# Is Offline Decision Making Possible with Only Few Samples? Reliable Decisions in Data Starved Bandits via Trust Region Enhance MENT

Anonymous authors

Paper under double-blind review

## ABSTRACT

What can an agent learn in a stochastic Multi-Armed Bandit (MAB) problem from a dataset that contains just a single sample for each arm? Surprisingly, in this work, we demonstrate that even in such a data-starved setting it may still be possible to find a policy competitive with the optimal one. This paves the way to reliable decision-making in settings where critical decisions must be made by relying only on a handful of samples. Our analysis reveals that *stochastic policies can be substantially better* than deterministic ones for offline decision-making. Focusing on offline multi-armed bandits, we design an algorithm called Trust Region of Uncertainty for Stochastic policy enhancemenT (TRUST) which is quite different from the predominant value-based lower confidence bound approach. Its design is enabled by localization laws, critical radii, and relative pessimism. We prove that its sample complexity is comparable to that of LCB on minimax problems while being substantially lower on problems with very few samples. Finally, we consider an application to offline reinforcement learning in the special case where the logging policies are known.

028 029

031

008

009

010 011 012

013

015

016

017

018

019

020

021

022

024

025

026

027

#### 1 INTRODUCTION

032 In several important problems, critical decisions must be made with just very few samples of pre-033 collected experience. For example, collecting samples in robotic manipulation may be slow and costly, 034 and the ability to learn from very few interactions is highly desirable (Hester & Stone, 2013; Liu et al., 2021). Likewise, in clinical trials and in personalized medical decisions, reliable decisions must be made by relying on very small datasets (Liu et al., 2017). Sample efficiency is also key in personalized 037 education (Bassen et al., 2020; Ruan et al., 2023). However, to achieve good performance, the state-038 of-the-art algorithms may require millions of samples (Fu et al., 2020). These empirical findings seem to be supported by the existing theories: the sample complexity bounds, even minimax optimal ones, can be large in practice due to the large constants and the warmup factors (Ménard et al., 2021; 040 Li et al., 2022; Azar et al., 2017; Zanette et al., 2019). 041

In this work, we study whether it is possible to make reliable decisions with only a few samples.
We focus on an offline Multi-Armed Bandit (MAB) problem, which is a foundation model for
decision-making (Lattimore & Szepesvári, 2020). In online MAB, an agent repeatedly chooses an
arm from a set of arms, each providing a stochastic reward. Offline MAB is a variant where the agent
cannot interact with the environment to gather new information and instead, it must make decisions
based on a pre-collected dataset without playing additional exploratory actions, aiming at identifying
the arm with the highest expected reward (Audibert et al., 2010; Garivier & Kaufmann, 2016; Russo,
2016; Ameko et al., 2020).

The standard approach to the problem is the Lower Confidence Bound (LCB) algorithm (Rashidinejad et al., 2021), a pessimistic variant of UCB (Auer et al., 2002) that involves selecting the arm with the highest lower bound on its performance. LCB encodes a principle called *pessimism under uncertainty*, which is the foundation principle for most algorithms for offline bandits and reinforcement learning (RL) (Jin et al., 2020; Zanette et al., 2020; Xie et al., 2021; Yin & Wang, 2021; Kumar et al.,

2020; Kostrikov et al., 2021). Unfortunately, the available methods that implement the principle of pessimism under uncertainty can fail in a data-starved regime because they rely on confidence intervals that are too loose when just a few samples are available. For example, even on a simple MAB instance with ten thousand arms, the best-known (Rashidinejad et al., 2021) performance bound for the LCB algorithm requires 24 samples per arm in order to provide meaningful guarantees, see Section 2. In more complex situations, such as in the sequential setting with function approximation, such a problem can become more severe due to the higher metric entropy of the function approximation class and the compounding of errors through time steps.

062 These considerations suggest that there is a "barrier of entry" to decision-making, both theoretically 063 and practically: one needs to have a substantial number of samples in order to make reliable decisions 064 even for settings as simple as offline MAB where the guarantees are tighter. Given the above technical 065 reasons, and the lack of good algorithms and guarantees for data-starved decision problems, it is 066 unclear whether it is even possible to find good decision rules with just a handful of samples.

In this paper, we make a substantial contribution towards lowering such barriers of entry. We discover that a carefully-designed algorithm tied to an advanced statistical analysis can substantially improve the sample complexity, both theoretically and practically, and enable reliable decision-making with just a handful of samples. More precisely, we focus on the offline MAB setting where we show that even if the dataset contains just a *single sample* in every arm, it may still be possible to compete with the optimal policy. This is remarkable, because with just one sample per arm—for example from a Bernoulli distribution—it is impossible to estimate the expected payoff of any of the arms! Our discovery is enabled by several key insights:

075 076

077

078

- We search over *stochastic* policies, which can yield better performance for offline-decision making;
- We use a *localized* notion of metric entropy to carefully control the size of the stochastic policy class that we search over;
- 079 080
- We implement a concept called *relative pessimism* to obtain sharper guarantees.

These considerations lead us to design a trust region policy optimization algorithm called Trust Region of Uncertainty for Stochastic policy enhancemenT (TRUST), one that offers superior theoretical as well as empirical performance compared to LCB in a data-scarce situation.

Moreover, we apply the algorithm to selected reinforcement learning problems from (Fu et al., 2020) in the special case where information about the logging policies is available. We do so by a simple reduction from reinforcement learning to bandits, by mapping policies and returns in the former to actions and rewards in the latter, thereby disregarding the sequential aspect of the problem. Although we rely on the information of the logging policies being available, the empirical study shows that our algorithm compares well with a strong deep reinforcement learning baseline (i.e, CQL from (Kumar et al., 2020)), without being sensitive to partial observability, sparse rewards, and hyper-parameters.

092 093

094

## 2 DATA-STARVED MULTI-ARMED BANDITS

In this section, we describe the MAB setting and give an example of a "data-starved" MAB instance
where prior methods (such as LCB) can fail. We informally say that an offline MAB is "data-starved"
if its dataset contains only very few samples in each arm.

**Notation.** We let  $[n] = \{1, 2, ..., n\}$  for a positive integer n. We let  $\|\cdot\|_2$  denote the Euclidean norm for vectors and the operator norm for matrices. We hide constants and logarithmic factors in the  $\widetilde{O}(\cdot)$ notation. We let  $\mathbb{B}_p^d(s) = \{x \in \mathbb{R}^d : \|x\|_p \leq s\}$  for any  $s \geq 0$  and  $p \geq 1$ .  $a \leq b$  ( $a \geq b$ ) means  $a \leq Cb$  ( $a \geq Cb$ ) for some numerical constant C.  $a \simeq b$  means that both  $a \leq b$  and  $b \leq a$  hold.

117

118

132 133

138

139

140

108 as  $a^* \in \arg \max_{a \in \mathcal{A}} [r(a)]$  and the single policy concentrability as  $C^* = 1/\mu(a^*)$  where  $\mu$  is the 109 distribution that generated the dataset. Without loss of generality, we assume the optimal arm is 110 unique. We also write  $r = (r_1, r_2, ..., r_d)^{\top}$ . Without loss of generality, we assume there is at least 111 one sample for each arm (such arm can otherwise be removed).

112 Lower confidence bound algorithm. One simple but effective method for the offline MAB prob-113 lem is the Lower Confidence Bound (LCB) algorithm, which is inspired by its online counterpart 114 (UCB) (Auer et al., 2002). Like UCB, LCB computes the empirical mean  $\hat{r}_i$  associated to the reward 115 of each arm i along with its half confidence width  $b_i$ . They are defined as 116

$$\widehat{r}_i := \frac{1}{N(a_i)} \sum_{k:x_k = a_i} x_k, \ b_i := \sqrt{\frac{2\sigma_i^2}{N(a_i)} \log\left(\frac{2d}{\delta}\right)}.$$
(1)

119 This definition ensures that each confidence interval brackets the corresponding expected reward with 120 probability  $1 - \delta$ : 121  $\widehat{r}$ 

$$\hat{a}_i - b_i \le r(a_i) \le \hat{r}_i + b_i \quad \forall i \in [d].$$
 (2)

122 The width of the confidence level depends on the noise level  $\sigma_i$ , which can be exploited by variance-123 aware methods (Zhang et al., 2021; Min et al., 2021; Yin et al., 2022; Dai et al., 2022). When the 124 true noise level is not accessible, we can replace it with the empirical standard deviation or with a 125 high-probability upper bound. For example, when the reward for each arm is restricted to be within 126 [0, 1], a simpler upper bound is  $\sigma_i^2 \leq 1/4$ .

127 Unlike UCB, the half-width of the confidence intervals for LCB is not added, but subtracted, from 128 the empirical mean, resulting in the lower bound  $l_i = \hat{r}_i - b_i$ . The action identified by LCB is then 129 the one that maximizes the resulting lower bound, thereby incorporating the principle of pessimism 130 under uncertainty (Jin et al., 2020; Kumar et al., 2020). Specifically, given the dataset  $\mathcal{D}$ , LCB selects 131 the arm using the following rule:

$$\widehat{a}_{\mathsf{LCB}} := \underset{a_i \in \mathcal{A}}{\operatorname{arg\,max}} l_i, \tag{3}$$

(Rashidinejad et al., 2021) analyzed the LCB strategy. Below we provide a modified version of their 134 theorem. 135

**Theorem 2.1** (LCB Performance). Suppose the noise of arm  $a_i$  is sub-Gaussian with proxy variance 136  $\sigma_i^2$ . Let  $\delta \in (0, 1/2)$ . Then, we have 137

- 1. (Comparison with any arm) With probability at least  $1 \delta$ , for any comparator policy  $a_i \in A$ , it holds that  $r(a_i) - r(\widehat{a}_{\mathsf{LCB}}) \leq \sqrt{8\sigma_i^2 \log(2d/\delta)/N(a_i)}$ .
- 2. (Comparison with the optimal arm) Assume  $\sigma_i = 1$  for any  $i \in [d]$  and  $N \ge 8C^* \log(1/\delta)$ . Then, 141 with probability at least  $1 - 2\delta$ , one has  $r(a^*) - r(\widehat{a}_{\mathsf{LCB}}) \leq \sqrt{4C^* \log(2d/\delta)/N}$ . 142

143 The statement of this theorem is slightly different from that in (Rashidinejad et al., 2021), in the sense 144 that their suboptimality is over  $\mathbb{E}_{\mathcal{D}}[r(a^*) - r(\hat{a}_{\mathsf{LCB}})]$  instead of a high-probability one. (Rashidinejad 145 et al., 2021) proved the minimax optimality of the algorithm when the single policy concentrability 146  $C^* \geq 2$  and the sample size  $N \geq O(C^*)$ . 147

A data-starved MAB problem and failure of LCB. In order to highlight the limitation of a strategy 148 such as LCB, let us describe a specific data-starved MAB instance, specifically one with d = 10000149 arms, equally partitioned into a set of good arms (i.e.,  $A_a$ ) and a set of bad arms (i.e.,  $A_b$ ). Each 150 good arm returns a reward following the uniform distribution over [0.5, 1.5], while each bad arm 151 returns a reward which follows  $\mathcal{N}(0, 1/4)$ . 152

Assume that we are given a dataset that contains only one sample per each arm. Instantiating the LCB 153 confidence interval in equation 2 with  $\sigma_i \leq 1/2$  and  $\delta = 0.1$ , one obtains  $\hat{r}_i - 2.5 \leq r(a_i) \leq \hat{r}_i + 2.5$ . 154 Such bound is uninformative, because the lower bound for the true reward mean is less than the 155 reward value of the worst arm. The performance bound for LCB confirms this intuition, because 156 Theorem 2.1 requires at least  $N(a_i) \ge [8 * \log(1/0.05)] = 24$  samples in each arm to provide any 157 guarantee with probability at least 0.9 (here  $C^* = d$ ). 158

159 **Can stochastic policies help?** At a first glance, extracting a good decision-making strategy for the problem discussed in Section 2 seems like a hopeless endeavor, because it is information-theoretically 160 impossible to reliably estimate the expected payoff of any of the arms with just a single sample on 161 each. In order to proceed, the key idea is to enlarge the search space to contain stochastic policies.

162 **Definition 2.2** (Stochastic Policies). A stochastic policy over a MAB is a probability distribution 163  $w \in \mathbb{R}^d, w_i \ge 0, \sum_{i=1}^d w_i = 1.$ 

165 To exemplify how stochastic policies can help, consider the *behavioral cloning* policy, which mimics the policy that generated the dataset for the offline MAB in Section 2. Such policy is stochastic, and 166 it plays all arms uniformly at random, thereby achieving a score around 0.5 with high probability. 167 The value of the behavioral cloning policy can be readily estimated using the Hoeffding bound (e.g., 168 Proposition 2.5 in (Wainwright, 2019)): with probability at least  $1 - \delta = 0.9$ , (here d = 10000 is the number of arms and  $\sigma = 1/2$  is the true standard deviation), the value of behavioral cloning 170 policy is greater or equal than  $1/2 - \sqrt{2\sigma^2 \log(2/\delta)}/d \approx 0.488$ . Such value is higher than the one 171 guaranteed for LCB by Theorem 2.1. Intuitively, a stochastic policy that selects multiple arms can 172 be evaluated more accurately because it averages the rewards experienced over different arms. This 173 consideration suggests optimizing over stochastic policies. 174

By optimizing a lower bound on the performance of the stochastic policies, it should be possible to find one with a provably high return. Such an idea leads to solving an offline *linear bandit* problem, as follows

$$\max_{\substack{\in \mathbb{R}^d, w_i \ge 0, \sum_{i=1}^d w_i = 1}} \sum_{i=1}^d w_i \widehat{r}_i - c(w) \tag{4}$$

where c(w) is a suitable confidence interval for the policy w and  $\hat{r}_i$  is the empirical reward for the *i*-th arm defined in equation 1. While this approach is appealing, enlarging the search space to include all stochastic policies brings an increase in the metric entropy of the function class, and concretely, a  $\sqrt{d}$ factor (Abbasi-Yadkori et al., 2011; Rusmevichientong & Tsitsiklis, 2010; Hazan & Karnin, 2016; Jun et al., 2017; Kim et al., 2022) in the confidence intervals c(w) (in equation 4), which negates all gains that arise from considering stochastic policies. In the next section, we propose an algorithm that bypasses the need for such  $\sqrt{d}$  factor by relying on a more careful analysis and optimization procedure.

