdot \mathbb{E}_\pi[\Gamma_h(s_h, a_h) \mid s_1 = s] \mid \mathfrak{F}^c \right] \cdot \mathbb{P}[\mathfrak{F}^c] \\ &\geq \mathbb{P} \left[\forall \pi, V_1^\pi(s) - V_1^{\widehat{\pi}}(s) \leq \sum_{h=1}^H 2 \cdot \mathbb{E}_\pi[\Gamma_h(s_h, a_h) \mid s_1 = s] \mid \mathfrak{F} \right] \cdot \mathbb{P}[\mathfrak{F}] = 1 \cdot \mathbb{P}[\mathfrak{F}] \geq 1 - \delta,\end{aligned}$$

735 which finishes the proof. \square

736 E.3.2 Bound the pessimistic penalty

737 By Lemma E.3, it remains to bound $|(\mathcal{T}_h \widetilde{V}_{h+1})(s, a) - (\widetilde{\mathcal{T}}_h \widetilde{V}_{h+1})(s, a)|$. Suppose w_h is the coefficient
738 corresponding to the $\mathcal{T}_h \widetilde{V}_{h+1}$ (such w_h exists by Lemma F.14), i.e. $\mathcal{T}_h \widetilde{V}_{h+1} = \phi^\top w_h$, and recall
739 $(\widetilde{\mathcal{T}}_h \widetilde{V}_{h+1})(s, a) = \phi(s, a)^\top \widetilde{w}_h$, then:

$$\begin{aligned}& \left(\mathcal{T}_h \widetilde{V}_{h+1} \right) (s, a) - \left(\widetilde{\mathcal{T}}_h \widetilde{V}_{h+1} \right) (s, a) = \phi(s, a)^\top (w_h - \widetilde{w}_h) \\ &= \phi(s, a)^\top w_h - \phi(s, a)^\top \widehat{\Lambda}_h^{-1} \left(\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \cdot (r_h^\tau + \widetilde{V}_{h+1}(s_{h+1}^\tau)) / \widetilde{\sigma}_h^2(s_h^\tau, a_h^\tau) + \phi_3 \right) \\ &= \underbrace{\phi(s, a)^\top w_h - \phi(s, a)^\top \widehat{\Lambda}_h^{-1} \left(\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \cdot (r_h^\tau + \widetilde{V}_{h+1}(s_{h+1}^\tau)) / \widetilde{\sigma}_h^2(s_h^\tau, a_h^\tau) \right)}_{(i)} \\ &\quad - \underbrace{\phi(s, a)^\top \widehat{\Lambda}_h^{-1} \phi_3 + \phi(s, a)^\top (\widehat{\Lambda}_h^{-1} - \widetilde{\Lambda}_h^{-1}) \left(\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \cdot (r_h^\tau + \widetilde{V}_{h+1}(s_{h+1}^\tau)) / \widetilde{\sigma}_h^2(s_h^\tau, a_h^\tau) + \phi_3 \right)}_{(ii)},\end{aligned}\tag{40}$$

740 where $\widehat{\Lambda}_h = \widetilde{\Lambda}_h - K_2 = \sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top / \widetilde{\sigma}_h^2(s_h^\tau, a_h^\tau) + \lambda I$.

741 Term (ii) can be handled by the following Lemma E.4

742 **Lemma E.4.** Recall κ in Assumption 2.2. Under the high probability event in Lemma E.2, suppose

743 $K \geq \max \left\{ \frac{512H^4 \cdot \log(\frac{2Hd}{\delta})}{\kappa^2}, \frac{4\lambda H^2}{\kappa} \right\}$, then with probability $1 - \delta$, for all $s, a, h \in \mathcal{S} \times \mathcal{A} \times [H]$, it

744 holds that

$$\left| \phi(s, a)^\top \widehat{\Lambda}_h^{-1} \phi_3 \right| \leq \frac{4H^2 L / \kappa}{K}.$$

745 *Proof of Lemma E.4.* Define $\widetilde{\Lambda}_h^p = \mathbb{E}_{\mu, h}[\widetilde{\sigma}_h^{-2}(s, a) \phi(s, a) \phi(s, a)^\top]$. Then because of Assumption
746 2.2 and $\widetilde{\sigma}_h \leq H$, it holds that $\lambda_{\min}(\widetilde{\Lambda}_h^p) \geq \frac{\kappa}{H^2}$. Therefore, due to Lemma F.13, we have with

747 probability $1 - \delta$,

$$\begin{aligned}
& \left| \phi(s, a)^\top \widehat{\Lambda}_h^{-1} \phi_3 \right| \leq \|\phi(s, a)\|_{\widehat{\Lambda}_h^{-1}} \cdot \|\phi_3\|_{\widehat{\Lambda}_h^{-1}} \\
& \leq \frac{4}{K} \|\phi(s, a)\|_{(\widetilde{\Lambda}_h^p)^{-1}} \cdot \|\phi_3\|_{(\widetilde{\Lambda}_h^p)^{-1}} \\
& \leq \frac{4L}{K} \|(\widetilde{\Lambda}_h^p)^{-1}\| \\
& \leq \frac{4H^2L/\kappa}{K}.
\end{aligned}$$

748 The first inequality is because of Cauchy-Schwarz inequality. The second inequality holds with
749 probability $1 - \delta$ due to Lemma F.13 and a union bound. The third inequality holds because
750 $\sqrt{a^\top \cdot A \cdot a} \leq \sqrt{\|a\|_2 \|A\|_2 \|a\|_2} = \|a\|_2 \sqrt{\|A\|_2}$. The last inequality arises from $\|(\widehat{\Lambda}_h^p)^{-1}\| =$
751 $\lambda_{\max}((\widetilde{\Lambda}_h^p)^{-1}) = \lambda_{\min}^{-1}(\widetilde{\Lambda}_h^p) \leq \frac{H^2}{\kappa}$. \square

752 The difference between $\widetilde{\Lambda}_h^{-1}$ and $\widehat{\Lambda}_h^{-1}$ can be bounded by the following Lemma E.5

753 **Lemma E.5.** Under the high probability event in Lemma E.2, suppose $K \geq \frac{128H^4 \log \frac{2dH}{\delta}}{\kappa^2}$, then with
754 probability $1 - \delta$, for all $h \in [H]$, it holds that $\|\widehat{\Lambda}_h^{-1} - \widetilde{\Lambda}_h^{-1}\| \leq \frac{4H^4E/\kappa^2}{K^2}$.

755 *Proof of Lemma E.5.* First of all, we have

$$\begin{aligned}
& \|\widehat{\Lambda}_h^{-1} - \widetilde{\Lambda}_h^{-1}\| = \|\widehat{\Lambda}_h^{-1} \cdot (\widehat{\Lambda}_h - \widetilde{\Lambda}_h) \cdot \widetilde{\Lambda}_h^{-1}\| \\
& \leq \|\widehat{\Lambda}_h^{-1}\| \cdot \|\widehat{\Lambda}_h - \widetilde{\Lambda}_h\| \cdot \|\widetilde{\Lambda}_h^{-1}\| \\
& \leq \lambda_{\min}^{-1}(\widehat{\Lambda}_h) \cdot \lambda_{\min}^{-1}(\widetilde{\Lambda}_h) \cdot E.
\end{aligned} \tag{41}$$

756 The first inequality is because $\|A \cdot B\| \leq \|A\| \cdot \|B\|$. The second inequality is due to Lemma E.2.

Let $\widehat{\Lambda}'_h = \frac{1}{K} \widehat{\Lambda}_h$, then because of Lemma F.12, with probability $1 - \delta$, it holds that for all $h \in [H]$,

$$\left\| \widehat{\Lambda}'_h - \mathbb{E}_{\mu, h}[\phi(s, a)\phi(s, a)^\top / \widetilde{\sigma}_h^2(s, a)] - \frac{\lambda}{K} I_d \right\| \leq \frac{4\sqrt{2}}{\sqrt{K}} \left(\log \frac{2dH}{\delta} \right)^{1/2},$$

which implies that when $K \geq \frac{128H^4 \log \frac{2dH}{\delta}}{\kappa^2}$, it holds that (according to Weyl's Inequality)

$$\lambda_{\min}(\widehat{\Lambda}'_h) \geq \lambda_{\min}(\mathbb{E}_{\mu, h}[\phi(s, a)\phi(s, a)^\top / \widetilde{\sigma}_h^2(s, a)]) + \frac{\lambda}{K} - \frac{\kappa}{2H^2} \geq \frac{\kappa}{2H^2}.$$

Under this high probability event, we have $\lambda_{\min}(\widehat{\Lambda}_h) \geq \frac{K\kappa}{2H^2}$ and therefore $\lambda_{\min}(\widetilde{\Lambda}_h) \geq \lambda_{\min}(\widehat{\Lambda}_h) \geq \frac{K\kappa}{2H^2}$. Plugging these two results into (41), we have

$$\|\widehat{\Lambda}_h^{-1} - \widetilde{\Lambda}_h^{-1}\| \leq \frac{4H^4E/\kappa^2}{K^2}.$$

757 \square

758 Then we can bound term (iii) by the following Lemma E.6

759 **Lemma E.6.** Suppose $K \geq \max\left\{\frac{128H^4 \log \frac{2dH}{\delta}}{\kappa^2}, \frac{\sqrt{2}L}{\sqrt{d\kappa}}\right\}$, under the high probability events in
760 Lemma E.2 and Lemma E.5, it holds that for all $s, a, h \in \mathcal{S} \times \mathcal{A} \times [H]$,

$$\left| \phi(s, a)^\top (\widehat{\Lambda}_h^{-1} - \widetilde{\Lambda}_h^{-1}) \left(\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \cdot (r_h^\tau + \widetilde{V}_{h+1}(s_{h+1}^\tau)) / \widetilde{\sigma}_h^2(s_h^\tau, a_h^\tau) + \phi_3 \right) \right| \leq \frac{4\sqrt{2}H^4E\sqrt{d}/\kappa^{3/2}}{K}.$$

Proof of Lemma E.6. First of all, the left hand side is bounded by

$$\left\| (\widehat{\Lambda}_h^{-1} - \widetilde{\Lambda}_h^{-1}) \left(\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \cdot (r_h^\tau + \widetilde{V}_{h+1}(s_{h+1}^\tau)) / \widetilde{\sigma}_h^2(s_h^\tau, a_h^\tau) \right) \right\|_2 + \frac{4H^4EL/\kappa^2}{K^2}$$

761 due to Lemma E.5. Then the left hand side can be further bounded by

$$\begin{aligned}
& H \sum_{\tau=1}^K \left\| \left(\widehat{\Lambda}_h^{-1} - \widetilde{\Lambda}_h^{-1} \right) \phi(s_h^\tau, a_h^\tau) / \widetilde{\sigma}_h(s_h^\tau, a_h^\tau) \right\|_2 + \frac{4H^4 EL / \kappa^2}{K^2} \\
& \leq H \sum_{\tau=1}^K \sqrt{\text{Tr} \left(\left(\widehat{\Lambda}_h^{-1} - \widetilde{\Lambda}_h^{-1} \right) \cdot \frac{\phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top}{\widetilde{\sigma}_h^2(s_h^\tau, a_h^\tau)} \cdot \left(\widehat{\Lambda}_h^{-1} - \widetilde{\Lambda}_h^{-1} \right) \right)} + \frac{4H^4 EL / \kappa^2}{K^2} \\
& \leq H \sqrt{K \cdot \text{Tr} \left(\left(\widehat{\Lambda}_h^{-1} - \widetilde{\Lambda}_h^{-1} \right) \cdot \widehat{\Lambda}_h \cdot \left(\widehat{\Lambda}_h^{-1} - \widetilde{\Lambda}_h^{-1} \right) \right)} + \frac{4H^4 EL / \kappa^2}{K^2} \\
& \leq H \sqrt{Kd \cdot \lambda_{\max} \left(\left(\widehat{\Lambda}_h^{-1} - \widetilde{\Lambda}_h^{-1} \right) \cdot \widehat{\Lambda}_h \cdot \left(\widehat{\Lambda}_h^{-1} - \widetilde{\Lambda}_h^{-1} \right) \right)} + \frac{4H^4 EL / \kappa^2}{K^2} \\
& = H \sqrt{Kd \cdot \left\| \left(\widehat{\Lambda}_h^{-1} - \widetilde{\Lambda}_h^{-1} \right) \cdot \widehat{\Lambda}_h \cdot \left(\widehat{\Lambda}_h^{-1} - \widetilde{\Lambda}_h^{-1} \right) \right\|_2} + \frac{4H^4 EL / \kappa^2}{K^2} \\
& \leq H \sqrt{Kd \cdot \left\| \widetilde{\Lambda}_h^{-1} \right\|_2 \cdot \left\| \widehat{\Lambda}_h - \widetilde{\Lambda}_h \right\|_2 \cdot \left\| \widehat{\Lambda}_h^{-1} - \widetilde{\Lambda}_h^{-1} \right\|_2} + \frac{4H^4 EL / \kappa^2}{K^2} \\
& \leq \frac{2\sqrt{2}H^4 E \sqrt{d} / \kappa^{3/2}}{K} + \frac{4H^4 EL / \kappa^2}{K^2} \\
& \leq \frac{4\sqrt{2}H^4 E \sqrt{d} / \kappa^{3/2}}{K}.
\end{aligned}$$

762 The first inequality is because $\|a\|_2 = \sqrt{a^\top a} = \sqrt{\text{Tr}(aa^\top)}$. The second inequality is due to
763 Cauchy-Schwarz inequality. The third inequality is because for positive definite matrix A , it holds
764 that $\text{Tr}(A) = \sum_{i=1}^d \lambda_i(A) \leq d \lambda_{\max}(A)$. The equation is because for symmetric, positive definite
765 matrix A , $\|A\|_2 = \lambda_{\max}(A)$. The fourth inequality is due to $\|A \cdot B\| \leq \|A\| \cdot \|B\|$. The fifth
766 inequality is because of Lemma E.2, Lemma E.5 and the statement in the proof of Lemma E.5 that
767 $\lambda_{\min}(\widetilde{\Lambda}_h) \geq \frac{K\kappa}{2H^2}$. The last inequality uses the assumption that $K \geq \frac{\sqrt{2}L}{\sqrt{d}\kappa}$. \square

768 Now the remaining part is term (i), we have

$$\begin{aligned}
& \underbrace{\phi(s, a)^\top w_h - \phi(s, a)^\top \widehat{\Lambda}_h^{-1} \left(\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \cdot \left(r_h^\tau + \widetilde{V}_{h+1}(s_{h+1}^\tau) \right) / \widetilde{\sigma}_h^2(s_h^\tau, a_h^\tau) \right)}_{(i)} \\
& = \underbrace{\phi(s, a)^\top w_h - \phi(s, a)^\top \widehat{\Lambda}_h^{-1} \left(\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \cdot \left(\mathcal{T}_h \widetilde{V}_{h+1} \right)(s_h^\tau, a_h^\tau) / \widetilde{\sigma}_h^2(s_h^\tau, a_h^\tau) \right)}_{(iv)} \\
& \quad - \underbrace{\phi(s, a)^\top \widehat{\Lambda}_h^{-1} \left(\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \cdot \left(r_h^\tau + \widetilde{V}_{h+1}(s_{h+1}^\tau) - \left(\mathcal{T}_h \widetilde{V}_{h+1} \right)(s_h^\tau, a_h^\tau) \right) / \widetilde{\sigma}_h^2(s_h^\tau, a_h^\tau) \right)}_{(v)}.
\end{aligned} \tag{42}$$

769 We are able to bound term (iv) by the following Lemma E.7.

770 **Lemma E.7.** Recall κ in Assumption 2.2. Under the high probability event in Lemma E.2, suppose
771 $K \geq \max \left\{ \frac{512H^4 \cdot \log(\frac{2Hd}{\delta})}{\kappa^2}, \frac{4\lambda H^2}{\kappa} \right\}$, then with probability $1 - \delta$, for all $s, a, h \in \mathcal{S} \times \mathcal{A} \times [H]$,

$$\left| \phi(s, a)^\top w_h - \phi(s, a)^\top \widehat{\Lambda}_h^{-1} \left(\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \cdot \left(\mathcal{T}_h \widetilde{V}_{h+1} \right)(s_h^\tau, a_h^\tau) / \widetilde{\sigma}_h^2(s_h^\tau, a_h^\tau) \right) \right| \leq \frac{8\lambda H^3 \sqrt{d} / \kappa}{K}.$$

772 *Proof of Lemma E.7.* Recall $\mathcal{T}_h \tilde{V}_{h+1} = \phi^\top w_h$ and apply Lemma F.13, we obtain with probability
 773 $1 - \delta$, for all $s, a, h \in \mathcal{S} \times \mathcal{A} \times [H]$,

$$\begin{aligned}
 & \left| \phi(s, a)^\top w_h - \phi(s, a)^\top \widehat{\Lambda}_h^{-1} \left(\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \cdot \left(\mathcal{T}_h \tilde{V}_{h+1} \right) (s_h^\tau, a_h^\tau) / \tilde{\sigma}_h^2(s_h^\tau, a_h^\tau) \right) \right| \\
 &= \left| \phi(s, a)^\top w_h - \phi(s, a)^\top \widehat{\Lambda}_h^{-1} \left(\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \cdot \phi(s_h^\tau, a_h^\tau)^\top w_h / \tilde{\sigma}_h^2(s_h^\tau, a_h^\tau) \right) \right| \\
 &= \left| \phi(s, a)^\top w_h - \phi(s, a)^\top \widehat{\Lambda}_h^{-1} (\widehat{\Lambda}_h - \lambda I) w_h \right| \\
 &= \left| \lambda \cdot \phi(s, a)^\top \widehat{\Lambda}_h^{-1} w_h \right| \\
 &\leq \lambda \|\phi(s, a)\|_{\widehat{\Lambda}_h^{-1}} \cdot \|w_h\|_{\widehat{\Lambda}_h^{-1}} \\
 &\leq \frac{4\lambda}{K} \|\phi(s, a)\|_{(\tilde{\Lambda}_h^p)^{-1}} \cdot \|w_h\|_{(\tilde{\Lambda}_h^p)^{-1}} \\
 &\leq \frac{4\lambda}{K} \cdot 2H\sqrt{d} \cdot \left\| (\tilde{\Lambda}_h^p)^{-1} \right\| \\
 &\leq \frac{8\lambda H^3 \sqrt{d} / \kappa}{K},
 \end{aligned}$$

774 where $\tilde{\Lambda}_h^p := \mathbb{E}_{\mu, h} [\tilde{\sigma}_h(s, a)^{-2} \phi(s, a) \phi(s, a)^\top]$. The first inequality applies Cauchy-Schwarz in-
 775 equality. The second inequality holds with probability $1 - \delta$ due to Lemma F.13 and a union bound.
 776 The third inequality uses $\sqrt{a^\top \cdot A \cdot a} \leq \sqrt{\|a\|_2 \|A\|_2 \|a\|_2} = \|a\|_2 \sqrt{\|A\|_2}$ and $\|w_h\| \leq 2H\sqrt{d}$. Fi-
 777 nally, as $\lambda_{\min}(\tilde{\Lambda}_h^p) \geq \frac{\kappa}{\max_{h, s, a} \tilde{\sigma}_h(s, a)^2} \geq \frac{\kappa}{H^2}$ implies $\left\| (\tilde{\Lambda}_h^p)^{-1} \right\| \leq \frac{H^2}{\kappa}$, the last inequality holds. \square

778 For term (v), denote: $x_\tau = \frac{\phi(s_h^\tau, a_h^\tau)}{\tilde{\sigma}_h(s_h^\tau, a_h^\tau)}$, $\eta_\tau = \left(r_h^\tau + \tilde{V}_{h+1}(s_{h+1}^\tau) - \left(\mathcal{T}_h \tilde{V}_{h+1} \right) (s_h^\tau, a_h^\tau) \right) / \tilde{\sigma}_h(s_h^\tau, a_h^\tau)$,
 779 then by Cauchy-Schwarz inequality, it holds that for all $h, s, a \in [H] \times \mathcal{S} \times \mathcal{A}$,

$$\begin{aligned}
 & \left| \phi(s, a)^\top \widehat{\Lambda}_h^{-1} \left(\sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \cdot \left(r_h^\tau + \tilde{V}_{h+1}(s_{h+1}^\tau) - \left(\mathcal{T}_h \tilde{V}_{h+1} \right) (s_h^\tau, a_h^\tau) \right) / \tilde{\sigma}_h^2(s_h^\tau, a_h^\tau) \right) \right| \\
 &\leq \sqrt{\phi(s, a)^\top \widehat{\Lambda}_h^{-1} \phi(s, a)} \cdot \left\| \sum_{\tau=1}^K x_\tau \eta_\tau \right\|_{\widehat{\Lambda}_h^{-1}}.
 \end{aligned} \tag{43}$$

780 We bound $\sqrt{\phi(s, a)^\top \widehat{\Lambda}_h^{-1} \phi(s, a)}$ by $\sqrt{\phi(s, a)^\top \tilde{\Lambda}_h^{-1} \phi(s, a)}$ using the following Lemma E.8.

Lemma E.8. Suppose $K \geq \max\left\{ \frac{128H^4 \log \frac{2dH}{\delta}}{\kappa^2}, \frac{\sqrt{2L}}{\sqrt{d\kappa}} \right\}$, under the high probability events in Lemma E.2 and Lemma E.5, it holds that for all $s, a, h \in \mathcal{S} \times \mathcal{A} \times [H]$,

$$\sqrt{\phi(s, a)^\top \widehat{\Lambda}_h^{-1} \phi(s, a)} \leq \sqrt{\phi(s, a)^\top \tilde{\Lambda}_h^{-1} \phi(s, a)} + \frac{2H^2 \sqrt{E} / \kappa}{K}.$$

Proof of Lemma E.8.

$$\begin{aligned}
 & \sqrt{\phi(s, a)^\top \widehat{\Lambda}_h^{-1} \phi(s, a)} = \sqrt{\phi(s, a)^\top \tilde{\Lambda}_h^{-1} \phi(s, a) + \phi(s, a)^\top (\widehat{\Lambda}_h^{-1} - \tilde{\Lambda}_h^{-1}) \phi(s, a)} \\
 &\leq \sqrt{\phi(s, a)^\top \tilde{\Lambda}_h^{-1} \phi(s, a)} + \left\| \widehat{\Lambda}_h^{-1} - \tilde{\Lambda}_h^{-1} \right\|_2 \\
 &\leq \sqrt{\phi(s, a)^\top \tilde{\Lambda}_h^{-1} \phi(s, a)} + \sqrt{\left\| \widehat{\Lambda}_h^{-1} - \tilde{\Lambda}_h^{-1} \right\|_2} \\
 &\leq \sqrt{\phi(s, a)^\top \tilde{\Lambda}_h^{-1} \phi(s, a)} + \frac{2H^2 \sqrt{E} / \kappa}{K}.
 \end{aligned} \tag{44}$$

781 The first inequality uses $|a^\top Aa| \leq \|a\|_2^2 \cdot \|A\|$. The second inequality is because for $a, b \geq 0$,
 782 $\sqrt{a} + \sqrt{b} \geq \sqrt{a+b}$. The last inequality uses Lemma E.5. \square

Remark E.9. Similarly, under the same assumption in Lemma E.8, we also have for all $s, a, h \in \mathcal{S} \times \mathcal{A} \times [H]$,

$$\sqrt{\phi(s, a)^\top \tilde{\Lambda}_h^{-1} \phi(s, a)} \leq \sqrt{\phi(s, a)^\top \hat{\Lambda}_h^{-1} \phi(s, a)} + \frac{2H^2 \sqrt{E}/\kappa}{K}.$$

783 E.3.3 An intermediate result: bounding the variance

784 Before we handle $\left\| \sum_{\tau=1}^K x_\tau \eta_\tau \right\|_{\tilde{\Lambda}_h^{-1}}$, we first bound $\sup_h \left\| \tilde{\sigma}_h^2 - \sigma_{\tilde{V}_{h+1}}^2 \right\|_\infty$ by the following
 785 Lemma E.10.