## 3 TRUST REGION OF UNCERTAINTY FOR STOCHASTIC POLICY ENHANCEMENT (TRUST)

w

In this section, we introduce our algorithm, called Trust Region of Uncertainty for Stochastic policy enhancemenT (TRUST). At a high level, the algorithm is a policy optimization algorithm based on a trust region centered around a reference policy. The size of the trust region determines the degree of pessimism, and its optimal problem-dependent size can be determined by analyzing the supremum of a problem-dependent empirical process. In the sequel, we describe 1) the decision variables, 2) the trust region optimization program, and 3) some techniques for its practical implementation.

#### 3.1 DECISION VARIABLES

The algorithm searches over the class of stochastic policies given by the weight vector  $w = (w_1, w_2, ..., w_d)^{\top}$  of Definition 2.2. Instead of directly optimizing over the weights of the stochastic policy, it is convenient to center w around a *reference stochastic policy*  $\hat{\mu}$  which is either known to perform well or is easy to estimate. In our theory and experiments, we consider a simple setup and use the behavioral cloning policy weighted by the noise levels  $\{\sigma_i\}$  if they are known. Namely, we consider

178 179

189

190

191 192

199

200

$$\widehat{\mu}_i = \frac{N_i / \sigma_i^2}{\sum_{j=1}^d N_j / \sigma_j^2} \quad \forall i \in [d].$$
(5)

209 When the size of the noise  $\sigma_i$  is constant across all arms, the policy  $\hat{\mu}$  is the behavioral cloning 210 policy; when  $\sigma_i$  differs across arms,  $\hat{\mu}$  minimizes the variance of the empirical reward  $\hat{\mu}$  = 211 arg min<sub> $w \in \mathbb{R}^d, w_i \ge 0, \sum_i w_i = 1$  Var  $(w^\top \cdot \hat{r})$ , where  $\hat{r} = (\hat{r}_1, ..., \hat{r}_d)^\top$  is defined in equation 1. Us-212 ing such definition, we define as *decision variable* the *policy improvement* vector  $\Delta := w - \hat{\mu}$ . This 213 preparatory step is key: it allows us to implement **relative pessimism**, namely pessimism on the 214 improvement—represented by  $\Delta$ —rather than on the absolute value of the policy w. Moreover, by 215 restricting the search space to a ball around  $\hat{\mu}$ , one can efficiently reduce the metric entropy of the 216 policy class and obtain tighter confidence intervals.</sub>

# 216 3.2 TRUST REGION OPTIMIZATION

218 Trust region. TRUST (Algorithm 1) returns the stochastic policy 219  $\pi_{TRUST} = \widehat{\Delta} + \widehat{\mu} \in \mathbb{R}^d$ , where  $\widehat{\mu}$  is 220 the reference policy defined in equa-221 222 tion 5 and  $\Delta$  is the policy improvement vector. In order to accurately 224 quantify the effect of the improvement vector  $\Delta$ , we constrain it to a trust re-225 gion C ( $\varepsilon$ ) centered around  $\hat{\mu}$  where 226  $\varepsilon > 0$  is the radius of the trust region. 227 More concretely, for a given radius 228  $\varepsilon > 0$ , the trust region is defined as 229

$$C(\varepsilon) := \left\{ \Delta : \Delta_i + \widehat{\mu}_i \ge 0, \\ \|\Delta + \widehat{\mu}\|_1 = 1, \\ \sum_{i=1}^d \Delta_i^2 \frac{\sigma_i^2}{N_i} \le \varepsilon^2 \right\}.$$
(6)

236The trust region above serves two<br/>purposes: it ensures that the policy<br/> $\widehat{\Delta} + \widehat{\mu}$  still represents a valid stochas-<br/>tic policy, and it regularizes the policy<br/>around the reference policy  $\widehat{\mu}$ . We<br/>then search for the best policy within<br/>C ( $\varepsilon$ ) by solving the optimization programFig.<br/>Fig.



Figure 1: A simple diagram for the trust regions on a 3-dim simplex. The central point is the reference (stochastic) policy, while red ellipses are trust regions around this reference policy.

$$\widehat{\Delta}_{\varepsilon} := \underset{\Delta \in \mathsf{C}(\varepsilon)}{\operatorname{arg\,max}} \Delta^{\top} \widehat{r}. \tag{7}$$

Computationally, the program equation 7 is a second-order cone program (Alizadeh & Goldfarb, 2003; Boyd & Vandenberghe, 2004), which can be solved efficiently with standard off-the shelf libraries (Diamond & Boyd, 2016).

When  $\varepsilon = 0$ , the trust region only includes the vector  $\Delta = 0$ , and the reference policy  $\hat{\mu}$  is the only feasible solution. When  $\varepsilon \to \infty$ , the search space includes all stochastic policies. In this latter case, the solution identified by the algorithm coincides with the greedy algorithm which chooses the arm with the highest empirical return. Rather than leaving  $\varepsilon$  as a hyper-parameter, in the following we highlight a selection strategy for  $\varepsilon$  based on localized Gaussian complexities.

**Critical radius.** The choice of  $\varepsilon$  is crucial to the performance of our algorithm because it balances optimization with regularization. Such consideration suggests that there is an optimal choice for the radius  $\varepsilon$  which balances searching over a larger space with keeping the metric entropy of such space under control. The optimal problem-dependent choice  $\hat{\varepsilon}_*$  can be found as a solution of a certain equation involving a problem-dependent supremum of an empirical process. More concretely, let *E* be the feasible set of  $\varepsilon$  (e.g.,  $E = \mathbb{R}^+$ ). We define the critical radius as

**Definition 3.1** (Critical Radius). The critical radius  $\hat{\varepsilon}_*$  of the trust region is the solution to the program

261 262

230 231

232

233

234

235

243

244

$$\widehat{\varepsilon}_{*} = \operatorname*{arg\,max}_{\varepsilon \in E} \left[ \widehat{\Delta}_{\varepsilon}^{\top} \cdot \widehat{r} - \mathcal{G}(\varepsilon) \right].$$
(8)

264

Such equation involves a quantile of the localized gaussian complexity  $\mathcal{G}(\varepsilon)$  of the stochastic policies identified by the trust region. Mathematically, this is defined as

**267 Definition 3.2** (Quantile of the supremum of Gaussian process). We denote the noise vector as  $\eta = \hat{r} - r$ , which by our assumption is coordinate-wise independent and satisfies  $\eta_i \sim \mathcal{N}\left(0, \sigma_i^2/N(a_i)\right)$ . **269** We define  $\mathcal{G}(\varepsilon)$  as the smallest quantity such that with probability at least  $1 - \delta$ , for any  $\varepsilon \in E$ , it holds that  $\sup_{\Delta \in C(\varepsilon)} \Delta^{\top} \eta \leq \mathcal{G}(\varepsilon)$ . 270 In essence,  $\mathcal{G}(\varepsilon)$  is an upper quantile of the supremum of the Gaussian process  $\sup_{\Delta \in \mathsf{C}(\varepsilon)} \Delta^{\perp} \eta$ 271 which holds uniformly for every  $\varepsilon \in E$ . We also remark that this quantity depends on the feasible set 272 E and the trust region  $C(\varepsilon)$ , and hence, is highly problem-dependent. 273

The critical radius plays a crucial role: it is the radius of the trust region that *optimally* balances 274 optimization with uncertainty. Enlarging  $\varepsilon$  enlarges the search space for  $\Delta$ , enabling the discovery of 275 policies with potentially higher return. However, this also brings an increase in the metric entropy 276 of the policy class encoded by  $\mathcal{G}(\varepsilon)$ , which means that each policy can be estimated less accurately. 277 The critical radius represents the optimal tradeoff between these two forces. The final improvement 278 vector that TRUST returns, which we denote as  $\Delta_*$ , is determined by solving equation 7 with the 279 critical radius  $\hat{\varepsilon}_*$  defined in equation 8. In mathematical terms, we express this as

281

285

286 287

289

290

291

292 293

297

299

301

303

305

307

311

Implementation details. Since it can be difficult to solve equation 8 for a continuous value of  $\varepsilon \in E = \mathbb{R}^+$ , we use a discretization argument by considering the following candidate subset:

 $\widehat{\Delta}_* := \underset{\Delta \in \mathsf{C}(\widehat{\varepsilon}_*)}{\arg\max} \Delta^\top \widehat{r}.$ 

$$E = \left\{ \varepsilon_0, \frac{\varepsilon_0}{\alpha}, ..., \frac{\varepsilon_0}{\alpha^{|E|-1}} \right\},\tag{10}$$

(9)

where  $\alpha > 1$  is the decaying rate and  $\varepsilon_0$  is the largest possible radius, which is the maximal weighted distance from the reference policy to any vertex. Mathematically, this is defined as  $\varepsilon_0 = \max_{i \in [d]} \sqrt{\sum_{j \neq i} \hat{\mu}_j^2 \sigma_j^2 / N_j + (1 - \hat{\mu}_i)^2 \sigma_i^2 / N_i}$ . Our analysis that leads to Theorem 4.1 takes into account such discretization argument.

Ì

In line 2 of Algorithm 1, the algorithm works by estimating the quantile of the supremum of the 295 localized Gaussian complexity  $\mathcal{G}(\varepsilon)$  that appears in Definition 3.2, and then choose the  $\varepsilon$  that 296 maximizes the objective function in equation 8. Although  $\mathcal{G}(\varepsilon)$  can be upper bounded analytically, in our experiments we aim to obtain tighter guarantees and so we estimate it via Monte-Carlo. This can 298 be achieved by 1) sampling independent noise vectors  $\eta$ , 2) solving  $\sup_{\Delta \in C(\varepsilon)} \Delta^+ \eta$  and 3) estimating the quantile via order statistics. More details can be found in Appendix D.

300 In summary, our practical algorithm can be seen as solving the optimization problem 302

$$(\widehat{\varepsilon}_*, \widehat{\Delta}_*) = \underset{\varepsilon \in E, \Delta \in \mathsf{C}(\varepsilon)}{\operatorname{arg\,max}} \left\{ \Delta^\top \widehat{r} - \widehat{\mathcal{G}}(\varepsilon) \right\}$$

where  $\widehat{r} \in \mathbb{R}^d$  is the empirical reward vec-306 tor with  $\hat{r}_i$  defined in equation 1. Here,  $\mathcal{G}(\varepsilon)$  is computed according to the Monte-308 Carlo method defined in Algorithm 2 in Ap-309 pendix D and E is the candidate subset for 310 radius defined in equation 10. This indicates a balance between the empirical reward of a 312 stochastic policy and the local entropy metric 313 it induces. 314

Algorithm 1 Trust Region of Uncertainty for Stochastic policy enhancemenT (TRUST) **Input:** Offline dataset  $\mathcal{D}$ , failure probability  $\delta$ , the

candidate set for the trust region widths E (in practice, this is chosen as equation 10).

1. For  $\varepsilon \in E$ , compute  $\Delta_{\varepsilon}$  from equation 7.

2. For  $\varepsilon \in E$ , estimate  $\mathcal{G}(\varepsilon)$  via Monte-Carlo method (see Algorithm 2 in Appendix D).

3. Solve equation 8 to obtain the critical radius  $\hat{\varepsilon}_*$ .

4. Compute the optimal improvement vector in  $C(\widehat{\varepsilon}_*)$  via equation 9, denoted as  $\widehat{\Delta}_*$ .

5. Return the stochastic policy  $\pi_{TRUST} = \hat{\mu} + \Delta_*$ .

315 316

317

#### THEORETICAL GUARANTEES 4

**Problem-dependent analysis** In this section, we provide some theoretical guarantees for the policy 318  $\pi_{TRUST}$  returned by TRUST. We present 1) an improvement over the reference policy  $\hat{\mu}$ , 2) a 319 sub-optimality gap with respect to any comparator policy  $\pi$  and 3) an actionable lower bound on 320 the performance of the output policy. Given a stochastic policy  $\pi$ , we let  $V^{\pi} = \mathbb{E}_{a \sim \pi}[r(a)]$  denote 321 its value function. Furthermore, we denote a comparator policy  $\pi$  by a triple  $(\varepsilon, \Delta, \pi)$  such that 322  $\varepsilon > 0, \Delta \in \mathsf{C}(\varepsilon), \pi = \widehat{\mu} + \Delta.$ 323

Theorem 4.1 (Main theorem). TRUST has the following properties.

*1.* With probability at least  $1 - \delta$ , the improvement over the behavioral policy is at least

$$V^{\pi_{TRUST}} - V^{\hat{\mu}} \ge \sup_{\varepsilon \le \varepsilon_0, \Delta \in \mathsf{C}(\varepsilon)} \left[ \Delta^\top r - 2\mathcal{G}\left( \lceil \varepsilon \rceil \right) \right], \quad \text{where} \quad \lceil \varepsilon \rceil = \inf\{\varepsilon' \in E, \varepsilon' \ge \varepsilon\}.$$
(11)

327 328

330

331 332 333

339

324

325 326

2. With probability at least  $1 - \delta$ , for any stochastic comparator policy  $(\varepsilon, \Delta, \pi)$ , the sub-optimality of the output policy can be upper bounded as

$$V^{\pi} - V^{\pi_{TRUST}} \le 2\mathcal{G}\left(\left[\varepsilon\right]\right). \tag{12}$$

3. With probability at least  $1 - 2\delta$ , the data-dependent lower bound on  $V^{\pi_{TRUST}}$  satisfies

$$V^{\pi_{TRUST}} \ge \pi_{TRUST}^{\top} \widehat{r} - \mathcal{G}\left(\left\lceil \widehat{\varepsilon}_* \right\rceil\right) - \sqrt{\frac{2\log(1/\delta)}{\sum_{j=1}^d N_j / \sigma_j^2}},\tag{13}$$

where  $\pi_{TBUST} = \hat{\mu} + \hat{\Delta}_*$  is the policy output by Algorithm 1.

340 The proof of Theorem 4.1 is deferred to Appendix B. A fine-grained analysis for the suboptimality 341 is contained in Appendix E. Our guarantees are problem-dependent as a function of the Gaussian 342 process  $\mathcal{G}(\cdot)$ ; in Section 5 we show how these can be instantiated on an actual problem, highlighting 343 the tightness of the analysis. Equation (11) highlights the improvement with respect to the behavioral 344 policy. It is expressed as a trade-off between maximizing the improvement  $\Delta^{+}r$  and minimizing its 345 uncertainty  $\mathcal{G}([\varepsilon])$ . The presence of the  $\sup_{\varepsilon}$  indicates that TRUST achieves an *optimal* balance between these two factors. The state of the art guarantees that we are aware of highlight a trade-off 346 between value and variance (Jin et al., 2021; Min et al., 2021). The novelty of our result lies in the 347 fact that TRUST optimally balances the uncertainty implicitly as a function of the 'coverage' as well 348 as the metric entropy of the search space. That is, TRUST selects the most appropriate search space 349 by trading off its metric entropy with the quality of the policies that it contains. The right-hand side in 350 Equation (13) gives actionable statistical guarantees on the quality of the final policy and it can be fully 351 computed from the available dataset; we give an example of the tightness of the analysis in Section 5. 352

353 354

# **Localized Gaussian complexity** $\mathcal{G}(\varepsilon)$ . In Theorem 4.1, we up-

355 per bound the suboptimality  $V^{\pi}$  – 356  $V^{\pi_{TRUST}}$  via a notion of local-357 ized metric entropy  $\mathcal{G}(\cdot)$ . It is 358 the quantile of the supremum of 359 a Gaussian process, which can 360 hardly be calculated analytically 361 but can be efficiently estimated 362 via Monte Carlo method (which does not collect additional samples, e.g., see Appendix D). It can 364 also be concentrated around its expectation, which is also localized 366 Gaussian width, a concept well-367 established in statistical learning 368 theory (Bellec, 2019; Wei et al., 369 2020; Wainwright, 2019). More 370 concretely, this is the localized 371 Gaussian width for an affine sim-



Figure 2: The upper bound for the localized Gaussian width over a shifted simplex on d = 10000 dimension. The shifted simplex is  $\{\Delta \in \mathbb{R}^d : \sum_{i=1}^d \Delta_i = 0\}$ . The two-staged upper bound is based on Theorem 1 in (Bellec, 2019)

372 plex:  $\mathbb{E}[\sup_{\Delta \in C(\varepsilon)} \Delta^{\top} \eta] = \mathbb{E}[\sup_{\mathbb{S}^{d-1} \cap \{\Delta : \|\Delta\|_{\Sigma} \le \varepsilon\}} \Delta^{\top} \eta]$ , where  $\mathbb{S}^{d-1}$  denotes the simplex in  $\mathbb{R}^d$ 373 and  $\Sigma := \operatorname{diag}\left(\frac{\sigma_1^2}{N_1}, \frac{\sigma_2^2}{N_2}, ..., \frac{\sigma_d^2}{N_d}\right)$  is the weighted matrix. Moreover, this localized Gaussian width 375 can be upper bound via

$$\mathbb{E}\left[\sup_{\Delta\in\mathsf{C}(\varepsilon)}\Delta^{\top}\eta\right]\lesssim\min\left\{\sqrt{\log\left(d\varepsilon^{2}\right)},\varepsilon\sqrt{d}\right\}.$$
(14)

378 To make it clearer, we plot this upper bound for localized Gaussian width in Figure 2. In equation 14, 379 the rate matches the minimax lower bound up to universal constant (Gordon et al., 2007; Lecué & 380 Mendelson, 2013; Bellec, 2019). To see the implication of the upper bound equation 14, let's consider 381 a simple example where the logging policy is uniform over all arms. We denote the optimal arm as 382  $a^*$  and define  $C^* := 1/\mu(a^*)$  as the concentrability coefficient. By applying equation 14 and some concentration techniques (see Wainwright, 2019), we can perform a fine-grained analysis for the suboptimality induced by  $\pi_{TRUST}$ . Specifically, with probability at least  $1 - \delta$ , one has 384

$$V^{\pi_*} - V^{\pi_{TRUST}} \lesssim \sqrt{C^* \log(2d|E|/\delta)/N}.$$
(15)

Note that, the high-probability upper bound here is minimax optimal up to constant and logarithmic 388 factor (Rashidinejad et al., 2021) when  $C^* \ge 2$ . Moreover, this example of uniform logging policy is 389 an instance where LCB achieves minimax sub-optimality (up to constant and log factors) (see the 390 proof of Theorem 2 in Rashidinejad et al., 2021). In this case, TRUST will achieve the same level of 391 guarantees for the suboptimality of the output policy. We also empirically show the effectiveness of 392 TRUST in Section 5. The full theorem for a fine-grained analysis for the suboptimality and its proof 393 are deferred to Appendix E. 394

Augmentation with LCB. Compared to classical LCB, Algorithm 1 considers a much larger searching 395 space, which encompasses not only the vertices of the simplex but the inner points as well. This 396 enlargement of searching space shows great advantage, but this also comes with the price of larger 397 uncertainty, especially when the width  $\varepsilon$  is large. In LCB, one considers the uncertainty by upper 398 bound the noise at each vertex uniformly, while in our case, the uniform upper bound for a sub-region 399 of the shifted simplex must be considered. When  $\varepsilon$  is large, the trust region method will induce 400 larger uncertainty and tend to select a more stochastic policy than LCB and hence, can achieve worse 401 performance. Moreover, when each arm has sufficiently many data samples to roughly estimate its 402 mean return to reasonable accuracy, LCB works well because it chooses the arm with a tight lower bound. However the current results for LCB do not cover the important case where only few samples 403 (e.g., less than 24 as described in Section 2) are available. Encouragingly our work shows strong 404 results in such settings. To determine the most effective final policy, one can always combine TRUST 405 (Algorithm 1) with LCB and select the better one between them based on the lower bound induced by 406 two algorithms. By comparing the lower bounds of LCB and TRUST, the value of the finally output 407 policy is guaranteed to outperform the lower bound for either LCB or TRUST with high probability. 408 We defer the detailed algorithm and its theoretical guarantees to Appendix G. 409

410 411

385 386 387

#### 5 **EXPERIMENTS**

412 413 414

415

417

We present simulated experiments where we show the failure of LCB and the strong performance of TRUST. Moreover, we also present an application of TRUST to offline reinforcement learning.

Simulated experiments: A data-starved MAB. We consider a data-starved MAB problem with 416 d = 10000 arms denoted by  $a_i, i \in [d]$ . The reward distributions are

418 419 420

$$r(a_i) \sim \text{Uniform}(0.5, 1.5) \text{ for } i \le 5000; \quad r(a_i) \sim \mathcal{N}(0, 1/4) \text{ for } i > 5000.$$
 (16)

Namely, the set of good arms have reward random variables from a uniform distribution over [0.5, 1.5]421 with unit mean, while the bad arms return a Gaussian reward with zero mean. We consider a dataset 422 that contains a single sample for each of these arms. 423

We test Algorithm 1 on this MAB instance with fixed variance level  $\sigma_i = 1/2$ . We set the reference 424 policy  $\hat{\mu}$  to be the behavioral cloning policy, which coincides with the uniform policy. We also test 425 LCB and the greedy method which simply chooses the policy with the highest empirical reward. 426

427 In this example, the greedy algorithm fails because it erroneously selects an arm with a reward 428 > 1.5, but such reward can only originate from a normal distribution with mean zero. Despite LCB 429 incorporates the principle of pessimism under uncertainty, it selects an arm with average return equal to zero; its performance lower bound given by the confidence intervals is -1.5, which is almost 430 vacuous and very uninformative. The behavioral cloning policy performs better, because it selects an 431 arm uniformly at random, achieving the score 0.5.

LCB

0

Greedy

0

/1	-2	- 1
	· • •	
	~	-
	- 4	e
	' <b>L</b>	

Behavior

Policy

0.5

Л	2	Л
7	0	-
	~	_
21		h

436

437

438

439 440 441

442

443

444

Table 1: Results of simulated experiments in a 10000-arm bandit. The reward distribution is described in equation 16. The offline dataset includes one sample for each arm. The greedy method chooses the arm with the highest empirical reward. LCR selects an arm based on equation 3. The lower bound for

Improvement

by TRUST

0.42

TRUST

0.92

TRUST

Lower Bound

0.6

LCB Lower

Bound

-1.5

arm with the highest empirical reward. LCB selects an arm based on equation 3. The lower bound for LCB and TRUST follow equation 2 and equation 13, respectively.

Algorithm 1 achieves the best performance: the value of the policy that it identifies is 0.92, which *almost matches the optimal policy*. The lower bound on its performance computed by instantiating the RHS in equation 13 is around 0.6, a guarantee much tighter than that for LCB.

In order to gain intuition on the learning mechanics of TRUST, in Figure 3 we progressively enlarge the radius of the trust region from zero to the largest possible radius (on the x axis) and plot the value of the policy that maximizes the linear objective  $\Delta^{T} \hat{r}$ ,  $\Delta \in C(\varepsilon)$  for each value of the radius  $\varepsilon$ . Note that we rescale the range of  $\varepsilon$  to make the largest possible  $\varepsilon$  be one. In the same figure we also plot the lower bound computed with the help of equation equation 13.

Initially, the value of the policy in-450 creases because the optimization in 451 equation 7 is performed over a larger 452 set of stochastic policies. However, 453 when  $\varepsilon$  approaches the maximal pos-454 sible radius, all stochastic policies are 455 included in the optimization program. 456 In this case, TRUST greedily selects 457 the arm with the highest empirical re-458 ward, which is from a normal distri-459 bution with a mean zero. The optimal balance between the size of the 460 policy search space and its metric en-461 tropy is given by the critical radius 462  $\varepsilon = 0.0116\varepsilon_0$ , which is the point 463 where the lower bound is the highest. 464

A more general data-starved MAB.
Besides the data-starved MAB we constructed, we also show that in general MABs, the performance of TRUST is on par with LCB, but



Figure 3: Policy values and their lower bounds for a datastarved MAB instance with 10000 arms whose reward distribution is described in equation 16.

470 TRUST will have a much tighter statistical guarantee, i.e., a larger lower bound for the value of the returned policy. We did experiments on a d = 1000-arm MAB where the reward distribution is 471  $r(a_i) \sim \mathcal{N}(i/1000, 1/4), \quad \forall i \in [d]$ . We ran TRUST Algorithm 1 and LCB over 8 different random 472 seeds. When we have a single sample for each arm, TRUST will get a similar score as LCB. However, 473 TRUST give a much tighter statistical guarantee than LCB, in the sense that the lower bound output by 474 TRUST is much higher than that output by LCB so that TRUST can output a policy that is guaranteed 475 to achieved a higher value. Moreover, we found the policies output from TRUST are much more 476 stable than those from LCB. In all runs, while the lowest value of the arm chosen by LCB is around 477 0.24, all policies returned by TRUST have values above 0.65 with a much smaller variance, as shown 478 in Table 2. 479

Offline reinforcement learning. In this section, we apply Algorithm 1 to the offline reinforcement learning (RL) setting under the assumption that the logging policies which generated the dataset are accessible. To be clear, our goal is not to exceed the performance of the state of the art deep RL algorithms—our algorithm is designed for bandit problems—but rather to illustrate the usefulness of our algorithm and theory.

485 Since our algorithm is designed for bandit problems, in order to apply it to the sequential setting, we map MDPs to MABs. Each policy in the MDP maps to an action in the MAB, and each trajectory

return in the MDP maps to an experienced return in the MAB setting. Notice that this reduction
disregards the sequential aspect of the problem and thus our algorithm cannot perform 'trajectory
stitching' (Levine et al., 2020; Kumar et al., 2020; Kostrikov et al., 2021). Furthermore, it can only
be applied under the assumption that the logging policies are known.

490 Specifically we consider a setting 491 where there are multiple known log-492 ging policies, each generating few tra-493 jectories. We test Algorithm 1 on 494 some selected environments from the 495 D4RL dataset (Fu et al., 2020) and 496 compare its performance to the (CQL) algorithm (Kumar et al., 2020), a 497 popular and strong baseline for of-498 fline RL algorithms. Since the D4RL 499 dataset does not directly include the 500 logging policies, we generate new 501 datasets by running Soft Actor Critic 502 (SAC) (Haarnoja et al., 2018) for 1000

	LCB	TRUST
mean reward	0.718	0.725
mean lower bound	0.156	0.544
variance	0.265	0.038
minimal reward	0.239	0.658

Table 2: Comparison between LCB and TRUST (Algorithm 1) on a data-starved MAB with 1000 arms whose reward distribution follows  $r(a_i) \sim \mathcal{N}(i/1000, 1/4)$ . Both methods are repeated on 8 random seeds.

episodes. We store 100 intermediate policies generated by SAC, and roll out 1 trajectory from each policy.

We use some default hyper-parameters for CQL.<sup>1</sup> We report the unnormalized scores in Table 3, each averaged over 4 random seeds. Algorithm 1 achieves a score on par with or higher than that of CQL, especially when the offline dataset is of poor quality and when there are very few—or just one—trajectory generated from each logging policy. Notice that while CQL is not guaranteed to outperform the behavioral policy, TRUST is backed by Theorem 4.1.

- 511 Additionally, while CQL took around
- 16-24 hours on one NVIDIA GeForce
  RTX 2080 Ti, TRUST only took 0.5-1
  hours on 10 CPUs. The experimental details are included in Appendix H.
  Moreover, while the performance of
  CQL is highly reliant on the choice of hyper-parameters, TRUST is essentially hyper-parameters free.

519 520 521

6 CONCLUSION

522
523 In this paper we make a substantial contribution towards sample efficient decision making, by designing a data-efficient policy optimization algorithm that leverages offline data for the MAB setting. The key intuition of

CQL TRUST 499 999 1-traj-low Hopper 2606 3437 1-traj-high 748 1-traj-low 763 Ant 1-traj-high 4115 4488 1-traj-low 311 346 Walker2d 4093 4097 1-traj-high 1-traj-low 5775 5473 HalfCheetah 1-traj-high 9067 10380

Table 3: Unnormalized score of CQL and TRUST in 4 environments from D4RL. In 1-traj-low case, we take the first 100 policies in the running of SAC. In 1-traj-high case, we take the (10x + 1)-th policy for  $x \in [100]$ . We sample one trajectory from each policy we take in all experiments.

this work is to search over stochastic policies, which can be estimated more easily than deterministic
ones. The design of our algorithm is enabled by a number of key insights, such as the use of the localized gaussian complexity which leads to the definition of the critical radius for the trust region. We
believe that these concepts can be used more broadly to help design truly sample efficient algorithms,
which can in turn enable the application of decision making to new settings where a high sample
efficiency is critical.

## References

Yasin Abbasi-Yadkori, David Pal, and Csaba Szepesvari. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems (NIPS)*, 2011.

534 535

536

<sup>538</sup> 539

<sup>&</sup>lt;sup>1</sup>We use the codebase and default hyper-parameters in https://github.com/young-geng/CQL.