786 **Lemma E.10** (Private version of Lemma C.7 in [Yin et al., 2022]). Recall the definition of $\tilde{\sigma}_h(\cdot, \cdot)^2 =$
 787 $\max\{1, \widetilde{\text{Var}}_h \tilde{V}_{h+1}(\cdot, \cdot)\}$ in Algorithm 2 where $[\widetilde{\text{Var}}_h \tilde{V}_{h+1}](\cdot, \cdot) = \langle \phi(\cdot, \cdot), \tilde{\beta}_h \rangle_{[0, (H-h+1)^2]} -$
 788 $[\langle \phi(\cdot, \cdot), \tilde{\theta}_h \rangle_{[0, H-h+1]}]^2$ ($\tilde{\beta}_h$ and $\tilde{\theta}_h$ are defined in Algorithm 2) and $\sigma_{\tilde{V}_{h+1}}(\cdot, \cdot)^2 :=$
 789 $\max\{1, \text{Var}_{P_h} \tilde{V}_{h+1}(\cdot, \cdot)\}$. Suppose $K \geq \max\left\{ \frac{512 \log(\frac{2Hd}{\delta})}{\kappa^2}, \frac{4\lambda}{\kappa}, \frac{128 \log \frac{2dH}{\delta}}{\kappa^2}, \frac{\sqrt{2L}}{H\sqrt{d\kappa}} \right\}$ and $K \geq$
 790 $\max\left\{ \frac{4L^2}{H^2 d^3 \kappa}, \frac{32E^2}{d^2 \kappa^2}, \frac{16\lambda^2}{d^2 \kappa} \right\}$, under the high probability event in Lemma E.2, it holds that with proba-
 791 bility $1 - 6\delta$,

$$\sup_h \left\| \tilde{\sigma}_h^2 - \sigma_{\tilde{V}_{h+1}}^2 \right\|_\infty \leq 36 \sqrt{\frac{H^4 d^3}{\kappa K} \log \left(\frac{(\lambda + K) 2K d H^2}{\lambda \delta} \right)}.$$

792 *Proof of Lemma E.10. Step 1:* The first step is to show for all $h, s, a \in [H] \times \mathcal{S} \times \mathcal{A}$, with probability
 793 $1 - 3\delta$,

$$\left| \langle \phi(s, a), \tilde{\beta}_h \rangle_{[0, (H-h+1)^2]} - \mathbb{P}_h(\tilde{V}_{h+1})^2(s, a) \right| \leq 12 \sqrt{\frac{H^4 d^3}{\kappa K} \log \left(\frac{(\lambda + K) 2K d H^2}{\lambda \delta} \right)}.$$

794 **Proof of Step 1.** We can bound the left hand side by the following decomposition:

$$\begin{aligned} & \left| \langle \phi(s, a), \tilde{\beta}_h \rangle_{[0, (H-h+1)^2]} - \mathbb{P}_h(\tilde{V}_{h+1})^2(s, a) \right| \leq \left| \langle \phi(s, a), \tilde{\beta}_h \rangle - \mathbb{P}_h(\tilde{V}_{h+1})^2(s, a) \right| \\ &= \left| \phi(s, a)^\top \tilde{\Sigma}_h^{-1} \left(\sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \cdot \tilde{V}_{h+1}(\bar{s}_{h+1}^\tau)^2 + \phi_1 \right) - \mathbb{P}_h(\tilde{V}_{h+1})^2(s, a) \right| \\ &\leq \underbrace{\left| \phi(s, a)^\top \tilde{\Sigma}_h^{-1} \left(\sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \cdot \tilde{V}_{h+1}(\bar{s}_{h+1}^\tau)^2 \right) - \mathbb{P}_h(\tilde{V}_{h+1})^2(s, a) \right|}_{(1)} + \underbrace{\left| \phi(s, a)^\top \tilde{\Sigma}_h^{-1} \phi_1 \right|}_{(2)} \\ &\quad + \underbrace{\left| \phi(s, a)^\top (\tilde{\Sigma}_h^{-1} - \bar{\Sigma}_h^{-1}) \left(\sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \cdot \tilde{V}_{h+1}(\bar{s}_{h+1}^\tau)^2 + \phi_1 \right) \right|}_{(3)}, \end{aligned}$$

795 where $\bar{\Sigma}_h = \tilde{\Sigma}_h - K_1 = \sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \phi(\bar{s}_h^\tau, \bar{a}_h^\tau)^\top + \lambda I$.

Similar to the proof in Lemma E.5, when $K \geq \max\left\{ \frac{128 \log \frac{2dH}{\delta}}{\kappa^2}, \frac{\sqrt{2L}}{H\sqrt{d\kappa}} \right\}$, it holds that with probability
 $1 - \delta$, for all $h \in [H]$,

$$\lambda_{\min}(\bar{\Sigma}_h) \geq \frac{K\kappa}{2}, \quad \lambda_{\min}(\tilde{\Sigma}_h) \geq \frac{K\kappa}{2}, \quad \left\| \tilde{\Sigma}_h^{-1} - \bar{\Sigma}_h^{-1} \right\|_2 \leq \frac{4E/\kappa^2}{K^2}.$$

796 (The only difference to Lemma E.5 is here $\mathbb{E}_{\mu, h}[\phi(s, a)\phi(s, a)^\top] \geq \kappa$.)

797 Under this high probability event, for term (2), it holds that for all $h, s, a \in [H] \times \mathcal{S} \times \mathcal{A}$,

$$\left| \phi(s, a)^\top \tilde{\Sigma}_h^{-1} \phi_1 \right| \leq \|\phi(s, a)\| \cdot \|\tilde{\Sigma}_h^{-1}\| \cdot \|\phi_1\| \leq \lambda_{\min}^{-1}(\bar{\Sigma}_h) \cdot HL \leq \frac{2HL/\kappa}{K}. \quad (45)$$

798 For term (3), similar to Lemma E.6, we have for all $h, s, a \in [H] \times \mathcal{S} \times \mathcal{A}$,

$$\left| \phi(s, a)^\top (\tilde{\Sigma}_h^{-1} - \bar{\Sigma}_h^{-1}) \left(\sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \cdot \tilde{V}_{h+1}(\bar{s}_{h+1}^\tau)^2 + \phi_1 \right) \right| \leq \frac{4\sqrt{2}H^2 E \sqrt{d} / \kappa^{3/2}}{K}. \quad (46)$$

799 (The only difference to Lemma E.6 is that here $\tilde{V}_{h+1}(s)^2 \leq H^2$, $\|\phi_1\|_2 \leq HL$, $\|\tilde{\Sigma}_h^{-1}\|_2 \leq \frac{2}{K\kappa}$ and
800 $\|\tilde{\Sigma}_h^{-1} - \bar{\Sigma}_h^{-1}\|_2 \leq \frac{4E/\kappa^2}{K^2}$.)

801 We further decompose term (1) as below.

$$\begin{aligned} (1) &= \left| \phi(s, a)^\top \bar{\Sigma}_h^{-1} \left(\sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \cdot \tilde{V}_{h+1}(\bar{s}_{h+1}^\tau)^2 \right) - \mathbb{P}_h(\tilde{V}_{h+1})^2(s, a) \right| \\ &= \left| \phi(s, a)^\top \bar{\Sigma}_h^{-1} \sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \cdot \tilde{V}_{h+1}(\bar{s}_{h+1}^\tau)^2 - \phi(s, a)^\top \bar{\Sigma}_h^{-1} \left(\sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \phi(\bar{s}_h^\tau, \bar{a}_h^\tau)^\top + \lambda I \right) \int_{\mathcal{S}} (\tilde{V}_{h+1})^2(s') d\nu_h(s') \right| \\ &\leq \underbrace{\left| \phi(s, a)^\top \bar{\Sigma}_h^{-1} \sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \cdot \left(\tilde{V}_{h+1}(\bar{s}_{h+1}^\tau)^2 - \mathbb{P}_h(\tilde{V}_{h+1})^2(\bar{s}_h^\tau, \bar{a}_h^\tau) \right) \right|}_{(4)} + \underbrace{\lambda \left| \phi(s, a)^\top \bar{\Sigma}_h^{-1} \int_{\mathcal{S}} (\tilde{V}_{h+1})^2(s') d\nu_h(s') \right|}_{(5)}. \end{aligned} \quad (47)$$

802 For term (5), because $K \geq \max \left\{ \frac{512 \log(\frac{2Hd}{\delta})}{\kappa^2}, \frac{4\lambda}{\kappa} \right\}$, by Lemma F.13 and a union bound, with
803 probability $1 - \delta$, for all $h, s, a \in [H] \times \mathcal{S} \times \mathcal{A}$,

$$\begin{aligned} \lambda \left| \phi(s, a)^\top \bar{\Sigma}_h^{-1} \int_{\mathcal{S}} (\tilde{V}_{h+1})^2(s') d\nu_h(s') \right| &\leq \lambda \|\phi(s, a)\|_{\bar{\Sigma}_h^{-1}} \left\| \int_{\mathcal{S}} (\tilde{V}_{h+1})^2(s') d\nu_h(s') \right\|_{\bar{\Sigma}_h^{-1}} \\ &\leq \lambda \frac{2}{\sqrt{K}} \|\phi(s, a)\|_{(\Sigma_h^p)^{-1}} \frac{2}{\sqrt{K}} \left\| \int_{\mathcal{S}} (\tilde{V}_{h+1})^2(s') d\nu_h(s') \right\|_{(\Sigma_h^p)^{-1}} \leq 4\lambda \|(\Sigma_h^p)^{-1}\| \frac{H^2 \sqrt{d}}{K} \leq 4\lambda \frac{H^2 \sqrt{d}}{\kappa K}, \end{aligned} \quad (48)$$

804 where $\Sigma_h^p = \mathbb{E}_{\mu, h}[\phi(s, a)\phi(s, a)^\top]$ and $\lambda_{\min}(\Sigma_h^p) \geq \kappa$.

805 For term (4), it can be bounded by the following inequality (because of Cauchy-Schwarz inequality).

$$(4) \leq \|\phi(s, a)\|_{\bar{\Sigma}_h^{-1}} \cdot \left\| \sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \cdot \left(\tilde{V}_{h+1}(\bar{s}_{h+1}^\tau)^2 - \mathbb{P}_h(\tilde{V}_{h+1})^2(\bar{s}_h^\tau, \bar{a}_h^\tau) \right) \right\|_{\bar{\Sigma}_h^{-1}}. \quad (49)$$

806 **Bounding using covering.** Note for any fix V_{h+1} , we can define $x_\tau = \phi(\bar{s}_h^\tau, \bar{a}_h^\tau)$ ($\|\phi\|_2 \leq 1$) and
807 $\eta_\tau = V_{h+1}(\bar{s}_{h+1}^\tau)^2 - \mathbb{P}_h(V_{h+1})^2(\bar{s}_h^\tau, \bar{a}_h^\tau)$ is H^2 -subgaussian, by Lemma F.9 (where $t = K$ and
808 $L = 1$), it holds that with probability $1 - \delta$,

$$\left\| \sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \cdot \left(V_{h+1}(\bar{s}_{h+1}^\tau)^2 - \mathbb{P}_h(V_{h+1})^2(\bar{s}_h^\tau, \bar{a}_h^\tau) \right) \right\|_{\bar{\Sigma}_h^{-1}} \leq \sqrt{8H^4 \cdot \frac{d}{2} \log \left(\frac{\lambda + K}{\lambda \delta} \right)}.$$

809 Let $\mathcal{N}_h(\epsilon)$ be the minimal ϵ -cover (with respect to the supremum norm) of

$$810 \mathcal{V}_h := \left\{ V_h : V_h(\cdot) = \max_{a \in \mathcal{A}} \left\{ \min \left\{ \phi(s, a)^\top \theta - C_1 \sqrt{d \cdot \phi(\cdot, \cdot)^\top \tilde{\Lambda}_h^{-1} \phi(\cdot, \cdot)} - C_2, H - h + 1 \right\}^+ \right\} \right\}.$$

811 That is, for any $V \in \mathcal{V}_h$, there exists a value function $V' \in \mathcal{N}_h(\epsilon)$ such that
812 $\sup_{s \in \mathcal{S}} |V(s) - V'(s)| < \epsilon$. Now by a union bound, we obtain with probability $1 - \delta$,

$$\sup_{V_{h+1} \in \mathcal{N}_{h+1}(\epsilon)} \left\| \sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \cdot \left(V_{h+1}(\bar{s}_{h+1}^\tau)^2 - \mathbb{P}_h(V_{h+1})^2(\bar{s}_h^\tau, \bar{a}_h^\tau) \right) \right\|_{\bar{\Sigma}_h^{-1}} \leq \sqrt{8H^4 \cdot \frac{d}{2} \log \left(\frac{\lambda + K}{\lambda \delta} |\mathcal{N}_{h+1}(\epsilon)| \right)}$$

813 which implies

$$\begin{aligned} & \left\| \sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \cdot \left(\tilde{V}_{h+1}(\bar{s}_{h+1}^\tau)^2 - \mathbb{P}_h(\tilde{V}_{h+1})^2(\bar{s}_h^\tau, \bar{a}_h^\tau) \right) \right\|_{\bar{\Sigma}_h^{-1}} \\ & \leq \sqrt{8H^4 \cdot \frac{d}{2} \log \left(\frac{\lambda + K}{\lambda \delta} |\mathcal{N}_{h+1}(\epsilon)| \right)} + 4H^2 \sqrt{\epsilon^2 K^2 / \lambda} \end{aligned}$$

814 choosing $\epsilon = d\sqrt{\lambda}/K$, applying Lemma B.3 of [Jin et al., 2021]¹³ to the covering number $\mathcal{N}_{h+1}(\epsilon)$
815 w.r.t. \mathcal{V}_{h+1} , we can further bound above by

$$\leq \sqrt{8H^4 \cdot \frac{d^3}{2} \log \left(\frac{\lambda + K}{\lambda \delta} 2dHK \right)} + 4H^2 \sqrt{d^2} \leq 6\sqrt{H^4 \cdot d^3 \log \left(\frac{\lambda + K}{\lambda \delta} 2dHK \right)}$$

816 Apply a union bound for $h \in [H]$, we have with probability $1 - \delta$, for all $h \in [H]$,

$$\left\| \sum_{\tau=1}^K \phi(\bar{s}_h^\tau, \bar{a}_h^\tau) \cdot \left(\tilde{V}_{h+1}(\bar{s}_{h+1}^\tau)^2 - \mathbb{P}_h(\tilde{V}_{h+1})^2(\bar{s}_h^\tau, \bar{a}_h^\tau) \right) \right\|_{\bar{\Sigma}_h^{-1}} \leq 6\sqrt{H^4 d^3 \log \left(\frac{(\lambda + K)2KdH^2}{\lambda \delta} \right)} \quad (50)$$

817 and similar to term (2), it holds that for all $h, s, a \in [H] \times \mathcal{S} \times \mathcal{A}$,

$$\|\phi(s, a)\|_{\bar{\Sigma}_h^{-1}} \leq \sqrt{\|\bar{\Sigma}_h^{-1}\|} \leq \sqrt{\frac{2}{\kappa K}}. \quad (51)$$

818 Combining (45), (46), (47), (48), (49), (50), (51) and the assumption that $K \geq$
819 $\max\{\frac{4L^2}{H^2 d^3 \kappa}, \frac{32E^2}{d^2 \kappa^2}, \frac{16\lambda^2}{d^2 \kappa}\}$, we obtain with probability $1 - 3\delta$ for all $h, s, a \in [H] \times \mathcal{S} \times \mathcal{A}$,

$$\left| \langle \phi(s, a), \tilde{\beta}_h \rangle_{[0, (H-h+1)^2]} - \mathbb{P}_h(\tilde{V}_{h+1})^2(s, a) \right| \leq 12\sqrt{\frac{H^4 d^3}{\kappa K} \log \left(\frac{(\lambda + K)2KdH^2}{\lambda \delta} \right)}.$$

820 **Step 2:** The second step is to show for all $h, s, a \in [H] \times \mathcal{S} \times \mathcal{A}$, with probability $1 - 3\delta$,

$$\left| \langle \phi(s, a), \tilde{\theta}_h \rangle_{[0, H-h+1]} - \mathbb{P}_h(\tilde{V}_{h+1})(s, a) \right| \leq 12\sqrt{\frac{H^2 d^3}{\kappa K} \log \left(\frac{(\lambda + K)2KdH^2}{\lambda \delta} \right)}. \quad (52)$$

821 The proof of Step 2 is nearly identical to Step 1 except \tilde{V}_h^2 is replaced by \tilde{V}_h .

822 **Step 3:** The last step is to prove $\sup_{P_h} \|\tilde{\sigma}_h^2 - \sigma_{\tilde{V}_{h+1}}^2\|_\infty \leq 36\sqrt{\frac{H^4 d^3}{\kappa K} \log \left(\frac{(\lambda + K)2KdH^2}{\lambda \delta} \right)}$ with high
823 probability.

824 **Proof of Step 3.** By (52),

$$\begin{aligned} & \left| \left[\langle \phi(\cdot, \cdot), \tilde{\theta}_h \rangle_{[0, H-h+1]} \right]^2 - \left[\mathbb{P}_h(\tilde{V}_{h+1})(s, a) \right]^2 \right| \\ & = \left| \langle \phi(s, a), \tilde{\theta}_h \rangle_{[0, H-h+1]} + \mathbb{P}_h(\tilde{V}_{h+1})(s, a) \right| \cdot \left| \langle \phi(s, a), \tilde{\theta}_h \rangle_{[0, H-h+1]} - \mathbb{P}_h(\tilde{V}_{h+1})(s, a) \right| \\ & \leq 2H \cdot \left| \langle \phi(s, a), \tilde{\theta}_h \rangle_{[0, H-h+1]} - \mathbb{P}_h(\tilde{V}_{h+1})(s, a) \right| \leq 24\sqrt{\frac{H^4 d^3}{\kappa K} \log \left(\frac{(\lambda + K)2KdH^2}{\lambda \delta} \right)}. \end{aligned}$$

825 Combining this with Step 1, we have with probability $1 - 6\delta$, $\forall h, s, a \in [H] \times \mathcal{S} \times \mathcal{A}$,

$$\left| \widetilde{\text{Var}}_h \tilde{V}_{h+1}(s, a) - \text{Var}_{P_h} \tilde{V}_{h+1}(s, a) \right| \leq 36\sqrt{\frac{H^4 d^3}{\kappa K} \log \left(\frac{(\lambda + K)2KdH^2}{\lambda \delta} \right)}.$$

826 Finally, by the non-expansiveness of operator $\max\{1, \cdot\}$, the proof is complete. \square

¹³Note that the conclusion in [Jin et al., 2021] hold here even though we have an extra constant C_2 .

827 **E.3.4 Validity of our pessimism**

828 Recall the definition $\widehat{\Lambda}_h = \sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top / \widetilde{\sigma}_h^2(s_h^\tau, a_h^\tau) + \lambda \cdot I$ and
 829 $\Lambda_h = \sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top / \sigma_{V_{h+1}}^2(s_h^\tau, a_h^\tau) + \lambda I$. Then we have the following lemma to bound
 830 the term $\sqrt{\phi(s, a)^\top \widehat{\Lambda}_h^{-1} \phi(s, a)}$ by $\sqrt{\phi(s, a)^\top \Lambda_h^{-1} \phi(s, a)}$.

831 **Lemma E.11** (Private version of lemma C.3 in [Yin et al., 2022]). *Denote the quantities $C_1 =$
 832 $\max\{2\lambda, 128 \log(2dH/\delta), \frac{128H^4 \log(2dH/\delta)}{\kappa^2}\}$ and $C_2 = \widetilde{O}(H^{12}d^3/\kappa^5)$. Suppose the number of
 833 episode K satisfies $K > \max\{C_1, C_2\}$ and the condition in Lemma E.10, under the high probability
 834 events in Lemma E.2 and Lemma E.10, it holds that with probability $1 - 2\delta$, for all $h, s, a \in$
 835 $[H] \times \mathcal{S} \times \mathcal{A}$,*

$$\sqrt{\phi(s, a)^\top \widehat{\Lambda}_h^{-1} \phi(s, a)} \leq 2\sqrt{\phi(s, a)^\top \Lambda_h^{-1} \phi(s, a)}.$$

836 *Proof of Lemma E.11.* By definition $\sqrt{\phi(s, a)^\top \widehat{\Lambda}_h^{-1} \phi(s, a)} = \|\phi(s, a)\|_{\widehat{\Lambda}_h^{-1}}$. Then denote

$$\widehat{\Lambda}'_h = \frac{1}{K} \widehat{\Lambda}_h, \quad \Lambda'_h = \frac{1}{K} \Lambda_h,$$

837 where $\Lambda_h = \sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top / \sigma_{V_{h+1}}^2(s_h^\tau, a_h^\tau) + \lambda I$. Under the assumption of K , by the
 838 conclusion in Lemma E.10, we have

$$\begin{aligned} \left\| \widehat{\Lambda}'_h - \Lambda'_h \right\| &\leq \sup_{s, a} \left\| \frac{\phi(s, a) \phi(s, a)^\top}{\widetilde{\sigma}_h^2(s, a)} - \frac{\phi(s, a) \phi(s, a)^\top}{\sigma_{V_{h+1}}^2(s, a)} \right\| \\ &\leq \sup_{s, a} \left| \frac{\widetilde{\sigma}_h^2(s, a) - \sigma_{V_{h+1}}^2(s, a)}{\widetilde{\sigma}_h^2(s, a) \cdot \sigma_{V_{h+1}}^2(s, a)} \right| \cdot \|\phi(s, a)\|^2 \\ &\leq \sup_{s, a} \left| \frac{\widetilde{\sigma}_h^2(s, a) - \sigma_{V_{h+1}}^2(s, a)}{1} \right| \cdot 1 \\ &\leq 36 \sqrt{\frac{H^4 d^3}{\kappa K} \log\left(\frac{(\lambda + K) 2K d H^2}{\lambda \delta}\right)}. \end{aligned} \tag{53}$$