549

565

576

580

581

582

583

- Farid Alizadeh and Donald Goldfarb. Second-order cone programming. *Mathematical programming*, 95(1):3–51, 2003.
- Mawulolo K Ameko, Miranda L Beltzer, Lihua Cai, Mehdi Boukhechba, Bethany A Teachman, and Laura E Barnes. Offline contextual multi-armed bandits for mobile health interventions: A case study on emotion regulation. In *Proceedings of the 14th ACM Conference on Recommender Systems*, pp. 249–258, 2020.
- Jean-Yves Audibert, Sébastien Bubeck, et al. Minimax policies for adversarial and stochastic bandits.
   In *COLT*, volume 7, pp. 1–122, 2009.
- Jean-Yves Audibert, Sébastien Bubeck, and Rémi Munos. Best arm identification in multi-armed bandits. In *COLT*, pp. 41–53, 2010.
- Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov):397–422, 2002.
- Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit
   problem. *Machine learning*, 47:235–256, 2002.
- Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for rein forcement learning. In *International Conference on Machine Learning*, pp. 263–272. PMLR, 2017.
- Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- Peter L Bartlett, Olivier Bousquet, and Shahar Mendelson. Local rademacher complexities. *The Annals of Statistics*, 33(4):1497–1537, 2005.
- Jonathan Bassen, Bharathan Balaji, Michael Schaarschmidt, Candace Thille, Jay Painter, Dawn Zimmaro, Alex Games, Ethan Fast, and John C Mitchell. Reinforcement learning for the adaptive scheduling of educational activities. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pp. 1–12, 2020.
- Pierre C Bellec. Localized gaussian width of m-convex hulls with applications to lasso and convex aggregation. 2019.
- 572 573 Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- Sébastien Bubeck, Nicolo Cesa-Bianchi, et al. Regret analysis of stochastic and nonstochastic
   multi-armed bandit problems. *Foundations and Trends*® *in Machine Learning*, 5(1):1–122, 2012.
- Wei Chen, Yajun Wang, Yang Yuan, and Qinshi Wang. Combinatorial multi-armed bandit and its
  extension to probabilistically triggered arms. *The Journal of Machine Learning Research*, 17(1): 1746–1778, 2016.
  - Ching-An Cheng, Tengyang Xie, Nan Jiang, and Alekh Agarwal. Adversarially trained actor critic for offline reinforcement learning. *arXiv preprint arXiv:2202.02446*, 2022.
  - Richard Combes, Mohammad Sadegh Talebi Mazraeh Shahi, Alexandre Proutiere, et al. Combinatorial bandits revisited. *Advances in neural information processing systems*, 28, 2015.
- Yan Dai, Ruosong Wang, and Simon S Du. Variance-aware sparse linear bandits. *arXiv preprint* arXiv:2205.13450, 2022.
- Rémy Degenne and Vianney Perchet. Anytime optimal algorithms in stochastic multi-armed bandits.
   In *International Conference on Machine Learning*, pp. 1587–1595. PMLR, 2016.
- Steven Diamond and Stephen Boyd. Cvxpy: A python-embedded modeling language for convex optimization. *The Journal of Machine Learning Research*, 17(1):2909–2913, 2016.
- 593 Yaqi Duan and Martin J Wainwright. Policy evaluation from a single path: Multi-step methods, mixing and mis-specification. *arXiv preprint arXiv:2211.03899*, 2022.

594 595	Yaqi Duan and Martin J Wainwright. A finite-sample analysis of multi-step temporal difference estimates. In <i>Learning for Dynamics and Control Conference</i> , pp. 612–624. PMLR, 2023.
598 597 598	Yaqi Duan, Mengdi Wang, and Martin J Wainwright. Optimal policy evaluation using kernel-based temporal difference methods. <i>arXiv preprint arXiv:2109.12002</i> , 2021.
599 600	Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven reinforcement learning. arXiv preprint arXiv:2004.07219, 2020.
601 602 603	Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. In <i>International conference on machine learning</i> , pp. 2052–2062. PMLR, 2019.
604 605	Aurélien Garivier and Emilie Kaufmann. Optimal best arm identification with fixed confidence. In <i>Conference on Learning Theory</i> , pp. 998–1027. PMLR, 2016.
606 607	Sara A Geer. Empirical Processes in M-estimation, volume 6. Cambridge university press, 2000.
608 609 610	Yehoram Gordon, Alexander E Litvak, Shahar Mendelson, and Alain Pajor. Gaussian averages of interpolated bodies and applications to approximate reconstruction. <i>Journal of Approximation Theory</i> , 149(1):59–73, 2007.
611 612 613 614	Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In <i>International conference on machine learning</i> , pp. 1861–1870. PMLR, 2018.
615 616	Botao Hao, Xiang Ji, Yaqi Duan, Hao Lu, Csaba Szepesvári, and Mengdi Wang. Bootstrapping statistical inference for off-policy evaluation. <i>arXiv preprint arXiv:2102.03607</i> , 2021.
617 618 619	Elad Hazan and Zohar Karnin. Volumetric spanners: an efficient exploration basis for learning. <i>Journal of Machine Learning Research</i> , 2016.
620 621	Todd Hester and Peter Stone. Texplore: real-time sample-efficient reinforcement learning for robots. <i>Machine learning</i> , 90:385–429, 2013.
622 623 624	Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl? <i>arXiv</i> preprint arXiv:2012.15085, 2020.
625 626	Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl? In <i>International Conference on Machine Learning</i> , pp. 5084–5096. PMLR, 2021.
627 628 629 630	Kwang-Sung Jun, Aniruddha Bhargava, Robert Nowak, and Rebecca Willett. Scalable generalized linear bandits: Online computation and hashing. <i>Advances in Neural Information Processing Systems</i> , 30, 2017.
631 632 633	Yeoneung Kim, Insoon Yang, and Kwang-Sung Jun. Improved regret analysis for variance-adaptive linear bandits and horizon-free linear mixture mdps. <i>Advances in Neural Information Processing Systems</i> , 35:1060–1072, 2022.
634 635	Vladimir Koltchinskii. Rademacher penalties and structural risk minimization. <i>IEEE Transactions</i> on Information Theory, 47(5):1902–1914, 2001.
637 638	Vladimir Koltchinskii. Local rademacher complexities and oracle inequalities in risk minimization. 2006.
639 640	Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-learning. <i>arXiv preprint arXiv:2110.06169</i> , 2021.
642 643	Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. <i>arXiv preprint arXiv:2006.04779</i> , 2020.
644 645 646	Branislav Kveton, Zheng Wen, Azin Ashkan, and Csaba Szepesvari. Tight regret bounds for stochastic combinatorial semi-bandits. In <i>Artificial Intelligence and Statistics</i> , pp. 535–543. PMLR, 2015.
0.47	Tze Leung Lai Adaptive treatment allocation and the multi-armed bandit problem. The annals of

647 Tze Leung Lai. Adaptive treatment allocation and the multi-armed bandit problem. *The annals of statistics*, pp. 1091–1114, 1987.

648 649	Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. <i>Advances in applied mathematics</i> , 6(1):4–22, 1985.
651 652	John Langford and Tong Zhang. The epoch-greedy algorithm for multi-armed bandits with side information. <i>Advances in neural information processing systems</i> , 20, 2007.
653 654	Tor Lattimore and Csaba Szepesvári. Bandit Algorithms. Cambridge University Press, 2020.
655 656	Guillaume Lecué and Shahar Mendelson. Learning subgaussian classes: Upper and minimax bounds. arXiv preprint arXiv:1305.4825, 2013.
657 658 659	Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. <i>arXiv preprint arXiv:2005.01643</i> , 2020.
660 661	Gen Li, Laixi Shi, Yuxin Chen, Yuejie Chi, and Yuting Wei. Settling the sample complexity of model-based offline reinforcement learning. <i>arXiv preprint arXiv:2204.05275</i> , 2022.
662 663 664	Zihao Li, Zhuoran Yang, and Mengdi Wang. Reinforcement learning with human feedback: Learning dynamic choices via pessimism. <i>arXiv preprint arXiv:2305.18438</i> , 2023.
665 666 667	Rongrong Liu, Florent Nageotte, Philippe Zanne, Michel de Mathelin, and Birgitta Dresp-Langley. Deep reinforcement learning for the control of robotic manipulation: a focussed mini-review. <i>Robotics</i> , 10(1):22, 2021.
668 669 670	Yi Liu and Veronika Ročková. Variable selection via thompson sampling. <i>Journal of the American Statistical Association</i> , 118(541):287–304, 2023.
671 672 673	Ying Liu, Brent Logan, Ning Liu, Zhiyuan Xu, Jian Tang, and Yangzhi Wang. Deep reinforcement learning for dynamic treatment regimes on medical registry data. In 2017 IEEE international conference on healthcare informatics (ICHI), pp. 380–385. IEEE, 2017.
674 675 676	Pierre Ménard, Omar Darwiche Domingues, Anders Jonsson, Emilie Kaufmann, Edouard Leurent, and Michal Valko. Fast active learning for pure exploration in reinforcement learning. In <i>International Conference on Machine Learning</i> , pp. 7599–7608. PMLR, 2021.
678 679 680	Yifei Min, Tianhao Wang, Dongruo Zhou, and Quanquan Gu. Variance-aware off-policy evaluation with linear function approximation. <i>Advances in neural information processing systems</i> , 34: 7598–7610, 2021.
681 682	Wenlong Mou, Martin J Wainwright, and Peter L Bartlett. Off-policy estimation of linear functionals: Non-asymptotic theory for semi-parametric efficiency. <i>arXiv preprint arXiv:2209.13075</i> , 2022.
684 685	Ofir Nachum, Bo Dai, Ilya Kostrikov, Yinlam Chow, Lihong Li, and Dale Schuurmans. Algaedice: Policy gradient from arbitrary experience. <i>arXiv preprint arXiv:1912.02074</i> , 2019.
686 687 688	Preetum Nakkiran, Behnam Neyshabur, and Hanie Sedghi. The deep bootstrap framework: Good online learners are good offline generalizers. <i>arXiv preprint arXiv:2010.08127</i> , 2020.
689 690	Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. <i>arXiv preprint arXiv:1910.00177</i> , 2019.
691 692 693	Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart Russell. Bridging offline rein- forcement learning and imitation learning: A tale of pessimism. <i>arXiv preprint arXiv:2103.12021</i> , 2021.
695 696 697	Sherry Ruan, Allen Nie, William Steenbergen, Jiayu He, JQ Zhang, Meng Guo, Yao Liu, Kyle Dang Nguyen, Catherine Y Wang, Rui Ying, et al. Reinforcement learning tutor better supported lower performers in a math task. <i>arXiv preprint arXiv:2304.04933</i> , 2023.
698 699 700	Paat Rusmevichientong and John N Tsitsiklis. Linearly parameterized bandits. <i>Mathematics of Operations Research</i> , 35(2):395–411, 2010.
700	Daniel Russo. Simple bayesian algorithms for best arm identification. In <i>Conference on Learning Theory</i> , pp. 1417–1418. PMLR, 2016.

702 703 704 705	Nian Si, Fan Zhang, Zhengyuan Zhou, and Jose Blanchet. Distributionally robust policy evaluation and learning in offline contextual bandits. In <i>International Conference on Machine Learning</i> , pp. 8884–8894. PMLR, 2020.
705	Richard S Sutton and Andrew G Barto. Reinforcement learning: An introduction. MIT Press, 2018.
707 708 709	Philip Thomas, Georgios Theocharous, and Mohammad Ghavamzadeh. High confidence policy improvement. In <i>International Conference on Machine Learning</i> , pp. 2380–2388. PMLR, 2015.
710	Roman Vershynin. High-dimensional probability. University of California, Irvine, 2020.
711 712 713	Martin J Wainwright. <i>High-dimensional statistics: A non-asymptotic viewpoint</i> , volume 48. Cambridge university press, 2019.
714 715 716	Kerong Wang, Hanye Zhao, Xufang Luo, Kan Ren, Weinan Zhang, and Dongsheng Li. Bootstrapped transformer for offline reinforcement learning. <i>Advances in Neural Information Processing Systems</i> , 35:34748–34761, 2022.
718 719	Siwei Wang and Wei Chen. Thompson sampling for combinatorial semi-bandits. In <i>International Conference on Machine Learning</i> , pp. 5114–5122. PMLR, 2018.
720 721 722	Yuting Wei, Billy Fang, and Martin J. Wainwright. From gauss to kolmogorov: Localized measures of complexity for ellipses. 2020.
723 724	Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning. arXiv preprint arXiv:1911.11361, 2019.
725 726 727	Tengyang Xie and Nan Jiang. Batch value-function approximation with only realizability. <i>arXiv</i> preprint arXiv:2008.04990, 2020.
728 729 730	Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-consistent pessimism for offline reinforcement learning. <i>arXiv preprint arXiv:2106.06926</i> , 2021.
731 732 733	Wei Xiong, Han Zhong, Chengshuai Shi, Cong Shen, Liwei Wang, and Tong Zhang. Nearly minimax optimal offline reinforcement learning with linear function approximation: Single-agent mdp and markov game. <i>arXiv preprint arXiv:2205.15512</i> , 2022.
734 735 736	Ming Yin and Yu-Xiang Wang. Towards instance-optimal offline reinforcement learning with pessimism. <i>arXiv preprint arXiv:2110.08695</i> , 2021.
737 738 739	Ming Yin, Yaqi Duan, Mengdi Wang, and Yu-Xiang Wang. Near-optimal offline reinforcement learning with linear representation: Leveraging variance information with pessimism. <i>arXiv</i> preprint arXiv:2203.05804, 2022.
740 741 742	Andrea Zanette, Emma Brunskill, and Mykel J. Kochenderfer. Almost horizon-free structure-aware best policy identification with a generative model. In <i>Advances in Neural Information Processing Systems</i> , 2019.
744 745 746	Andrea Zanette, Alessandro Lazaric, Mykel J Kochenderfer, and Emma Brunskill. Provably effi- cient reward-agnostic navigation with linear value iteration. In <i>Advances in Neural Information</i> <i>Processing Systems</i> , 2020.
747 748 749	Andrea Zanette, Martin J Wainwright, and Emma Brunskill. Provable benefits of actor-critic methods for offline reinforcement learning. <i>arXiv preprint arXiv:2108.08812</i> , 2021.
750 751 752	Ruiqi Zhang, Xuezhou Zhang, Chengzhuo Ni, and Mengdi Wang. Off-policy fitted q-evaluation with differentiable function approximators: Z-estimation and inference theory. In <i>International Conference on Machine Learning</i> , pp. 26713–26749. PMLR, 2022.
753 754 755	Zihan Zhang, Jiaqi Yang, Xiangyang Ji, and Simon S Du. Improved variance-aware confidence sets for linear bandits and linear mixture mdp. <i>Advances in Neural Information Processing Systems</i> , 34:4342–4355, 2021.