839 Next by Lemma F.12 (with ϕ to be $\phi/\sigma_{V_{h+1}}$ and therefore $C = 1$) and a union bound, it holds with
 840 probability $1 - \delta$, for all $h \in [H]$,

$$\left\| \Lambda'_h - \left(\mathbb{E}_{\mu, h}[\phi(s, a) \phi(s, a)^\top / \sigma_{V_{h+1}}^2(s, a)] + \frac{\lambda}{K} I_d \right) \right\| \leq \frac{4\sqrt{2}}{\sqrt{K}} \left(\log \frac{2dH}{\delta} \right)^{1/2}.$$

841 Therefore by Weyl's inequality and the assumption that K satisfies that

842 $K > \max\{2\lambda, 128 \log(2dH/\delta), \frac{128H^4 \log(2dH/\delta)}{\kappa^2}\}$, the above inequality leads to

$$\begin{aligned} \|\Lambda'_h\| &= \lambda_{\max}(\Lambda'_h) \leq \lambda_{\max} \left(\mathbb{E}_{\mu, h}[\phi(s, a) \phi(s, a)^\top / \sigma_{V_{h+1}}^2(s, a)] \right) + \frac{\lambda}{K} + \frac{4\sqrt{2}}{\sqrt{K}} \left(\log \frac{2dH}{\delta} \right)^{1/2} \\ &= \left\| \mathbb{E}_{\mu, h}[\phi(s, a) \phi(s, a)^\top / \sigma_{V_{h+1}}^2(s, a)] \right\|_2 + \frac{\lambda}{K} + \frac{4\sqrt{2}}{\sqrt{K}} \left(\log \frac{2dH}{\delta} \right)^{1/2} \\ &\leq \|\phi(s, a)\|^2 + \frac{\lambda}{K} + \frac{4\sqrt{2}}{\sqrt{K}} \left(\log \frac{2dH}{\delta} \right)^{1/2} \leq 1 + \frac{\lambda}{K} + \frac{4\sqrt{2}}{\sqrt{K}} \left(\log \frac{2dH}{\delta} \right)^{1/2} \leq 2, \\ \lambda_{\min}(\Lambda'_h) &\geq \lambda_{\min} \left(\mathbb{E}_{\mu, h}[\phi(s, a) \phi(s, a)^\top / \sigma_{V_{h+1}}^2(s, a)] \right) + \frac{\lambda}{K} - \frac{4\sqrt{2}}{\sqrt{K}} \left(\log \frac{2dH}{\delta} \right)^{1/2} \\ &\geq \lambda_{\min} \left(\mathbb{E}_{\mu, h}[\phi(s, a) \phi(s, a)^\top / \sigma_{V_{h+1}}^2(s, a)] \right) - \frac{4\sqrt{2}}{\sqrt{K}} \left(\log \frac{2dH}{\delta} \right)^{1/2} \\ &\geq \frac{\kappa}{H^2} - \frac{4\sqrt{2}}{\sqrt{K}} \left(\log \frac{2dH}{\delta} \right)^{1/2} \geq \frac{\kappa}{2H^2}. \end{aligned}$$

843 Hence with probability $1 - \delta$, $\|\Lambda'_h\| \leq 2$ and $\|\Lambda_h'^{-1}\| = \lambda_{\min}^{-1}(\Lambda'_h) \leq \frac{2H^2}{\kappa}$. Similarly, one can show
 844 $\|\widehat{\Lambda}_h'^{-1}\| \leq \frac{2H^2}{\kappa}$ with probability $1 - \delta$ using identical proof.

845 Now apply Lemma F.11 and a union bound to $\widehat{\Lambda}_h'$ and Λ_h' , we obtain with probability $1 - \delta$, for all
 846 $h, s, a \in [H] \times \mathcal{S} \times \mathcal{A}$,

$$\begin{aligned} \|\phi(s, a)\|_{\widehat{\Lambda}_h'^{-1}} &\leq \left[1 + \sqrt{\|\Lambda_h'^{-1}\| \cdot \|\Lambda_h'\| \cdot \|\widehat{\Lambda}_h'^{-1}\| \cdot \|\widehat{\Lambda}_h' - \Lambda_h'\|} \right] \cdot \|\phi(s, a)\|_{\Lambda_h'^{-1}} \\ &\leq \left[1 + \sqrt{\frac{2H^2}{\kappa} \cdot 2 \cdot \frac{2H^2}{\kappa} \cdot \|\widehat{\Lambda}_h' - \Lambda_h'\|} \right] \cdot \|\phi(s, a)\|_{\Lambda_h'^{-1}} \\ &\leq \left[1 + \sqrt{\frac{288H^4}{\kappa^2} \left(\sqrt{\frac{H^4 d^3}{\kappa K} \log \left(\frac{(\lambda + K)2KdH^2}{\lambda\delta} \right)} \right)} \right] \cdot \|\phi(s, a)\|_{\Lambda_h'^{-1}} \\ &\leq 2 \|\phi(s, a)\|_{\Lambda_h'^{-1}} \end{aligned}$$

847 where the third inequality uses (53) and the last inequality uses $K > \widetilde{O}(H^{12}d^3/\kappa^5)$. Note the
 848 conclusion can be derived directly by the above inequality multiplying $1/\sqrt{K}$ on both sides. \square

849 In order to bound $\left\| \sum_{\tau=1}^K x_\tau \eta_\tau \right\|_{\widehat{\Lambda}_h'^{-1}}$, we apply the following Lemma E.12.

850 **Lemma E.12** (Lemma C.4 in [Yin et al., 2022]). Recall $x_\tau = \frac{\phi(s_h^\tau, a_h^\tau)}{\sigma_h(s_h^\tau, a_h^\tau)}$ and

851 $\eta_\tau = \left(r_h^\tau + \widetilde{V}_{h+1}(s_h^\tau, a_h^\tau) - (\mathcal{T}_h \widetilde{V}_{h+1})(s_h^\tau, a_h^\tau) \right) / \widetilde{\sigma}_h(s_h^\tau, a_h^\tau)$. Denote

$$\xi := \sup_{V \in [0, H], s' \sim P_h(s, a), h \in [H]} \left| \frac{r_h + V(s') - (\mathcal{T}_h V)(s, a)}{\sigma_V(s, a)} \right|.$$

852 Suppose $K \geq \widetilde{O}(H^{12}d^3/\kappa^5)$ ¹⁴, then with probability $1 - \delta$,

$$\left\| \sum_{\tau=1}^K x_\tau \eta_\tau \right\|_{\widehat{\Lambda}_h'^{-1}} \leq \widetilde{O} \left(\max \{ \sqrt{d}, \xi \} \right),$$

853 where \widetilde{O} absorbs constants and Polylog terms.

854 Now we are ready to prove the following key lemma, which gives a high probability bound for
 855 $\left| (\mathcal{T}_h \widetilde{V}_{h+1} - \widetilde{\mathcal{T}}_h \widetilde{V}_{h+1})(s, a) \right|$.

856 **Lemma E.13.** Assume $K > \max\{\mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_3, \mathcal{M}_4\}$, for any $0 < \lambda < \kappa$, suppose $\sqrt{d} > \xi$,

857 where $\xi := \sup_{V \in [0, H], s' \sim P_h(s, a), h \in [H]} \left| \frac{r_h + V(s') - (\mathcal{T}_h V)(s, a)}{\sigma_V(s, a)} \right|$. Then with probability $1 - \delta$, for all

858 $h, s, a \in [H] \times \mathcal{S} \times \mathcal{A}$,

$$\left| (\mathcal{T}_h \widetilde{V}_{h+1} - \widetilde{\mathcal{T}}_h \widetilde{V}_{h+1})(s, a) \right| \leq \widetilde{O} \left(\sqrt{d} \sqrt{\phi(s, a)^\top \widetilde{\Lambda}_h^{-1} \phi(s, a)} \right) + \frac{D}{K},$$

where $\widetilde{\Lambda}_h = \sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top / \widetilde{\sigma}_h^2(s_h^\tau, a_h^\tau) + \lambda I + K_2$,

$$D = \widetilde{O} \left(\frac{H^2 L}{\kappa} + \frac{H^4 E \sqrt{d}}{\kappa^{3/2}} + H^3 \sqrt{d} + \frac{H^2 \sqrt{E d}}{\kappa} \right) = \widetilde{O} \left(\frac{H^2 L}{\kappa} + \frac{H^4 E \sqrt{d}}{\kappa^{3/2}} + H^3 \sqrt{d} \right)$$

859 and \widetilde{O} absorbs constants and Polylog terms.

¹⁴Note that here the assumption is stronger than the assumption in [Yin et al., 2022], therefore the conclusion of Lemma C.4 holds.

860 *Proof of Lemma E.13.* The proof is by combining (40), (42), Lemma E.4, Lemma E.6, Lemma E.7,
 861 Lemma E.8, Lemma E.12 and a union bound. \square

862 **Remark E.14.** Under the same assumption of Lemma E.13, because of Remark E.9 and Lemma E.11,
 863 we have with probability $1 - \delta$, for all $h, s, a \in [H] \times \mathcal{S} \times \mathcal{A}$,

$$\begin{aligned} & \left| (\mathcal{T}_h \tilde{V}_{h+1} - \tilde{\mathcal{T}}_h \tilde{V}_{h+1})(s, a) \right| \leq \tilde{O} \left(\sqrt{d} \sqrt{\phi(s, a)^\top \tilde{\Lambda}_h^{-1} \phi(s, a)} \right) + \frac{D}{K} \\ & \leq \tilde{O} \left(\sqrt{d} \sqrt{\phi(s, a)^\top \hat{\Lambda}_h^{-1} \phi(s, a)} \right) + \frac{2D}{K} \\ & \leq \tilde{O} \left(2\sqrt{d} \sqrt{\phi(s, a)^\top \Lambda_h^{-1} \phi(s, a)} \right) + \frac{2D}{K}. \end{aligned} \quad (54)$$

864 Because $D = \tilde{O} \left(\frac{H^2 L}{\kappa} + \frac{H^4 E \sqrt{d}}{\kappa^{3/2}} + H^3 \sqrt{d} \right)$ and \tilde{O} absorbs constant, we will write as below for
 865 simplicity:

$$\left| (\mathcal{T}_h \tilde{V}_{h+1} - \tilde{\mathcal{T}}_h \tilde{V}_{h+1})(s, a) \right| \leq \tilde{O} \left(\sqrt{d} \sqrt{\phi(s, a)^\top \Lambda_h^{-1} \phi(s, a)} \right) + \frac{D}{K}. \quad (55)$$

866 E.3.5 Finalize the proof of the first part

867 We are ready to prove the first part of Theorem 4.1.

868 **Theorem E.15** (First part of Theorem 4.1). Let K be the number of episodes. Suppose $\sqrt{d} > \xi$, where

$$869 \xi := \sup_{V \in [0, H], s' \sim P_h(s, a), h \in [H]} \left| \frac{r_h + V(s') - (\mathcal{T}_h V)(s, a)}{\sigma_V(s, a)} \right| \text{ and } K > \max\{\mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_3, \mathcal{M}_4\}.$$

870 Then for any $0 < \lambda < \kappa$, with probability $1 - \delta$, for all policy π simultaneously, the output $\hat{\pi}$
 871 of Algorithm 2 satisfies

$$v^\pi - v^{\hat{\pi}} \leq \tilde{O} \left(\sqrt{d} \cdot \sum_{h=1}^H \mathbb{E}_\pi \left[\left(\phi(\cdot, \cdot)^\top \Lambda_h^{-1} \phi(\cdot, \cdot) \right)^{1/2} \right] \right) + \frac{DH}{K},$$

872 where $\Lambda_h = \sum_{\tau=1}^K \frac{\phi(s_h^\tau, a_h^\tau) \cdot \phi(s_h^\tau, a_h^\tau)^\top}{\sigma_{V_{h+1}}^2(s_h^\tau, a_h^\tau)} + \lambda I_d$, $D = \tilde{O} \left(\frac{H^2 L}{\kappa} + \frac{H^4 E \sqrt{d}}{\kappa^{3/2}} + H^3 \sqrt{d} \right)$ and \tilde{O} absorbs
 873 constants and Polylog terms.

874 *Proof of Theorem E.15.* Combining Lemma E.3 and Remark E.14, we have with probability $1 - \delta$,
 875 for all policy π simultaneously,

$$V_1^\pi(s) - V_1^{\hat{\pi}}(s) \leq \tilde{O} \left(\sqrt{d} \cdot \sum_{h=1}^H \mathbb{E}_\pi \left[\left(\phi(\cdot, \cdot)^\top \Lambda_h^{-1} \phi(\cdot, \cdot) \right)^{1/2} \middle| s_1 = s \right] \right) + \frac{DH}{K}, \quad (56)$$

876 now the proof is complete by taking the initial distribution d_1 on both sides. \square

877 E.3.6 Finalize the proof of the second part

878 To prove the second part of Theorem 4.1, we begin with a crude bound on $\sup_h \|V_h^* - \tilde{V}_h\|_\infty$.

879 **Lemma E.16** (Private version of Lemma C.8 in [Yin et al., 2022]). Suppose $K \geq$
 880 $\max\{\mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_3, \mathcal{M}_4\}$, under the high probability event in Lemma E.13, with probability at
 881 least $1 - \delta$,

$$\sup_h \|V_h^* - \tilde{V}_h\|_\infty \leq \tilde{O} \left(\frac{H^2 \sqrt{d}}{\sqrt{\kappa K}} \right).$$

882 *Proof of Lemma E.16. Step 1:* The first step is to show with probability at least $1 - \delta$,

$$883 \sup_h \|V_h^* - V_h^{\hat{\pi}}\|_\infty \leq \tilde{O} \left(\frac{H^2 \sqrt{d}}{\sqrt{\kappa K}} \right).$$

884 Indeed, combine Lemma E.3 and Lemma E.13, similar to the proof of Theorem E.15, we directly
 885 have with probability $1 - \delta$, for all policy π simultaneously, and for all $s \in \mathcal{S}$, $h \in [H]$,

$$V_h^\pi(s) - V_h^{\hat{\pi}}(s) \leq \tilde{O} \left(\sqrt{d} \cdot \sum_{t=h}^H \mathbb{E}_\pi \left[(\phi(\cdot, \cdot)^\top \Lambda_t^{-1} \phi(\cdot, \cdot))^{1/2} \middle| s_h = s \right] \right) + \frac{DH}{K}, \quad (57)$$

886 Next, since $K \geq \max \left\{ \frac{512 \log(\frac{2Hd}{\delta})}{\kappa^2}, \frac{4\lambda}{\kappa} \right\}$, by Lemma F.13 and a union bound over $h \in [H]$, with
 887 probability $1 - \delta$,

$$\sup_{s,a} \|\phi(s, a)\|_{\Lambda_h^{-1}} \leq \frac{2}{\sqrt{K}} \sup_{s,a} \|\phi(s, a)\|_{(\Lambda_h^p)^{-1}} \leq \frac{2}{\sqrt{K}} \sqrt{\lambda_{\min}^{-1}(\Lambda_h^p)} \leq \frac{2H}{\sqrt{\kappa K}}, \quad \forall h \in [H],$$

888 where $\Lambda_h^p = \mathbb{E}_{\mu, h} [\sigma_{\tilde{V}_{h+1}}^{-2}(s, a) \phi(s, a) \phi(s, a)^\top]$ and $\lambda_{\min}(\Lambda_h^p) \geq \frac{\kappa}{H^2}$.

889 Lastly, taking $\pi = \pi^*$ in (57) to obtain

$$\begin{aligned} 0 \leq V_h^{\pi^*}(s) - V_h^{\hat{\pi}}(s) &\leq \tilde{O} \left(\sqrt{d} \cdot \sum_{t=h}^H \mathbb{E}_{\pi^*} \left[(\phi(\cdot, \cdot)^\top \Lambda_t^{-1} \phi(\cdot, \cdot))^{1/2} \middle| s_h = s \right] \right) + \frac{DH}{K} \\ &\leq \tilde{O} \left(\frac{H^2 \sqrt{d}}{\sqrt{\kappa K}} \right) + \tilde{O} \left(\frac{H^3 L / \kappa}{K} + \frac{H^5 E \sqrt{d} / \kappa^{3/2}}{K} + \frac{H^4 \sqrt{d}}{K} \right). \end{aligned} \quad (58)$$

890 This implies by using the condition $K > \max\{\frac{H^2 L^2}{d\kappa}, \frac{H^6 E^2}{\kappa^2}, H^4 \kappa\}$, we finish the proof of Step 1.

891 **Step 2:** The second step is to show with probability $1 - \delta$, $\sup_h \|\tilde{V}_h - V_h^{\hat{\pi}}\|_\infty \leq \tilde{O} \left(\frac{H^2 \sqrt{d}}{\sqrt{\kappa K}} \right)$.

892 Indeed, applying Lemma F.7 with $\pi = \pi' = \hat{\pi}$, then with probability $1 - \delta$, for all s, h

$$\begin{aligned} \left| \tilde{V}_h(s) - V_h^{\hat{\pi}}(s) \right| &= \left| \sum_{t=h}^H \mathbb{E}_{\hat{\pi}} \left[\hat{Q}_h(s_h, a_h) - (\mathcal{T}_h \tilde{V}_{h+1})(s_h, a_h) \middle| s_h = s \right] \right| \\ &\leq \sum_{t=h}^H \left\| (\tilde{\mathcal{T}}_h \tilde{V}_{h+1} - \mathcal{T}_h \tilde{V}_{h+1})(s, a) \right\|_\infty + H \cdot \|\Gamma_h(s, a)\|_\infty \\ &\leq \tilde{O} \left(H \sqrt{d} \left\| \sqrt{\phi(s, a)^\top \Lambda_h^{-1} \phi(s, a)} \right\|_\infty \right) + \tilde{O} \left(\frac{DH}{K} \right) \\ &\leq \tilde{O} \left(\frac{H^2 \sqrt{d}}{\sqrt{\kappa K}} \right), \end{aligned}$$

893 where the second inequality uses Lemma E.13, Remark E.14 and the last inequality holds due to the
 894 same reason as Step 1.

895 **Step 3:** The proof of the lemma is complete by combining Step 1, Step 2, triangular inequality and a
 896 union bound.

897 □

898 Then we can give a high probability bound of $\sup_h \|\sigma_{\tilde{V}_{h+1}}^2 - \sigma_h^{*2}\|_\infty$.

899 **Lemma E.17** (Private version of Lemma C.10 in [Yin et al., 2022]). Recall $\sigma_{\tilde{V}_{h+1}}^2 =$
 900 $\max \left\{ 1, \text{Var}_{P_h} \tilde{V}_{h+1} \right\}$ and $\sigma_h^{*2} = \max \left\{ 1, \text{Var}_{P_h} V_{h+1}^* \right\}$. Suppose $K \geq \max\{\mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_3, \mathcal{M}_4\}$,
 901 then with probability $1 - \delta$,

$$\sup_h \|\sigma_{\tilde{V}_{h+1}}^2 - \sigma_h^{*2}\|_\infty \leq \tilde{O} \left(\frac{H^3 \sqrt{d}}{\sqrt{\kappa K}} \right).$$

902 *Proof of Lemma E.17.* By definition and the non-expansiveness of $\max\{1, \cdot\}$, we have

$$\begin{aligned}
& \left\| \sigma_{\tilde{V}_{h+1}}^2 - \sigma_h^{*2} \right\|_{\infty} \leq \left\| \text{Var} \tilde{V}_{h+1} - \text{Var} V_{h+1}^* \right\|_{\infty} \\
& \leq \left\| \mathbb{P}_h \left(\tilde{V}_{h+1}^2 - V_{h+1}^{*2} \right) \right\|_{\infty} + \left\| (\mathbb{P}_h \tilde{V}_{h+1})^2 - (\mathbb{P}_h V_{h+1}^*)^2 \right\|_{\infty} \\
& \leq \left\| \tilde{V}_{h+1}^2 - V_{h+1}^{*2} \right\|_{\infty} + \left\| (\mathbb{P}_h \tilde{V}_{h+1} + \mathbb{P}_h V_{h+1}^*) (\mathbb{P}_h \tilde{V}_{h+1} - \mathbb{P}_h V_{h+1}^*) \right\|_{\infty} \\
& \leq 2H \left\| \tilde{V}_{h+1} - V_{h+1}^* \right\|_{\infty} + 2H \left\| \mathbb{P}_h \tilde{V}_{h+1} - \mathbb{P}_h V_{h+1}^* \right\|_{\infty} \\
& \leq \tilde{O} \left(\frac{H^3 \sqrt{d}}{\sqrt{\kappa K}} \right).
\end{aligned}$$

903 The second inequality is because of the definition of variance. The last inequality comes from
904 Lemma E.16. \square

905 We transfer $\sqrt{\phi(s, a)^\top \Lambda_h^{-1} \phi(s, a)}$ to $\sqrt{\phi(s, a)^\top \Lambda_h^{*-1} \phi(s, a)}$ by the following Lemma E.18.

906 **Lemma E.18** (Private version of Lemma C.11 in [Yin et al., 2022]). *Suppose $K \geq$*
907 *$\max\{\mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_3, \mathcal{M}_4\}$, then with probability $1 - \delta$,*

$$\sqrt{\phi(s, a)^\top \Lambda_h^{-1} \phi(s, a)} \leq 2 \sqrt{\phi(s, a)^\top \Lambda_h^{*-1} \phi(s, a)}, \quad \forall h, s, a \in [H] \times \mathcal{S} \times \mathcal{A},$$

908 *Proof of Lemma E.18.* By definition $\sqrt{\phi(s, a)^\top \Lambda_h^{-1} \phi(s, a)} = \|\phi(s, a)\|_{\Lambda_h^{-1}}$. Then denote

$$\Lambda'_h = \frac{1}{K} \Lambda_h, \quad \Lambda_h^{*'} = \frac{1}{K} \Lambda_h^*,$$

909 where $\Lambda_h^* = \sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau) \phi(s_h^\tau, a_h^\tau)^\top / \sigma_{V_{h+1}^*}^2(s_h^\tau, a_h^\tau) + \lambda I$. Under the condition of K , by
910 Lemma E.17, with probability $1 - \delta$, for all $h \in [H]$,

$$\begin{aligned}
& \left\| \Lambda_h^{*'} - \Lambda'_h \right\| \leq \sup_{s, a} \left\| \frac{\phi(s, a) \phi(s, a)^\top}{\sigma_h^{*2}(s, a)} - \frac{\phi(s, a) \phi(s, a)^\top}{\sigma_{\tilde{V}_{h+1}}^2(s, a)} \right\| \\
& \leq \sup_{s, a} \left| \frac{\sigma_h^{*2}(s, a) - \sigma_{\tilde{V}_{h+1}}^2(s, a)}{\sigma_h^{*2}(s, a) \cdot \sigma_{\tilde{V}_{h+1}}^2(s, a)} \right| \cdot \|\phi(s, a)\|^2 \\
& \leq \sup_{s, a} \left| \frac{\sigma_h^{*2}(s, a) - \sigma_{\tilde{V}_{h+1}}^2(s, a)}{1} \right| \cdot 1 \\
& \leq \tilde{O} \left(\frac{H^3 \sqrt{d}}{\sqrt{\kappa K}} \right).
\end{aligned} \tag{59}$$