## A ADDITIONAL RELATED WORK

758 Multi-armed bandit (MAB) is a classical decision-making framework (Lattimore & Szepesvári, 759 2020; Lai & Robbins, 1985; Lai, 1987; Langford & Zhang, 2007; Auer, 2002; Bubeck et al., 2012; 760 Audibert et al., 2009; Degenne & Perchet, 2016). The natural approach in offline MABs is the LCB 761 algorithm (Ameko et al., 2020; Si et al., 2020), an offline variant of the classical UCB method (Auer 762 et al., 2002) which is minimax optimal (Rashidinejad et al., 2021). The optimization over stochastic policies is also considered in combinatorial multi-armed bandits (CMAB) (Combes et al., 2015). 763 Most works on CMAB focus on variants of the UCB algorithm (Kveton et al., 2015; Combes et al., 764 2015; Chen et al., 2016) or of Thompson sampling (Wang & Chen, 2018; Liu & Ročková, 2023), 765 and they are generally online. Our framework can also be applied to offline reinforcement learning 766 (RL) (Sutton & Barto, 2018) whenever the logging policies are accessible. There exist a lot of 767 practical algorithms for offline RL (Fujimoto et al., 2019; Peng et al., 2019; Wu et al., 2019; Kumar 768 et al., 2020; Kostrikov et al., 2021). Theory has also been investigated extensively in tabular domain 769 and function approximation setting (Nachum et al., 2019; Xie & Jiang, 2020; Zanette et al., 2021; 770 Xie et al., 2021; Yin et al., 2022; Xiong et al., 2022). Some works also tried to establish general 771 guarantees for deep RL algorithms via sophisticated statistical tools, such as bootstrapping (Thomas 772 et al., 2015; Nakkiran et al., 2020; Hao et al., 2021; Wang et al., 2022; Zhang et al., 2022).

773 We rely on the notion of pessimism, which is a key concept in offline bandits and RL. While most prior 774 works focused on the so-called absolute pessimism (Jin et al., 2020; Xie et al., 2021; Yin et al., 2022; 775 Rashidinejad et al., 2021; Li et al., 2023), the work of (Cheng et al., 2022) applied pessimism not on 776 the policy value but on the difference (or improvement) between policies. However, their framework 777 is very different from ours. We make extensive use of two key concepts, namely localization laws and 778 critical radii (Wainwright, 2019), which control the relative scale of the signal and uncertainty. The 779 idea of localization plays a critical role in the theory of empirical process (Geer, 2000) and statistical learning theory (Koltchinskii, 2001; 2006; Bartlett & Mendelson, 2002; Bartlett et al., 2005). The concept of critical radius or critical inequality is used in non-parametric regression (Wainwright, 781 2019) and in off-policy evaluation (Duan et al., 2021; Duan & Wainwright, 2022; 2023; Mou et al., 782 2022). 783

784 785

786

799 800

803 804

809

## B PROOF OF THEOREM 4.1

787 To prove Lemma E.3, we first define the following event

$$\mathcal{E} := \left\{ \sup_{\Delta \in \mathsf{C}(\varepsilon)} \Delta^{\top} \eta \leq \mathcal{G}(\varepsilon) \quad \forall \varepsilon \in E \right\}.$$
(17)

792 When  $\mathcal{E}$  happens, the quantity  $\mathcal{G}(\varepsilon)$  can upper bound the supremum of the Gaussian process we care 793 about, and hence, we can effectively upper bound the uncertainty for any stochastic policy using 794  $\mathcal{G}(\cdot)$ . It follows from the Definition 3.2 that the event  $\mathcal{E}$  happens with probability a leas  $1 - \delta$ .

We can now prove all the claims in the theorem, starting from the first and the second. A comparator policy  $\pi$  identifies a weight vector w, an improvement vector  $\Delta$  and a radius  $\varepsilon$  such that  $w = \hat{\mu} + \Delta$ and  $\Delta \in C(\varepsilon)$ . In fact, we can always take  $\varepsilon$  to be the minimal value such that  $\Delta \in C(\varepsilon)$ . The first claim in Equation (11) can be proved by establishing that with probability at least  $1 - \delta$ 

$$w^{\top}r - \pi_{TRUST}^{\top}r = \Delta^{\top}r - \widehat{\Delta}_{*}^{\top}r \le 2\mathcal{G}\left(\left\lceil \varepsilon \right\rceil\right), \tag{18}$$

801 where  $\pi_{TRUST}$  is the policy weight returned by Algorithm 1. In order to show Equation (18), we can 802 decompose  $\widehat{\Delta}_*^\top r$  using the fact that  $\widehat{\varepsilon}_* \in E$  and  $\widehat{\Delta}_* \in C(\widehat{\varepsilon}_*)$  to obtain

$$\widehat{\Delta}_{*}^{\top}r = \widehat{\Delta}_{*}^{\top}\widehat{r} - \widehat{\Delta}_{*}^{\top}\eta \ge \widehat{\Delta}_{*}^{\top}\widehat{r} - \mathcal{G}\left(\widehat{\varepsilon}_{*}\right) = \widehat{\Delta}_{*}^{\top}\widehat{r} - \mathcal{G}\left(\left\lceil\widehat{\varepsilon}_{*}\right\rceil\right).$$
<sup>(19)</sup>

To further lower bound the RHS above, we have the following lemma, which shows that Algorithm 1 can be written in an equivalent way.

Lemma B.1. *The output of Algorithm 1 satisfies* 

 $\left(\widehat{\varepsilon}_{*},\widehat{\Delta}_{*}\right) = \underset{\varepsilon \leq \varepsilon_{0}, \Delta \in \mathsf{C}(\varepsilon)}{\arg\max} \left[\Delta^{\top}\widehat{r} - \mathcal{G}\left(\left\lceil \varepsilon \right\rceil\right)\right].$ (20)

This shows that Algorithm 1 optimizes over an objective function which consists of a signal term (i.e.,  $\Delta^{\top} \hat{r}$ ) minus a noise term (i.e.,  $\mathcal{G}(\lceil \varepsilon \rceil)$ ). Applying this lemma to equation 19, we know

$$\widehat{\Delta}_{*}^{\top} r \ge \Delta^{\top} \widehat{r} - \mathcal{G}\left(\left\lceil \varepsilon \right\rceil\right) = \Delta^{\top} r + \Delta^{\top} \eta - \mathcal{G}\left(\left\lceil \varepsilon \right\rceil\right).$$
(21)

After recalling that under  $\mathcal{E}$ 

$$\Delta^{\top} \eta \leq \sup_{\Delta \in \mathsf{C}(\varepsilon)} \Delta^{\top} \eta \leq \sup_{\Delta \in \mathsf{C}(\lceil \varepsilon \rceil)} \Delta^{\top} \eta \leq \mathcal{G}\left(\lceil \varepsilon \rceil\right),\tag{22}$$

plugging the equation 22 back into equation 21 concludes the bound in Equation (18), which also
 proves our first claim. Rearranging the terms in Equation (18) and taking supremum over all
 comparator policies, we obtain

$$\widehat{\Delta}_{*}^{\top} r \geq \sup_{\varepsilon \leq \varepsilon_{0}, \Delta \in \mathsf{C}(\varepsilon)} \left[ \Delta^{\top} r - 2\mathcal{G}\left( \left\lceil \varepsilon \right\rceil \right) \right],$$
(23)

which proves the first claim since  $V^{\pi_{TRUST}} - V^{\hat{\mu}} = \widehat{\Delta}_*^\top r$ .

In order to prove the last claim, it suffices to lower bound the policy value of the reference policy  $\hat{\mu}$ . From equation 5, we have  $\hat{\mu}(\hat{r}-r) \sim \mathcal{N}(0, 1/[\sum_{i=1}^{d} N_i/\sigma_i^2])$ , which implies with probability at least  $1 - \delta$ ,

$$\widehat{\mu}\left(\widehat{r}-r\right) \le \sqrt{\frac{2\log(1/\delta)}{\sum_{i=1}^{d} N_i/\sigma_i^2}} \tag{24}$$

from the standard Hoeffding inequality (e.g., Prop 2.5 in (Wainwright, 2019)). Combining equation 19 and equation 24, we obtain

$$\pi_{TRUST}^{\top} r = \widehat{\mu}^{\top} r + \widehat{\Delta}_{*}^{\top} r \ge \widehat{\mu}^{\top} \widehat{r} + \widehat{\mu}^{\top} (r - \widehat{r}) + \widehat{\Delta}_{*}^{\top} \widehat{r} - \mathcal{G}(\widehat{\varepsilon}_{*})$$
 (From equation 19)

$$\geq \pi_{TRUST}^{\top} \hat{r} - \mathcal{G}\left(\hat{\varepsilon}_{*}\right) - \sqrt{\frac{2\log(1/\delta)}{\sum_{i=1}^{d} N_{i}/\sigma_{i}^{2}}}$$
(From equation 24)

with probability at least  $1 - 2\delta$ . Therefore, we conclude.

### C ONE-SAMPLE CASE WITH STRONG SIGNALS

In this section, we give a simple example of one-sample-per-arm case. This can be view as a special case of data-starved MAB and Theorem 4.1 can be applied to get a non-trivial guarantees. Specifically, consider an MAB with 2d arms. Assume the true mean reward vector is  $r = (1, 1, ..., 1, 0, 0, ..., 0)^{\top}$  and the noise vector is  $\eta \sim \mathcal{N}(0, \sigma^2 I_{2d})$  That is, the first d arms have rewards independently sampled from  $\mathcal{N}(1, 1)$  and the rewards for other d arms are independently sampled from  $\mathcal{N}(0, 0)$ . The stochastic reference policy is set to the uniform one, i.e.,  $\hat{\mu} = (\frac{1}{d}, \frac{1}{d}, ..., \frac{1}{d})^{\top}$ .

We apply Algorithm 1 to this MAB instance. In the next theorem, we will show that for a specific  $\varepsilon$ , the optimal improvement in C ( $\varepsilon$ ) (denoted as  $\widehat{\Delta}_{\varepsilon}$  in equation 7) can achieve an improved reward value of constant level.

**Proposition C.1.** Assume  $r = (1, 1, ..., 1, 0, 0, ..., 0)^{\top}$  and noise  $\eta \sim \mathcal{N}(0, I_{2d})$ . For any  $0 \leq \varepsilon \leq \frac{1}{\sqrt{d}}$ , with probability at least  $1 - \delta$ , the improvement of policy value can be lower bounded by

$$\widehat{\Delta}_{\varepsilon}^{\top} r \geq \varepsilon \sqrt{d} \left[ \frac{1}{2} - \sigma \left( 1 + \sqrt{\frac{8 \log\left(2/\delta\right)}{d}} \right) \right],$$

where the improvement vector in  $C(\varepsilon)$  is defined in equation 7. Therefore, for  $\varepsilon = \frac{1}{\sqrt{d}}$  and  $d \ge 8 \log(2/\delta)$ , with probability at least  $1 - \delta$ , we can get a constant policy improvement

$$\widehat{\Delta}_{\varepsilon}^{\top} r \ge \frac{1}{2} - 2\sigma.$$

*Proof.* We define the optimal improvement vector as

$$\Delta_{\varepsilon}^* := \underset{\Delta \in \mathsf{C}(\varepsilon)}{\arg\max} \, \Delta^\top r.$$

Then, from the definition of  $\widehat{\Delta}_{\varepsilon}$ , we have

$$\widehat{\Delta}_{\varepsilon}^{\top} r = \widehat{\Delta}_{\varepsilon}^{\top} \widehat{r} - \widehat{\Delta}_{\varepsilon}^{\top} \eta \ge (\Delta_{\varepsilon}^{*})^{\top} \widehat{r} - \widehat{\Delta}_{\varepsilon}^{\top} \eta = (\Delta_{\varepsilon}^{*})^{\top} r + (\Delta_{\varepsilon}^{*})^{\top} \eta - \widehat{\Delta}_{\varepsilon}^{\top} \eta \ge \underbrace{(\Delta_{\varepsilon}^{*})^{\top} r}_{\text{signal}} - \underbrace{\left[\sup_{\Delta \in \mathsf{C}(\varepsilon)} \Delta^{\top} \eta - (\Delta_{\varepsilon}^{*})^{\top} \eta\right]}_{(25)}$$

In order to lower bound the policy value improvement, it suffices to lower bound the signal part and upper bound the noise. We denote  $\mathcal{H} = \{x \in \mathbb{R}^d : \sum_{i=1}^d x_i = 0\}$  as a hyperplane in  $\mathbb{R}^d$ . To deal with the signal part, it suffices to notice that

$$\mathbb{C}(\varepsilon) \subset \mathcal{H} \cap \mathbb{B}_2^d(\varepsilon).$$

We denote  $r_{\parallel}$  as the orthogonal projection of r on the  $\mathcal{H}$  and  $r_{\perp} = r - r_{\parallel}$ . In the strong signal case, we have

$$r_{\parallel} = \left(\frac{1}{2}, \frac{1}{2}, \dots, \frac{1}{2}, -\frac{1}{2}, -\frac{1}{2}, \dots, -\frac{1}{2}\right)^{\top}, \quad r_{\perp} = \left(\frac{1}{2}, \frac{1}{2}, \dots, \frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \dots, \frac{1}{2}\right)^{\top}.$$

Then, the signal part satisfies

$$\sup_{\Delta \in \mathsf{C}(\varepsilon)} \Delta^{\top} r = \sup_{\Delta \in \mathsf{C}(\varepsilon)} \Delta^{\top} r_{\parallel} \le \sup_{\Delta \in \mathcal{H} \cap \mathbb{B}_{2}^{d}(\varepsilon)} \Delta^{\top} r_{\parallel} = \left(\varepsilon \cdot \frac{r_{\parallel}}{\|r_{\parallel}\|_{2}}\right)^{\top} r_{\parallel} = \varepsilon \|r_{\parallel}\|_{2} = \frac{\varepsilon \sqrt{d}}{2}.$$
 (26)

On the other hand, we notice that when  $\varepsilon \leq \frac{1}{\sqrt{d}}$ ,

$$\varepsilon \cdot \frac{r_{\parallel}}{\left\|r_{\parallel}\right\|_{2}} = \left(\frac{\varepsilon}{\sqrt{d}}, \frac{\varepsilon}{\sqrt{d}}, ..., \frac{\varepsilon}{\sqrt{d}}, -\frac{\varepsilon}{\sqrt{d}}, -\frac{\varepsilon}{\sqrt{d}}, ..., -\frac{\varepsilon}{\sqrt{d}}\right)^{\top} \in \mathsf{C}\left(\varepsilon\right)$$

So actually the inequality in the equation 26 should be an equation, which implies

$$\sup_{\Delta \in \mathsf{C}(\varepsilon)} \Delta^{\top} = \frac{\varepsilon \sqrt{d}}{2}.$$
(27)