911 Next by Lemma F.12 (with ϕ to be $\phi / \sigma_{V_{h+1}^*}$ and $C = 1$), it holds with probability $1 - \delta$,

$$\left\| \Lambda_h^{*'} - \left(\mathbb{E}_{\mu, h} [\phi(s, a) \phi(s, a)^\top / \sigma_{V_{h+1}^*}^2(s, a)] + \frac{\lambda}{K} I_d \right) \right\| \leq \frac{4\sqrt{2}}{\sqrt{K}} \left(\log \frac{2dH}{\delta} \right)^{1/2}.$$

912 Therefore by Weyl's inequality and the condition $K > \max\{2\lambda, 128 \log(\frac{2dH}{\delta}), \frac{128H^4 \log(2dH/\delta)}{\kappa^2}\}$,
 913 the above inequality implies

$$\begin{aligned}
 \|\Lambda_h^{*\prime}\| &= \lambda_{\max}(\Lambda_h^{*\prime}) \leq \lambda_{\max}\left(\mathbb{E}_{\mu,h}[\phi(s,a)\phi(s,a)^\top / \sigma_{V_{h+1}^*}^2(s,a)]\right) + \frac{\lambda}{K} + \frac{4\sqrt{2}}{\sqrt{K}} \left(\log \frac{2dH}{\delta}\right)^{1/2} \\
 &\leq \left\|\mathbb{E}_{\mu,h}[\phi(s,a)\phi(s,a)^\top / \sigma_{V_{h+1}^*}^2(s,a)]\right\| + \frac{\lambda}{K} + \frac{4\sqrt{2}}{\sqrt{K}} \left(\log \frac{2dH}{\delta}\right)^{1/2} \\
 &\leq \|\phi(s,a)\|^2 + \frac{\lambda}{K} + \frac{4\sqrt{2}}{\sqrt{K}} \left(\log \frac{2dH}{\delta}\right)^{1/2} \leq 1 + \frac{\lambda}{K} + \frac{4\sqrt{2}}{\sqrt{K}} \left(\log \frac{2dH}{\delta}\right)^{1/2} \leq 2, \\
 \lambda_{\min}(\Lambda_h^{*\prime}) &\geq \lambda_{\min}\left(\mathbb{E}_{\mu,h}[\phi(s,a)\phi(s,a)^\top / \sigma_{V_{h+1}^*}^2(s,a)]\right) + \frac{\lambda}{K} - \frac{4\sqrt{2}}{\sqrt{K}} \left(\log \frac{2dH}{\delta}\right)^{1/2} \\
 &\geq \lambda_{\min}\left(\mathbb{E}_{\mu,h}[\phi(s,a)\phi(s,a)^\top / \sigma_{V_{h+1}^*}^2(s,a)]\right) - \frac{4\sqrt{2}}{\sqrt{K}} \left(\log \frac{2dH}{\delta}\right)^{1/2} \\
 &\geq \frac{\kappa}{H^2} - \frac{4\sqrt{2}}{\sqrt{K}} \left(\log \frac{2dH}{\delta}\right)^{1/2} \geq \frac{\kappa}{2H^2}.
 \end{aligned}$$

914 Hence with probability $1 - \delta$, $\|\Lambda_h^{*\prime}\| \leq 2$ and $\|\Lambda_h^{*\prime-1}\| = \lambda_{\min}^{-1}(\Lambda_h^{*\prime}) \leq \frac{2H^2}{\kappa}$. Similarly, we can show
 915 that $\|\Lambda_h^{\prime-1}\| \leq \frac{2H^2}{\kappa}$ holds with probability $1 - \delta$ by using identical proof.

916 Now apply Lemma F.11 and a union bound to $\Lambda_h^{*\prime}$ and Λ_h^{\prime} , we obtain with probability $1 - \delta$, for all
 917 $h, s, a \in [H] \times \mathcal{S} \times \mathcal{A}$,

$$\begin{aligned}
 \|\phi(s,a)\|_{\Lambda_h^{\prime-1}} &\leq \left[1 + \sqrt{\|\Lambda_h^{*\prime-1}\| \cdot \|\Lambda_h^{*\prime}\| \cdot \|\Lambda_h^{\prime-1}\| \cdot \|\Lambda_h^{*\prime} - \Lambda_h^{\prime}\|}\right] \cdot \|\phi(s,a)\|_{\Lambda_h^{*\prime-1}} \\
 &\leq \left[1 + \sqrt{\frac{2H^2}{\kappa} \cdot 2 \cdot \frac{2H^2}{\kappa} \cdot \|\Lambda_h^{*\prime} - \Lambda_h^{\prime}\|}\right] \cdot \|\phi(s,a)\|_{\Lambda_h^{*\prime-1}} \\
 &\leq \left[1 + \sqrt{\frac{H^4}{\kappa^2} \left[\tilde{O}\left(\frac{H^3\sqrt{d}}{\sqrt{\kappa K}}\right)\right]}\right] \cdot \|\phi(s,a)\|_{\Lambda_h^{*\prime-1}} \\
 &\leq 2 \|\phi(s,a)\|_{\Lambda_h^{*\prime-1}}
 \end{aligned}$$

918 where the third inequality uses (59) and the last inequality uses $K \geq \tilde{O}(H^{14}d/\kappa^5)$. The conclusion
 919 can be derived directly by the above inequality multiplying $1/\sqrt{K}$ on both sides. \square

920 Finally, the second part of Theorem 4.1 can be proven by combining Theorem E.15 (with $\pi = \pi^*$)
 921 and Lemma E.18.

922 E.4 Put everything together

923 Combining Lemma E.1, Theorem E.15, and the discussion above, the proof of Theorem 4.1 is
 924 complete.

925 F Assisting technical lemmas

Lemma F.1 (Multiplicative Chernoff bound [Chernoff et al., 1952]). *Let X be a Binomial random variable with parameter p, n . For any $1 \geq \theta > 0$, we have that*

$$\mathbb{P}[X < (1 - \theta)pn] < e^{-\frac{\theta^2 pn}{2}}, \quad \text{and} \quad \mathbb{P}[X \geq (1 + \theta)pn] < e^{-\frac{\theta^2 pn}{3}}$$

Lemma F.2 (Hoeffding's Inequality [Sridharan, 2002]). *Let x_1, \dots, x_n be independent bounded random variables such that $\mathbb{E}[x_i] = 0$ and $|x_i| \leq \xi_i$ with probability 1. Then for any $\epsilon > 0$ we have*

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n x_i \geq \epsilon\right) \leq e^{-\frac{2n^2 \epsilon^2}{\sum_{i=1}^n \xi_i^2}}.$$

Lemma F.3 (Bernstein's Inequality). *Let x_1, \dots, x_n be independent bounded random variables such that $\mathbb{E}[x_i] = 0$ and $|x_i| \leq \xi$ with probability 1. Let $\sigma^2 = \frac{1}{n} \sum_{i=1}^n \text{Var}[x_i]$, then with probability $1 - \delta$ we have*

$$\frac{1}{n} \sum_{i=1}^n x_i \leq \sqrt{\frac{2\sigma^2 \cdot \log(1/\delta)}{n}} + \frac{2\xi}{3n} \log(1/\delta).$$

Lemma F.4 (Empirical Bernstein's Inequality [Maurer and Pontil, 2009]). *Let x_1, \dots, x_n be i.i.d random variables such that $|x_i| \leq \xi$ with probability 1. Let $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ and $\widehat{V}_n = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$, then with probability $1 - \delta$ we have*

$$\left| \frac{1}{n} \sum_{i=1}^n x_i - \mathbb{E}[x] \right| \leq \sqrt{\frac{2\widehat{V}_n \cdot \log(2/\delta)}{n}} + \frac{7\xi}{3n} \log(2/\delta).$$

926 **Lemma F.5** (Lemma I.8 in [Yin and Wang, 2021b]). *Let $n \geq 2$ and $V \in \mathbb{R}^S$ be any function with*
 927 $\|V\|_\infty \leq H$, P be any S -dimensional distribution and \widehat{P} be its empirical version using n samples.
 928 *Then with probability $1 - \delta$,*

$$\left| \sqrt{\text{Var}_{\widehat{P}}(V)} - \sqrt{\frac{n-1}{n} \text{Var}_P(V)} \right| \leq 2H \sqrt{\frac{\log(2/\delta)}{n-1}}.$$

929 **Lemma F.6** (Claim 2 in [Vietri et al., 2020]). *Let $y \in \mathbb{R}$ be any positive real number. Then for all*
 930 $x \in \mathbb{R}$ with $x \geq 2y$, *it holds that $\frac{1}{x-y} \leq \frac{1}{x} + \frac{2y}{x^2}$.*

931 F.1 Extended Value Difference

932 **Lemma F.7** (Extended Value Difference (Section B.1 in [Cai et al., 2020])). *Let $\pi = \{\pi_h\}_{h=1}^H$ and*
 933 $\pi' = \{\pi'_h\}_{h=1}^H$ *be two arbitrary policies and let $\{\widehat{Q}_h\}_{h=1}^H$ be any given Q -functions. Then define*
 934 $\widehat{V}_h(s) := \langle \widehat{Q}_h(s, \cdot), \pi_h(\cdot | s) \rangle$ *for all $s \in \mathcal{S}$. Then for all $s \in \mathcal{S}$,*

$$\begin{aligned} \widehat{V}_1(s) - V_1^{\pi'}(s) &= \sum_{h=1}^H \mathbb{E}_{\pi'} \left[\langle \widehat{Q}_h(s_h, \cdot), \pi_h(\cdot | s_h) - \pi'_h(\cdot | s_h) \rangle \mid s_1 = s \right] \\ &\quad + \sum_{h=1}^H \mathbb{E}_{\pi'} \left[\widehat{Q}_h(s_h, a_h) - (\mathcal{T}_h \widehat{V}_{h+1})(s_h, a_h) \mid s_1 = s \right] \end{aligned} \quad (60)$$

935 *where $(\mathcal{T}_h V)(\cdot, \cdot) := r_h(\cdot, \cdot) + (P_h V)(\cdot, \cdot)$ for any $V \in \mathbb{R}^S$.*

936 **Lemma F.8** (Lemma I.10 in [Yin and Wang, 2021b]). *Let $\widehat{\pi} = \{\widehat{\pi}_h\}_{h=1}^H$ and $\widehat{Q}_h(\cdot, \cdot)$ be the*
 937 *arbitrary policy and Q -function and also $\widehat{V}_h(s) = \langle \widehat{Q}_h(s, \cdot), \widehat{\pi}_h(\cdot | s) \rangle \forall s \in \mathcal{S}$, and $\xi_h(s, a) =$
 938 $(\mathcal{T}_h \widehat{V}_{h+1})(s, a) - \widehat{Q}_h(s, a)$ *element-wisely. Then for any arbitrary π , we have**

$$\begin{aligned} V_1^\pi(s) - V_1^{\widehat{\pi}}(s) &= \sum_{h=1}^H \mathbb{E}_\pi [\xi_h(s_h, a_h) \mid s_1 = s] - \sum_{h=1}^H \mathbb{E}_{\widehat{\pi}} [\xi_h(s_h, a_h) \mid s_1 = s] \\ &\quad + \sum_{h=1}^H \mathbb{E}_\pi \left[\langle \widehat{Q}_h(s_h, \cdot), \pi_h(\cdot | s_h) - \widehat{\pi}_h(\cdot | s_h) \rangle \mid s_1 = s \right] \end{aligned}$$