For the noise part, we decompose the noise as  $\eta = \eta_{\perp} + \eta_{\parallel}$ , where  $\eta_{\parallel}$  is the orthogonal projection of  $\eta$  on  $\mathcal{H}$ . Then, from  $C(\varepsilon) \subset \mathcal{H} \cap \mathbb{B}_2^d(\varepsilon)$ , one has

$$\sup_{\Delta \in \mathsf{C}(\varepsilon)} \Delta^{\top} \eta = \sup_{\Delta \in \mathsf{C}(\varepsilon)} \Delta^{\top} \left( \eta_{\parallel} + \eta_{\perp} \right) = \sup_{\Delta \in \mathsf{C}(\varepsilon)} \Delta^{\top} \eta_{\parallel} \leq \sup_{\Delta \in \mathcal{H} \cap \mathbb{B}_{2}^{d}(\varepsilon)} \Delta^{\top} \eta_{\parallel}$$
$$= \left( \varepsilon \cdot \frac{\eta_{\parallel}}{\|\eta_{\parallel}\|_{2}} \right)^{\top} \eta_{\parallel} = \varepsilon \|\eta_{\parallel}\|_{2} \leq \varepsilon \|\eta\|_{2}.$$

This implies  $\widehat{\Delta}_{\varepsilon}^{\top} r \ge (\Delta_{\varepsilon}^{*})^{\top} r - [\varepsilon ||\eta||_{2} - (\Delta_{\varepsilon}^{*})^{\top} \eta]$ . From our assumption,  $\frac{1}{\sigma^{2}} ||\eta||_{2}^{2}$  is a chi-square random variable with degree d, so from the Example 2.11 in (Wainwright, 2019), we know with probability at least  $1 - \delta/2$ , one has

$$\frac{\|\eta\|_2^2}{d\sigma^2} \le 1 + \sqrt{\frac{8\log\left(2/\delta\right)}{d}}$$

This implies

$$\left\|\eta\right\|_{2} \leq \sqrt{d\sigma^{2}\left(1 + \sqrt{\frac{8\log\left(2/\delta\right)}{d}}\right)} \leq \sqrt{d}\sigma\left(1 + \sqrt{\frac{2\log\left(2/\delta\right)}{d}}\right).$$

The last inequality comes from  $\sqrt{1+u} \le 1 + \frac{u}{2}$  for positive u. Moreover, since  $\Delta_{\varepsilon}^*$  is a fixed vector, we know  $(\Delta_{\varepsilon}^*)^{\top} \eta \sim \mathcal{N}\left(0, \sigma^2 \|\Delta_{\varepsilon}^*\|_2^2\right)$ . So with probability at least  $1 - \delta/2$ , one has

916  
917 
$$(\Delta_{\varepsilon}^{*})^{\top} \eta \ge -\sigma \|\Delta_{\varepsilon}^{*}\|_{2} \sqrt{2\log\left(\frac{2}{\delta}\right)} \ge -\sigma\varepsilon \sqrt{2\log\left(\frac{2}{\delta}\right)}$$

Combining the two terms above, one has with probability at least  $1 - \delta$ , it holds

$$\varepsilon \|\eta\|_2 - (\Delta_{\varepsilon}^*)^{\top} \eta \le \varepsilon \sqrt{d}\sigma \left(1 + \sqrt{\frac{2\log\left(2/\delta\right)}{d}}\right) + \sigma \varepsilon \sqrt{2\log\left(\frac{2}{\delta}\right)} = \varepsilon \sqrt{d}\sigma \left(1 + \sqrt{\frac{8\log\left(2/\delta\right)}{d}}\right). \tag{28}$$

Combining equation 25, equation 27 and equation 28, we finish the proof.

#### D MONTE CARLO COMPUTATION

918

925 926 927

928 929

930

931

932

933

934

935

936 937

939

941

945

950 951 952

953

954

960

961

962 963 964

965

966 967 968

969

Algorithm 2 Monte-Carlo method for computing  $\mathcal{G}(\varepsilon)$ 

**Input:** Offline dataset  $\mathcal{D}$ , the radius value  $\varepsilon \in E$ , the total sample size M and threshold  $M_0$ . Independently sample M noise vectors, denoted as  $\eta_i$  for  $i \in [M]$ , where  $\eta_i \sim$ 1.  $\mathcal{N}(0,\sigma_i^2/N(a_i),\sigma_i^2)$  is the noise variance for the *i*-th arm and  $N(a_i)$  is the sample size for  $a_i$ in  $\mathcal{D}$ . 2. Solve  $X_i := \sup_{\Delta \in \mathsf{C}(\varepsilon)} \Delta^{\top} \eta_i$  for  $i \in [M]$  and order them as  $X_{(1)} \leq X_{(2)} \leq ... \leq X_{(M)}$ .

3. **Return**  $X_{(M-M_0+1)}$  as an estimate of  $\mathcal{G}(\varepsilon)$  defined in Definition 3.2.

As discussed in Section 3, we can estimate  $\mathcal{G}(\varepsilon)$  using classical Monte Carlo method. In this 938 section, we illustrate the detailed implementation. We first sample M i.i.d. noise and then solve  $\sup_{\Delta \in \mathsf{C}(\varepsilon)} \Delta^{\perp} \eta$  for each to get M suprema. We eventually select the  $M_0$ -th largest values of all 940 suprema as our estimate for the bonus function, where  $M_0$  is a pre-computed integer dependent on M and the pre-determined failure probability  $\delta > 0$ . Here, the program  $\sup_{\Delta \in C(\varepsilon)} \Delta^{\top} \eta$  is a 942 second-order cone program and can be efficiently solved via standard off-the shelf libraries (Alizadeh 943 & Goldfarb, 2003; Boyd & Vandenberghe, 2004; Diamond & Boyd, 2016). The pseudocode for the 944 Monte-Carlo sampling is in Algorithm 2.

To determine  $M_0$ , we denote  $\eta_i$  as the i.i.d. noise vector for  $i \in [M]$  and  $X_i = \sup_{\Delta \in \mathsf{C}(\varepsilon)} \Delta^{\top} \eta$ . 946 We denote the order statistics of  $X_i$ -s as  $X_{(1)} \leq X_{(2)} \leq ... \leq X_{(M)}$ . Suppose the cumulative 947 distribution function of  $X_i$  is F(x), then from the property of the order statistics, we know the 948 cumulative distribution function of  $X_{(M-M_0+1)}$  is 949

$$F_{X_{(M-M_0+1)}}(x) = \sum_{j=M-M_0+1}^{M} C_M^j \left(F(x)\right)^j \left(1 - F(x)\right)^{M-j}.$$

We denote  $q_{1-\delta}$  as the  $(1 - \delta)$ -lower quantile of the random variable X, then we have  $F_{X_{(M-M_0+1)}}(q_{1-\delta}) = \sum_{j=M-M_0+1}^{M} C_M^j (1-\delta)^j (\delta)^{M-j}$ . For integer M and  $\delta > 0$ , we define  $Q(M, \delta)$  as the maximal integer  $M_0$  such that  $\sum_{j=M-M_0+1}^{M} C_M^j (1-\delta)^j (\delta)^{M-j} \leq \delta$ . With this definition, we take a fixed M and a total failure tolerance  $\delta$  for all  $\varepsilon \in E$ , then we take

$$M_0 = Q\left(M, \frac{\delta}{2|E|}\right)$$

as the threshold number. Under this choice, for any  $\varepsilon \in E$ , with probability at least  $1 - \delta/2|E|$ , it holds  $X_{(M-M_0+1)} > q_{1-\delta/2|E|}$ . On the other hand, with probability  $1 - \delta/2|E|$ , it holds that  $\sup_{\Delta \in \mathcal{C}(\varepsilon)} \Delta^{\top} \eta \leq q_{1-\delta/2|E|}$  This implies

$$\sup_{\Delta \in \mathcal{C}(\varepsilon)} \Delta^{\top} \eta \le q_{1-\delta/2|E|} < X_{(M-M_0+1)}$$

with probability at least  $1 - \delta/|E|$ . From a union bound, we know with probability at least  $1 - \delta$ , the bound above holds for any  $\varepsilon \in E$ .

#### Ε A FINE-GRAINED ANALYSIS TO THE SUBOPTIMALITY

We have shown a problem-dependent upper bound for the suboptimality in equation 12. In this 970 section, we will give a further upper bound for  $\mathcal{G}(\varepsilon)$  and hence, for the suboptimality. We have the 971 following theorem. The proof is deferred to Appendix E.1.

972 **Theorem E.1.** For a policy  $\pi$  (deterministic or stochastic), we denote its reward value as  $V^{\pi}$ . TRUST has the following properties. 974

1. We denote a comparator policy as a triple  $(\varepsilon, \Delta, \pi)$  such that  $\varepsilon = \sum_{i=1}^{d} \frac{\sigma_i^2 \Delta_i^2}{N_i}, \pi = \hat{\mu} + \Delta$ . We take the discrete candidate set *E* defined in equation 10. With probability at least  $1 - \delta$ , for any stochastic comparator policy  $(\varepsilon, \Delta, \pi)$ , the sub-optimality of the output policy of Algorithm 1 can be upper bounded as

$$V^{\pi} - V^{\pi_{TRUST}} \le 2\sqrt{2\sum_{i=1}^{d} \frac{\alpha \Delta_{i}^{2} \sigma_{i}^{2}}{N_{i}} \log\left(\frac{2|E|}{\delta}\right)} + 2\min\left\{\sqrt{\sum_{i=1}^{d} \frac{\alpha \Delta_{i}^{2} \sigma_{i}^{2}}{N_{i}}}, 4D\sqrt{\log_{+}\left(\frac{4ed\sum_{i=1}^{d} \frac{\alpha \Delta_{i}^{2} \sigma_{i}^{2}}{N_{i}}}{D^{2}}\right)}\right\}$$

$$(29)$$

where D is defined as any quantity satisfying

$$D \ge \sqrt{\max_{i \in [d]} \left[\frac{\sigma_i^2}{N_i} - \frac{2\sigma_i^2}{N}\right]} + \frac{\sum_{j=1}^d N_j \sigma_j^2}{N^2}.$$
 (30)

 $\alpha$  is the decaying rate defined in equation 10,  $\log_+(a) = \max(1, \log(a))$ .

2. (Comparison with the optimal policy) We further assume  $\sigma_i = 1$  for  $i \in [d]$  and assume the offline dataset is generated from the policy  $\mu(\cdot)$  with  $\min_{i \in [d]} \mu(a_i) > 0$ . Without loss of generality we assume  $a_1$  is the optimal arm and denote the optimal policy as  $\pi_*$ . We write

$$C^* := \frac{1}{\mu(a_1)}, \quad C_{\min} := \frac{1}{\min_{i \in [d]} \mu(a_i)}.$$
 (31)

When  $N \ge 8C_{\min} \log(d/\delta)$ , with probability at least  $1 - 2\delta$ , one has

$$V^{\pi_*} - V^{\pi_{TRUST}} \lesssim \sqrt{\frac{C_{\min}}{N} \log_+\left(\frac{dC^*}{C_{\min}}\right)} + \sqrt{\frac{C^*}{N} \log\left(\frac{2|E|}{\delta}\right)}.$$
 (32)

Specially, when  $C_{\min} \simeq C^*$ , we have with probability at least  $1 - 2\delta$ ,

$$V^{\pi_*} - V^{\pi_{TRUST}} \lesssim \sqrt{\frac{C^*}{N} \log\left(\frac{2d|E|}{\delta}\right)}.$$
(33)

1007

1003 1004

975 976

977

978

991

992

993

We remark that equation 29 is problem-dependent, and it gives an explicit upper bound for  $\mathcal{G}([\varepsilon])$ in equation 12. This is derived by first concentrating  $\mathcal{G}(\varepsilon)$  around  $\mathbb{E}\sup_{\Delta \in \mathsf{C}(\varepsilon)} \Delta^{\top} \eta$ , which is well-defined as localized Gaussian width or local Gaussian complexity (Koltchinskii, 2006), and then upper bounding the localized Gaussian width of a convex hull via tools in convex analysis (Bellec, 2019). Different from Theorem 2.1, when  $\pi = a_i$  represents a single arm, equation 29 relies not only on  $\sigma_i^2/N_i$ , but on  $\sigma_j^2/N_j$  for  $j \neq i$  as well, since the size of trust regions depend on  $\sigma_i^2/N_i$  for all  $i \in [d]$ .

1015 Notably, equation 33 gives an analogous upper bound depending on  $\mu(\cdot)$  and N, which is comparable 1016 to the bound for LCB in Theorem 2.1 up to constant and logarithmic factors. This indicates that, when 1017 behavioral cloning policy is not too imbalanced, TRUST is guaranteed to achieve the same level of 1018 performance as LCB. In fact, this improvement is remarkable since TRUST is exploring a much larger 1019 policy searching space than LCB, which encompasses all stochastic policies (the whole simplex) rather than the set of all single arms only. We also remark that both the bound in Theorem 2.1 and 1020 in equation 46 are worst-case upper bound, and in practice, we will show in Section 5 that in some 1021 settings, TRUST can achieve good performance while LCB fails completely. 1022

1023

Is TRUST minimax-optimal? We consider the hard cases in MAB (Rashidinejad et al., 2021)
 where LCB achieves the minimax-optimal upper bound and we show for these hard cases, TRUST will achieve the same sample complexity as LCB up to log and constant factors. More specifically, we

1027 consider a two-arm MAB  $\mathcal{A} = \{1, 2\}$  and the uniform behavioral cloning policy  $\mu(1) = \mu(2) = 1/2$ . For  $\delta$  in[0, 1/4], we define  $\mathcal{M}_1$  and  $\mathcal{M}_2$  are two MDPs whose reward distributions are as follows.

1029  
1030
$$\mathcal{M}_1: r(1) \sim \mathsf{Bernoulli}\left(\frac{1}{2}\right), \ r(2) \sim \mathsf{Bernoulli}\left(\frac{1}{2} + \delta\right)$$
1031
(1)
(1)

 $\mathcal{M}_2: r(1) \sim \mathsf{Bernoulli}\left(rac{1}{2}
ight), \ r(2) \sim \mathsf{Bernoulli}\left(rac{1}{2} - \delta
ight),$ 

where Bernoulli (p) is the Bernoulli distribution with probability p. The next result is a corollary from Theorem E.1.

**Corollary E.2.** We define  $\mathcal{M}_1$ ,  $mdp_2$  as above for  $\delta \in [0, 1/4]$ . Assume  $N \ge \widetilde{O}(1)$ . Then, we have 1037

1. The minimax optimal lower bound for the suboptimality of LCB is

$$\inf_{\widehat{a}_{\mathsf{LCB}}\in\mathcal{A}} \sup_{\mathcal{M}\in\{\mathcal{M}_1,\mathcal{M}_2\}} \mathbb{E}_{\mathcal{D}}\left[r(a^*) - r(\widehat{a}_{\mathsf{LCB}})\right] \gtrsim \sqrt{\frac{C^*}{N}},\tag{34}$$

where  $\mathbb{E}_{\mathcal{D}}[\cdot]$  is the expectation over the offline dataset  $\mathcal{D}$ .

1044 2. The upper bound for suboptimality of TRUST mathces the lower bound above up to log factor. 1045 Namely, for any  $\mathcal{M} \in {\mathcal{M}_1, \mathcal{M}_2}$ , one has

$$\mathbb{E}_{\mathcal{D}}\left[r(a^*) - V^{\pi_{TRUST}}\right] \lesssim \sqrt{\frac{C^* \log(dN)}{N}}.$$
(35)

The first claim comes from Theorem 2 of (Rashidinejad et al., 2021), while the second claim is adirect corollary to Theorem E.1.

1053 E.1 PROOF OF THEOREM E.1

*Proof.* Recall from Theorem 4.1 that for any comparator policy  $(\varepsilon, \Delta, \pi)$  defined above, one has 

 $V^{\pi} - V^{\pi_{TRUST}} \le 2\mathcal{G}\left(\left\lceil \varepsilon \right\rceil\right),$ 

where  $\lceil \varepsilon \rceil := \inf \{ \varepsilon' \in E : \varepsilon \le \varepsilon' \}$ . The following lemma upper bounds the quantile of Gaussian suprema  $\mathcal{G}(\varepsilon)$  for each  $\varepsilon \in E$ . The proof is deferred to Appendix E.2.

**Lemma E.3.** For  $\varepsilon \in E$ , one can upper bound  $\mathcal{G}(\varepsilon)$  as follows.

$$\mathcal{G}(\varepsilon) \le \min\left\{\varepsilon \cdot \sqrt{d} , \ 4D\sqrt{\log_+\left(\frac{4ed\varepsilon^2}{D^2}\right)}\right\} + \sqrt{2\varepsilon^2 \log\left(\frac{2|E|}{\delta}\right)}$$
(36)

where  $\log_+(a) = \max(1, \log(a))$  and D is a quantity satisfying

$$D \ge \sqrt{\max_{i \in [d]} \left[\frac{\sigma_i^2}{N_i} - \frac{2\sigma_i^2}{N}\right] + \frac{\sum_{j=1}^d N_j \sigma_j^2}{N^2}}.$$
(37)

1071 Applying Lemma E.3 to  $[\varepsilon] \in E$ , we obtain

$$V^{\pi} - V^{\pi_{TRUST}} \le 2\min\left\{\left\lceil\varepsilon\right\rceil \cdot \sqrt{d} , \ 4D\sqrt{\log_{+}\left(\frac{4ed\left\lceil\varepsilon\right\rceil^{2}}{D^{2}}\right)}\right\} + 2\sqrt{2\left\lceil\varepsilon\right\rceil^{2}\log\left(\frac{2|E|}{\delta}\right)}.$$
(38)

Since  $\varepsilon = \sum_{i=1}^{d} \frac{\sigma_i^2 \Delta_i^2}{N_i}$ , we know from our discretization scheme in equation 10

$$\lceil \varepsilon \rceil \le \alpha \cdot \sum_{i=1}^{d} \frac{\sigma_i^2 \Delta_i^2}{N_i}.$$
(39)

Bridging equation 39 into equation 38, we obtain our first claim. In order to get the second claim, we take  $\sigma_i = 1$  for  $i \in [d]$  and  $\Delta = \pi_* - \hat{\mu}$ , which is the vector pointing at the vertex corresponding to the optimal arm from the uniform reference policy  $\hat{\mu}$  defined in equation 5. Then, we have

$$\sum_{i=1}^{d} \frac{\Delta_i^2 \sigma_i^2}{N_i} = \frac{1}{N_1} - \frac{2}{N} + \frac{1}{N} = \frac{1}{N_1} - \frac{1}{N} \le \frac{1}{N_1},$$

where  $N_1$  is the sample size for the optimal arm  $a_1$ . Therefore, we can further bound equation 38 as 

$$V^{\pi_*} - V^{\pi_{TRUST}} \le 4D\sqrt{\log_+\left(\frac{4\alpha ed}{N_1 D^2}\right)} + 2\sqrt{\frac{2\alpha}{N_1}\log\left(\frac{2|E|}{\delta}\right)}.$$
(40)

Finally, we take a specific value of D and lower bound  $N_1$  via Chernoff bound in Lemma E.7. From Lemma E.7, we know that when  $N \ge 8C_{\min}\log(d/\delta)$ , with probability at least  $1 - \delta$ , we have 

$$N_i \ge \frac{1}{2} N \mu(a_i) \tag{41}$$

for any  $i \in [d]$ . Recall the definition of D in equation 37, we know that D can be arbitrary value greater than  $\sqrt{\max_{i \in [d]} \left[\frac{\sigma_i^2}{N_i} - \frac{2\sigma_i^2}{N}\right] + \frac{\sum_{j=1}^d N_j \sigma_j^2}{N}}$ . Then, when  $\sigma_i = 1$ , one has

 $\sqrt{\max_{i \in [d]} \left[\frac{\sigma_i^2}{N_i} - \frac{2\sigma_i^2}{N}\right] + \frac{\sum_{j=1}^d N_j \sigma_j^2}{N^2}} \le \sqrt{\frac{1}{\min_{i \in [d]} N_i}}.$ 

We denote  $N_j = \min_{i \in [d]} N_i$  (when there are multiple minimizers, we arbitrarily pick one). Then, we have 

$$\sqrt{\max_{i \in [d]} \left[\frac{\sigma_i^2}{N_i} - \frac{2\sigma_i^2}{N}\right] + \frac{\sum_{j=1}^d N_j \sigma_j^2}{N^2}} \le \sqrt{\frac{1}{N_j}} \le \sqrt{\frac{2}{N\mu(a_j)}} \le \sqrt{\frac{2}{N \cdot \min_{i \in [d]} \mu(a_i)}} = \sqrt{\frac{2C_{\min}}{N}}$$

Therefore, we take  $D = \sqrt{\frac{2C_{\min}}{N}}$  in equation 40 and apply  $N_1 \ge \frac{1}{2}N\mu(a_i)$  to obtain 

$$V^{\pi_*} - V^{\pi_{TRUST}} \le 4\sqrt{\frac{2C_{\min}}{N}\log_+\left(\frac{4\alpha edC^*}{C_{\min}}\right)} + 4\sqrt{\frac{\alpha C^*}{N}\log\left(\frac{2|E|}{\delta}\right)},\tag{42}$$

which proves equation 32. Finally, when  $C^* \simeq C_{\min}$ , one has 

$$V^{\pi_*} - V^{\pi_{TRUST}} \lesssim \sqrt{\frac{C^*}{N} \log\left(\frac{2d|E|}{\delta}\right)}$$

Therefore, we conclude. 

#### E.2 PROOF OF LEMMA E.3

*Proof.* Recall that  $\Delta = (\Delta_1, \Delta_2, ..., \Delta_d)^{\top}$  is the improvement vector,  $\eta = (\eta_1, \eta_2, ..., \eta_d)^{\top}$  is the noise vector, where entries are independent and  $\eta_i \sim \mathcal{N}(0, \sigma_i^2/N_i)$  and  $N_i$  is the sample size of arm  $a_i$  in the offline dataset. To proceed with the proof, let's further define 

$$\widetilde{\eta} = \left(\widetilde{\eta}_1, \widetilde{\eta}_2, ..., \widetilde{\eta}_d\right)^\top, \quad \widetilde{\Delta} = \left(\widetilde{\Delta}_1, \widetilde{\Delta}_2, ..., \widetilde{\Delta}_d\right)^\top, \quad \text{where} \quad \widetilde{\eta}_i = \eta_i \frac{\sqrt{N_i}}{\sigma_i}, \ \widetilde{\Delta}_i = \frac{\Delta_i \sigma_i}{\sqrt{N_i}}.$$
(43)

With this notation, one has 

$$\mathcal{N}(0, I_d), \quad \eta^{\top} \Delta = \widetilde{\eta}^{\top} \widetilde{\Delta}$$

We also write the equivalent trust region (for  $\widetilde{\Delta}$ ) as 

 $\widetilde{\eta} \sim$ 

$$\widetilde{\mathsf{C}}\left(\varepsilon\right) = \left\{\widetilde{\Delta} \in \mathbb{R}^{d} : \frac{\sqrt{N_{i}}}{\sigma_{i}}\widetilde{\Delta}_{i} + \widehat{\mu}_{i} \ge 0, \quad \sum_{i=1}^{d} \left[\frac{\sqrt{N_{i}}}{\sigma_{i}}\widetilde{\Delta}_{i} + \widehat{\mu}_{i}\right] = 1, \quad \left\|\widetilde{\Delta}\right\|_{2} \le \varepsilon\right\}, \tag{44}$$

where  $\hat{\mu} = (\hat{\mu}_1, \hat{\mu}_2, ..., \hat{\mu}_d)^{\top}$  is the policy weight for the reference policy. From the definition above, one has for any  $\varepsilon > 0$ ,

$$\Delta\in\mathsf{C}\left(\varepsilon\right)\ \Leftrightarrow\ \widetilde{\Delta}\in\widetilde{\mathsf{C}}\left(\varepsilon\right).$$

Then, we apply Lemma E.4 to  $\sup_{\Delta \in C(\varepsilon)} \Delta^{\top} \eta$  for a  $\varepsilon \in E$ . One has with probability at least  $1 - \frac{\delta}{|E|}$ ,

$$\left|\sup_{\Delta\in\mathsf{C}(\varepsilon)}\Delta^{\top}\eta - \mathbb{E}\sup_{\Delta\in\mathsf{C}(\varepsilon)}\Delta^{\top}\eta\right| \leq \sqrt{2\varepsilon^{2}\log\left(\frac{2|E|}{\delta}\right)}$$

From a union bound, one immediately has with probability at least  $1 - \delta$ , for any  $\varepsilon \in E$ , it holds that

$$\sup_{\Delta \in \mathsf{C}(\varepsilon)} \Delta^{\top} \eta \leq \mathbb{E} \sup_{\Delta \in \mathsf{C}(\varepsilon)} \Delta^{\top} \eta + \sqrt{2\varepsilon^2 \log\left(\frac{2|E|}{\delta}\right)}.$$
(45)

From the definition of  $\mathcal{G}(\varepsilon)$  in equation 3.2, we know that  $\mathcal{G}(\varepsilon)$  is the minimal quantity that satisfy equation 45 with probability at least  $1 - \delta$ . Therefore, one has

$$\begin{array}{l} 1148\\ 1149\\ 1150 \end{array} \quad \mathcal{G}\left(\varepsilon\right) \leq \mathbb{E} \sup_{\Delta \in \mathsf{C}(\varepsilon)} \Delta^{\top} \eta + \sqrt{2\varepsilon^{2} \log\left(\frac{2|E|}{\delta}\right)} = \mathbb{E}_{\widetilde{\eta} \sim \mathcal{N}(0,I_{d})} \left[ \sup_{\widetilde{\Delta} \in \widetilde{\mathsf{C}}(\varepsilon)} \widetilde{\Delta}^{\top} \widetilde{\eta} \right] + \sqrt{2\varepsilon^{2} \log\left(\frac{2|E|}{\delta}\right)} \quad \forall \varepsilon \in E$$

$$\begin{array}{l} (46) \end{array}$$

Note that, the first term in the RHS of equation 46 is well-defined as localized Gaussian width over the convex hull defined by the trust region C ( $\varepsilon$ ) (or equivalently,  $\tilde{C}(\varepsilon)$ ). We denote

$$T := \left\{ \widetilde{\Delta} \in \mathbb{R}^d : \frac{\sqrt{N_i}}{\sigma_i} \widetilde{\Delta}_i + \widehat{\mu}_i \ge 0, \quad \sum_{i=1}^d \left[ \frac{\sqrt{N_i}}{\sigma_i} \widetilde{\Delta}_i + \widehat{\mu}_i \right] = 1 \right\}.$$
(47)

1156 We immediately have that T is a convex hull of d points in  $\mathbb{R}^d$  and the vertices of this convex hull are 1157 the vertices of the simplex in  $\mathbb{R}^d$  shifted by the reference policy  $\hat{\mu}$ . In what follows, we plan to apply 1158 Lemma E.5 to the localized Gaussian width of  $T \cap \varepsilon \mathbb{B}_2$ . However, T is not subsumed by the unit ball 1159 in  $\mathbb{R}^d$ , so we need to do some additional scaling. Note that, the zero vector is included in T. Let's 1160 compute the farthest distance for the vertices of T. We denote the *i*-th vertex of T as

$$\widetilde{\Delta} = \left( -\frac{\sigma_1}{\sqrt{N_1}} \widehat{\mu}_1, \dots, -\frac{\sigma_{i-1}}{\sqrt{N_{i-1}}} \widehat{\mu}_{i-1}, \frac{\sigma_i}{\sqrt{N_i}} \left( 1 - \widehat{\mu}_i \right), -\frac{\sigma_{i+1}}{\sqrt{N_{i+1}}} \widehat{\mu}_{i+1}, \dots, -\frac{\sigma_d}{\sqrt{N_d}} \widehat{\mu}_d \right).$$
(48)

1164 The  $\ell_2$ -norm of this improvement vector is

$$\left\|\widetilde{\Delta}\right\|_2 = \sqrt{\frac{\sigma_i^2}{N_i} - \frac{2\sigma_i^2}{N} + \frac{\sum_{i=1}^d N_i \sigma_i^2}{N^2}},$$

1167 where N is the total sample size of the offline dataset. Therefore, the maximal radius of T can be 1168 upper bounded by D, where D is any quantity that satisfies

$$D \ge \sqrt{\max_{i \in [d]} \left[\frac{\sigma_i^2}{N_i} - \frac{2\sigma_i^2}{N}\right] + \frac{\sum_{j=1}^d N_j \sigma_j^2}{N^2}}.$$
(49)