939 *where the expectation are taken over s_h, a_h .*

940 F.2 Assisting lemmas for linear MDP setting

941 **Lemma F.9** (Hoeffding inequality for self-normalized martingales [Abbasi-Yadkori et al., 2011]).
 942 *Let $\{\eta_t\}_{t=1}^\infty$ be a real-valued stochastic process. Let $\{\mathcal{F}_t\}_{t=0}^\infty$ be a filtration, such that η_t is \mathcal{F}_t -*
 943 *measurable. Assume η_t also satisfies η_t given \mathcal{F}_{t-1} is zero-mean and R -subgaussian, i.e.*

$$\forall \lambda \in \mathbb{R}, \quad \mathbb{E} [e^{\lambda \eta_t} \mid \mathcal{F}_{t-1}] \leq e^{\lambda^2 R^2 / 2}.$$

944 Let $\{x_t\}_{t=1}^\infty$ be an \mathbb{R}^d -valued stochastic process where x_t is \mathcal{F}_{t-1} measurable and $\|x_t\| \leq L$. Let
 945 $\Lambda_t = \lambda I_d + \sum_{s=1}^t x_s x_s^\top$. Then for any $\delta > 0$, with probability $1 - \delta$, for all $t > 0$,

$$\left\| \sum_{s=1}^t x_s \eta_s \right\|_{\Lambda_t^{-1}}^2 \leq 8R^2 \cdot \frac{d}{2} \log \left(\frac{\lambda + tL}{\lambda\delta} \right).$$

946 **Lemma F.10** (Bernstein inequality for self-normalized martingales [Zhou et al., 2021]). Let $\{\eta_t\}_{t=1}^\infty$
 947 be a real-valued stochastic process. Let $\{\mathcal{F}_t\}_{t=0}^\infty$ be a filtration, such that η_t is \mathcal{F}_t -measurable.
 948 Assume η_t also satisfies

$$|\eta_t| \leq R, \mathbb{E}[\eta_t | \mathcal{F}_{t-1}] = 0, \mathbb{E}[\eta_t^2 | \mathcal{F}_{t-1}] \leq \sigma^2.$$

949 Let $\{x_t\}_{t=1}^\infty$ be an \mathbb{R}^d -valued stochastic process where x_t is \mathcal{F}_{t-1} measurable and $\|x_t\| \leq L$. Let
 950 $\Lambda_t = \lambda I_d + \sum_{s=1}^t x_s x_s^\top$. Then for any $\delta > 0$, with probability $1 - \delta$, for all $t > 0$,

$$\left\| \sum_{s=1}^t x_s \eta_s \right\|_{\Lambda_t^{-1}} \leq 8\sigma \sqrt{d \log \left(1 + \frac{tL^2}{\lambda d} \right) \cdot \log \left(\frac{4t^2}{\delta} \right) + 4R \log \left(\frac{4t^2}{\delta} \right)}$$

951 **Lemma F.11** (Lemma H.4 in [Yin et al., 2022]). Let Λ_1 and $\Lambda_2 \in \mathbb{R}^{d \times d}$ be two positive semi-definite
 952 matrices. Then:

$$\|\Lambda_1^{-1}\| \leq \|\Lambda_2^{-1}\| + \|\Lambda_1^{-1}\| \cdot \|\Lambda_2^{-1}\| \cdot \|\Lambda_1 - \Lambda_2\|$$

953 and

$$\|\phi\|_{\Lambda_1^{-1}} \leq \left[1 + \sqrt{\|\Lambda_2^{-1}\| \cdot \|\Lambda_2\| \cdot \|\Lambda_1^{-1}\| \cdot \|\Lambda_1 - \Lambda_2\|} \right] \cdot \|\phi\|_{\Lambda_2^{-1}}.$$

954 for all $\phi \in \mathbb{R}^d$.

955 **Lemma F.12** (Lemma H.4 in [Min et al., 2021]). Let $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ satisfies $\|\phi(s, a)\| \leq C$ for
 956 all $s, a \in \mathcal{S} \times \mathcal{A}$. For any $K > 0, \lambda > 0$, define $\tilde{G}_K = \sum_{k=1}^K \phi(s_k, a_k) \phi(s_k, a_k)^\top + \lambda I_d$ where
 957 (s_k, a_k) 's are i.i.d samples from some distribution ν . Then with probability $1 - \delta$,

$$\left\| \frac{\tilde{G}_K}{K} - \mathbb{E}_\nu \left[\frac{\tilde{G}_K}{K} \right] \right\| \leq \frac{4\sqrt{2}C^2}{\sqrt{K}} \left(\log \frac{2d}{\delta} \right)^{1/2}.$$

958 **Lemma F.13** (Lemma H.5 in [Min et al., 2021]). Let $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ be a bounded function s.t.
 959 $\|\phi\|_2 \leq C$. Define $\tilde{G}_K = \sum_{k=1}^K \phi(s_k, a_k) \phi(s_k, a_k)^\top + \lambda I_d$ where (s_k, a_k) 's are i.i.d samples from
 960 some distribution ν . Let $G = \mathbb{E}_\nu[\phi(s, a) \phi(s, a)^\top]$. Then for any $\delta \in (0, 1)$, if K satisfies

$$K \geq \max \left\{ 512C^4 \|\mathbf{G}^{-1}\|^2 \log \left(\frac{2d}{\delta} \right), 4\lambda \|\mathbf{G}^{-1}\| \right\}.$$

961 Then with probability at least $1 - \delta$, it holds simultaneously for all $u \in \mathbb{R}^d$ that

$$\|u\|_{\tilde{G}_K^{-1}} \leq \frac{2}{\sqrt{K}} \|u\|_{G^{-1}}.$$

962 **Lemma F.14** (Lemma H.9 in [Yin et al., 2022]). For a linear MDP, for any $0 \leq V(\cdot) \leq H$, there
 963 exists a $w_h \in \mathbb{R}^d$ s.t. $\mathcal{T}_h V = \langle \phi, w_h \rangle$ and $\|w_h\|_2 \leq 2H\sqrt{d}$ for all $h \in [H]$. Here $\mathcal{T}_h(V)(s, a) =$
 964 $r_h(x, a) + (P_h V)(s, a)$. Similarly, for any π , there exists $w_h^\pi \in \mathbb{R}^d$, such that $Q_h^\pi = \langle \phi, w_h^\pi \rangle$ with
 965 $\|w_h^\pi\|_2 \leq 2(H - h + 1)\sqrt{d}$.

966 F.3 Assisting lemmas for differential privacy

967 **Lemma F.15** (Converting zCDP to DP [Bun and Steinke, 2016]). If M satisfies ρ -zCDP then M
 968 satisfies $(\rho + 2\sqrt{\rho \log(1/\delta)}, \delta)$ -DP.

969 **Lemma F.16** (zCDP Composition [Bun and Steinke, 2016]). Let $M : \mathcal{U}^n \rightarrow \mathcal{Y}$ and $M' : \mathcal{U}^n \rightarrow \mathcal{Z}$
 970 be randomized mechanisms. Suppose that M satisfies ρ -zCDP and M' satisfies ρ' -zCDP. Define
 971 $M'' : \mathcal{U}^n \rightarrow \mathcal{Y} \times \mathcal{Z}$ by $M''(U) = (M(U), M'(U))$. Then M'' satisfies $(\rho + \rho')$ -zCDP.

972 **Lemma F.17** (Adaptive composition and Post processing of zCDP [Bun and Steinke, 2016]). *Let*
 973 $M : \mathcal{X}^n \rightarrow \mathcal{Y}$ *and* $M' : \mathcal{X}^n \times \mathcal{Y} \rightarrow \mathcal{Z}$. *Suppose* M *satisfies* ρ -zCDP *and* M' *satisfies* ρ' -zCDP
 974 *(as a function of its first argument). Define* $M'' : \mathcal{X}^n \rightarrow \mathcal{Z}$ *by* $M''(x) = M'(x, M(x))$. *Then* M''
 975 *satisfies* $(\rho + \rho')$ -zCDP.

976 **Definition F.18** (ℓ_1 sensitivity). *Define the* ℓ_1 *sensitivity of a function* $f : \mathbb{N}^{\mathcal{X}} \mapsto \mathbb{R}^d$ *as*

$$\Delta_1(f) = \sup_{\text{neighboring } U, U'} \|f(U) - f(U')\|_1.$$

977 **Definition F.19** (Laplace Mechanism [Dwork et al., 2014]). *Given any function* $f : \mathbb{N}^{\mathcal{X}} \mapsto \mathbb{R}^d$, *the*
 978 *Laplace mechanism is defined as:*

$$\mathcal{M}_L(x, f, \epsilon) = f(x) + (Y_1, \dots, Y_d),$$

979 *where* Y_i *are i.i.d. random variables drawn from* $\text{Lap}(\Delta_1(f)/\epsilon)$.

980 **Lemma F.20** (Privacy guarantee of Laplace Mechanism [Dwork et al., 2014]). *The Laplace mecha-*
 981 *nism preserves* $(\epsilon, 0)$ -*differential privacy. For simplicity, we say* ϵ -DP.

982 G Details for the Evaluation part

983 In the Evaluation part, we apply a synthetic linear MDP case that is similar to [Min et al., 2021, Yin
 984 et al., 2022] but with some modifications for our evaluation task. The linear MDP example we use
 985 consists of $|\mathcal{S}| = 2$ states and $|\mathcal{A}| = 100$ actions, while the feature dimension $d = 10$. We denote
 986 $\mathcal{S} = \{0, 1\}$ and $\mathcal{A} = \{0, 1, \dots, 99\}$ respectively. For each action $a \in \{0, 1, \dots, 99\}$, we obtain a
 987 vector $\mathbf{a} \in \mathbb{R}^8$ via binary encoding. More specifically, each coordinate of \mathbf{a} is either 0 or 1.

988 First, we define the following indicator function $\delta(s, a) = \begin{cases} 1 & \text{if } \mathbb{1}\{s = 0\} = \mathbb{1}\{a = 0\} \\ 0 & \text{otherwise} \end{cases}$, then
 989 our non-stationary linear MDP example can be characterized by the following parameters.

990

The feature map ϕ is:

$$\phi(s, a) = (\mathbf{a}^\top, \delta(s, a), 1 - \delta(s, a))^\top \in \mathbb{R}^{10}.$$

The unknown measure ν_h is:

$$\begin{aligned} \nu_h(0) &= (0, \dots, 0, \alpha_{h,1}, \alpha_{h,2}), \\ \nu_h(1) &= (0, \dots, 0, 1 - \alpha_{h,1}, 1 - \alpha_{h,2}), \end{aligned}$$

where $\{\alpha_{h,1}, \alpha_{h,2}\}_{h \in [H]}$ is a sequence of random values sampled uniformly from $[0, 1]$.

The unknown vector θ_h is:

$$\theta_h = (r_h/8, 0, r_h/8, 1/2 - r_h/2, r_h/8, 0, r_h/8, 0, r_h/2, 1/2 - r_h/2) \in \mathbb{R}^{10},$$

991 where r_h is also sampled uniformly from $[0, 1]$. Therefore, the transition kernel follows $P_h(s'|s, a) =$
 992 $\langle \phi(s, a), \nu_h(s') \rangle$ and the expected reward function $r_h(s, a) = \langle \phi(s, a), \theta_h \rangle$.

993 Finally, the behavior policy is to always choose action $a = 0$ with probability p , and other actions
 994 uniformly with probability $(1 - p)/99$. Here we choose $p = 0.6$. The initial distribution is a uniform
 995 distribution over $\mathcal{S} = \{0, 1\}$.