We denote  $S = \frac{1}{D} \cdot T := \left\{ \frac{1}{D} \cdot x : x \in T \right\}$ . Then, from Lemma E.5, one has

$$\mathbb{E}_{\widetilde{\eta} \sim \mathcal{N}(0, I_d)} \left[ \sup_{\widetilde{\Delta} \in \widetilde{\mathsf{C}}(\varepsilon)} \widetilde{\Delta}^\top \widetilde{\eta} \right] = \mathbb{E}_{\widetilde{\eta} \sim \mathcal{N}(0, I_d)} \left[ \sup_{\widetilde{\Delta} \in T \cap \varepsilon \mathbb{B}_2} \widetilde{\Delta}^\top \widetilde{\eta} \right]$$
$$= D \cdot \mathbb{E}_{\widetilde{\eta} \sim \mathcal{N}(0, I_d)} \left[ \sup_{\widetilde{\lambda} \in \mathcal{N}(0, I_d)} \widetilde{\Delta}^\top \widetilde{\eta} \right]$$

1139 1140 1141

1143

1144

1154 1155

1161 1162 1163

1165 1166

1170

1171

$$(S \cap \frac{\varepsilon}{D} \cdot \mathbb{B}_2 \text{ can be got by scaling } T \cap \varepsilon \mathbb{B}_2 \text{ bt } \frac{1}{D})$$

$$(S + \frac{1}{D} - \frac{1}{D^2} \operatorname{can} \operatorname{bc} \operatorname{got} \operatorname{by} \operatorname{scanng} T + \operatorname{can} \operatorname{can} \operatorname{scanng} T + \operatorname{sc$$

1183  
1184 
$$\begin{bmatrix} V & (D^2 & f) \\ Take s = \frac{\varepsilon}{D} \text{ and } M = d \text{ in Lemma E.5} \end{bmatrix}$$

1185 
$$\left[\left(\begin{array}{c} \left( \left( \varepsilon^2 \right) \right) \right) \right] \left( \varepsilon \right] \right]$$

1186  
1187
$$= D \cdot \left[ \left( 4\sqrt{\log_+ \left( 4ed\left(\frac{\sigma}{D^2}\right) \right)} \right) \wedge \left(\frac{\sigma}{D}\sqrt{d}\right) \right].$$
( $\varepsilon \le D$  for any  $\varepsilon \in E$ )

1188 This finishes the proof.

1191 E.3 AUXILIARY LEMMAS

**Lemma E.4** (Concentration of Gaussian suprema, Exercise 5.10 in (Wainwright, 2019)). Let  $\{X_{\theta}, \theta \in \mathbb{T}\}$  be a zero-mean Gaussian process, and define  $Z = \sup_{\theta \in \mathbb{T}} X_{\theta}$ . Then, we have

$$\mathbb{P}[|Z - \mathbb{E}[Z]| \ge \delta] \le 2 \exp\left(-\frac{\delta^2}{2\sigma^2}\right),$$

1198 where  $\sigma^2 := \sup_{\theta \in \mathbb{T}} \operatorname{var} (X_{\theta})$  is the maximal variance of the process.

**Lemma E.5** (Localized Gaussian Width of a Convex Hull, Proposition 1 in (Bellec, 2019)). Let  $d \ge 1, M \ge 2$  and T be the convex hull of M points in  $\mathbb{R}^d$ . We write  $\mathbb{B}_2 = \{x \in \mathbb{R}^d : ||x||_2 \le 1\}$ and  $s\mathbb{B}_2 = \{s \cdot x : x \in \mathbb{R}^d, ||x||_2 \le 1\}$ . Assume  $T \subset \mathbb{B}_2^d(1)$ . Let  $g \in \mathbb{R}^d$  be a standard Gaussian vector. Then, for all s > 0, one has

$$\mathbb{E}\left[\sup_{x\in T\cap s\mathbb{B}_{2}}x^{\top}g\right] \leq \left(4\sqrt{\log_{+}\left(4eM\left(s^{2}\wedge1\right)\right)}\right)\wedge\left(s\sqrt{d\wedge M}\right),\tag{50}$$

1208 where  $\log_+(a) = \max(1, \log(a)), a \wedge b = \min\{a, b\}$ . 

1211 Lemma E.6 (Chernoff bound for binomial random variables, Theorem 2.3.1 in (Vershynin, 2020)). 1212 Let  $X_i$  be independent Bernoulli random variables with parameters  $p_i$ . Consider their sum  $S_N = \sum_{i=1}^{N} X_i$  and denote its mean by  $\mu = \mathbb{E}S_N$ . Then, for any  $t > \mu$ , we have

$$\mathbb{P}\left\{S_N \ge t\right\} \le e^{-\mu} \left(\frac{e\mu}{t}\right)^t.$$

1219 Lemma E.7 (Chernoff bound for offline MAB). Under the setting in Theorem E.1, we have

$$\mathbb{P}\left(N_i \ge \frac{1}{2}N\mu(a_i) \quad \forall i \in [d]\right) \le 1 - d\exp\left(-\frac{N \cdot \min_{j \in [d]} \mu(a_j)}{8}\right),$$

*Proof.* For arm  $i \in [d]$ , we take  $\mu = N\mu(a_i)$  and  $t = \frac{1}{2}M\mu(a_i)$  in Lemma E.6 and obtain

$$\mathbb{P}\left(N_i \ge \frac{1}{2}N\mu(a_i)\right) \le \exp\left(-N\mu(a_i)\right) \cdot \left(\frac{eN\mu(a_i)}{\frac{1}{2}N\mu(a_i)}\right)^{\frac{1}{2}N\mu(a_i)} = \exp\left(N\mu(a_i)\left[-1 + \frac{1}{2}\log(2e)\right]\right) \le \exp\left(-\frac{N\mu(a_i)}{8}\right).$$

1229 We finish the proof by a union bound for all arms.

## <sup>1231</sup> F PROOF OF LEMMA B.1

*Proof.* Recall the definition of  $[\varepsilon]$ :

$$\left[\varepsilon\right] := \inf\left\{\varepsilon' \in E : \varepsilon' \ge \varepsilon\right\}.$$
(51)

1237 We additionally define

$$\lfloor \varepsilon \rfloor := \sup \left\{ \varepsilon' \in E : \varepsilon' < \varepsilon \right\}.$$
(52)

1239 Specially, if there is no  $\varepsilon' \in E$  such that  $\varepsilon' < \varepsilon$ , then we define  $\lfloor \varepsilon \rfloor = 0$ . Then we know for any  $\varepsilon \le \varepsilon_0 \in E$  ( $\varepsilon_0$  is the largest possible radius) and a finite set E, it holds that 1241

$$[\varepsilon] < \varepsilon \le [\varepsilon], \text{ and } \varepsilon = [\varepsilon] \text{ if and only if } \varepsilon \in E.$$
 (53)

For any  $\varepsilon \in E$ , recall  $\widehat{\Delta}_{\varepsilon}$  is the optimal improvement vector within  $C(\varepsilon)$  defined in equation 7. It holds that 

$$\begin{aligned} \widehat{\Delta}_{\varepsilon} &:= \underset{\Delta \in \mathsf{C}(\varepsilon)}{\arg\max} \, \Delta^{\top} \widehat{r} = \underset{\Delta \in \mathsf{C}(\varepsilon)}{\arg\max} \left[ \Delta^{\top} \widehat{r} - \mathcal{G}\left(\varepsilon\right) \right] & (\text{since } \mathcal{G}\left(\varepsilon\right) \text{ does not depend on } \Delta) \\ &= \underset{\Delta \in \mathsf{C}(\varepsilon)}{\arg\max} \left[ \Delta^{\top} \widehat{r} - \mathcal{G}\left(\left\lceil \varepsilon \right\rceil\right) \right] & (\varepsilon \in E, \text{ so } \left\lceil \varepsilon \right\rceil = \varepsilon) \end{aligned}$$

$$\leq \mathop{\arg\max}_{\boldsymbol{\varepsilon}' \in ([\boldsymbol{\varepsilon}], [\boldsymbol{\varepsilon}]], \Delta \in \mathsf{C}(\boldsymbol{\varepsilon}')} \left[ \Delta^\top \widehat{r} - \mathcal{G}\left([\boldsymbol{\varepsilon}']\right) \right]$$

On the other hand, when  $\varepsilon \in E$  and  $\varepsilon' \in (|\varepsilon|, [\varepsilon]]$ , one has  $[\varepsilon'] = [\varepsilon] = \varepsilon$ , so

$$\begin{array}{ll} 1252 & \arg\max_{\varepsilon' \in (\lfloor \varepsilon \rfloor, \lceil \varepsilon \rceil], \Delta \in \mathsf{C}(\varepsilon')} \left[ \Delta^{\top} \widehat{r} - \mathcal{G}\left( \lceil \varepsilon' \rceil \right) \right] = \arg\max_{\varepsilon' \in (\lfloor \varepsilon \rfloor, \lceil \varepsilon \rceil], \Delta \in \mathsf{C}(\varepsilon')} \left[ \Delta^{\top} \widehat{r} - \mathcal{G}\left( \lceil \varepsilon \rceil \right) \right] \leq \arg\max_{\Delta \in \mathsf{C}(\varepsilon)} \left[ \Delta^{\top} \widehat{r} - \mathcal{G}\left( \lceil \varepsilon \rceil \right) \right]$$

where the last inequality comes from the fact that  $C(\varepsilon') \subset C(\varepsilon)$  when  $\varepsilon' \leq [\varepsilon] = \varepsilon$  by definition of the trust region in equation 6. Combining two inequalities above, we have for any  $\varepsilon \in E$ , 

$$\left(\varepsilon, \widehat{\Delta}_{\varepsilon}\right) = \arg\max_{\varepsilon' \in \left(\lfloor\varepsilon\rfloor, \lceil\varepsilon\rceil\right], \Delta \in \mathsf{C}(\varepsilon')} \left[\Delta^{\top} \widehat{r} - \mathcal{G}\left(\lceil\varepsilon'\rceil\right)\right],\tag{54}$$

where the variables in RHS above are  $\varepsilon'$  and  $\Delta$ , and Therefore, from the definition of we have

$$\left(\widehat{\varepsilon}_{*},\widehat{\Delta}_{*}\right) = \underset{\varepsilon \in E}{\operatorname{arg\,max}} \underset{\varepsilon' \in (\lfloor \varepsilon \rfloor, \lceil \varepsilon \rceil], \Delta \in \mathsf{C}(\varepsilon')}{\operatorname{arg\,max}} \left[\Delta^{\top}\widehat{r} - \mathcal{G}\left(\lceil \varepsilon' \rceil\right)\right] = \underset{\varepsilon \leq \varepsilon_{0}, \Delta \in \mathsf{C}(\varepsilon)}{\operatorname{arg\,max}} \left[\Delta^{\top}\widehat{r} - \mathcal{G}\left(\lceil \varepsilon \rceil\right)\right].$$

This finishes the proof.

#### AUGMENTATION WITH LCB Gì

To determine the most effective final policy, we can compare the outputs of the LCB and Algorithm 1 and combine both policies, based on the relative magnitude of their corresponding lower bounds. Specifically, the combined policy is 

$$\pi_{\rm combined} =$$

$$\begin{cases} \widehat{a}_{\mathsf{LCB}} \text{ If } \max_{a_i \in \mathcal{A}} l_i \geq w_{\mathsf{TR}}^\top \widehat{r} - \mathcal{G}\left(\lceil \widehat{\varepsilon}_* \rceil\right) - \sqrt{\frac{2\log(1/\delta)}{\sum_{j=1}^d N_j/\sigma_j^2}},\\ w_{\mathsf{TR}} \text{ If } \max_{a_i \in \mathcal{A}} l_i < w_{\mathsf{TR}}^\top \widehat{r} - \mathcal{G}\left(\lceil \widehat{\varepsilon}_* \rceil\right) - \sqrt{\frac{2\log(1/\delta)}{\sum_{j=1}^d N_j/\sigma_j^2}}, \end{cases}$$
(55)

where  $l_i = \hat{r}_i - b_i$  is defined in equation 1 and  $\mathcal{G}(\varepsilon)$  is defined in Definition 3.2. This combined policy will perform at least as well as LCB with high probability. More specifically, we have **Corollary G.1.** We denote the arm chosen by LCB as  $\hat{a}_{LCB}$ . We also denote  $r(\cdot)$  as the true reward

of a policy (deterministic or stochastic). With probability at least  $1-3\delta$ , one has 

$$V^{\pi_{\text{combined}}} \ge \max_{a_i \in \mathcal{A}} l_i.$$
(56)

*Proof.* We denote  $\hat{r}(\hat{a}_{LCB}) = r_{\hat{a}_{LCB}}$  and  $\hat{r}(w_{TR})$  as the empirical reward of the policy returned by LCB and Algorithm 1, respectively. Recall the uncertainty term of LCB in equation 1 and of Algorithm 1 in equation 55, we write  $b(\hat{a}_{LCB}) = b_{\hat{a}_{LCB}}$  and  $b(w_{TR}) = \mathcal{G}(\lceil \hat{\varepsilon}_* \rceil) + \sqrt{2\log(1/\delta)/[\sum_{j=1}^d N_j/\sigma_j^2]}$ . Then, from Theorem 4.1, equation 2 and a union bound, we know with probability at least  $1 - 3\delta$ , it holds that 

$$r(\widehat{a}_{\mathsf{LCB}}) \ge \widehat{r}(\widehat{a}_{\mathsf{LCB}}) - b(\widehat{a}_{\mathsf{LCB}}), \ r(w_{\mathsf{TR}}) \ge \widehat{r}(w_{\mathsf{TR}}) - b(w_{\mathsf{TR}}),$$

which implies 

1292  
1293  
1294  
1295  

$$V^{\pi_{\text{combined}}} \ge \hat{r}(\pi_{\text{combined}}) - b(\pi_{\text{combined}})$$
  
 $\ge \hat{r}(\hat{a}_{\text{LCB}}) - b(\hat{a}_{\text{LCB}})$  (By equation 55)  
 $= \max_{a_i \in \mathcal{A}} l_i.$  (By the definition of  $\hat{a}_{\text{LCB}}$  in equation 3)

Therefore, we conclude.

#### Η **EXPERIMENT DETAILS**

We did experiments on Mujoco environment in the D4RL dataset (Fu et al., 2020). All environments we test on are v3. Since the original D4RL dataset does not include the exact form of logging policies, we retrain SAC (Haarnoja et al., 2018) on these environment for 1000 episodes and keep record of the policy in each episode. We test 4 environments in two settings, denoted as '1-traj-low' and '1-traj-high'. In either setting, the offline dataset is generated from 100 policies with one trajectory from each. In the '1-traj-low' setting, the data is generated from the first 100 policies in the training process of SAC, while in the '1-traj-high' setting, it is generated from the policy in (10x + 1)-th episodes in the training process.

For all experiments on Mujoco, we average the results over 4 random seeds (from 2023 to 2026), and to run CQL, we use default hyper-parameters in https://github.com/young-geng/CQL to run 2000 episodes. For TRUST, we run it using a fixed standard deviation level  $\sigma_i = 150$  for all experiments.