# **Entropic Risk Optimization in Discounted MDPs:** Sample Complexity Bounds with a Generative Model

# Oliver Mortensen

University of Copenhagen olmo@di.ku.dk

## Sadegh Talebi

University of Copenhagen sadegh.talebi@di.ku.dk

## **Abstract**

In this paper, we analyze the sample complexities of learning the optimal stateaction value function  $Q^*$  and an optimal policy  $\pi^*$  in a finite discounted Markov decision process (MDP) where the agent has recursive entropic risk-preferences with risk-parameter  $\beta \neq 0$  and where a generative model of the MDP is available. We provide and analyze a simple model based approach which we call modelbased risk-sensitive Q-value-iteration (MB-RS-QVI) which leads to  $(\varepsilon, \delta)$ -PACbounds on  $||Q^* - Q^k||$ , and  $||V^* - V^{\pi_k}||$  where  $Q_k$  is the output of MB-RS-QVI after k iterations and  $\pi_k$  is the greedy policy with respect to  $Q_k$ . Both PAC-bounds have exponential dependence on the effective horizon  $\frac{1}{1-\gamma}$  and the strength of this dependence grows with the learners risk-sensitivity  $|\hat{\beta}|$ . We also provide two lower bounds which shows that exponential dependence on  $|\beta|\frac{1}{1-\gamma}$  is unavoidable in both cases. The lower bounds reveal that the PAC-bounds are tight in the parameters  $S, A, \delta, \varepsilon$  and that unlike in the classical setting it is not possible to have polynomial dependence in all model parameters.

## Introduction

In reinforcement learning (RL), the aim of the agent is to conventionally optimize the expected return, which is defined in terms of a (discounted) sum of rewards [50]. In the majority of RL literature, the environment is modeled via the Markov Decision Process (MDP) framework [40], wherein efficient computation of an optimal policy, thanks to optimal Bellman equations, renders possible. However, as a risk-neutral objective, the expected return fails to capture the true needs of many high-stake applications arising in, e.g., medical treatment [19], finance [43, 9], and operations research [16]. Decision making in such applications must take into account the variability of returns, and risks thereof. To account for this, one may opt to to maximize a risk measure of the return distribution, while another approach could be to consider the entire distribution of return, as is done under the distributional RL framework [8] that has received much attention over the last decade.

Within the first approach, the risk is quantified via concave risk measures, which yield well-defined mathematical optimization frameworks. Notably, they include value-at-risk (VaR), Conditional VaR (CVaR) [44], entropic risk [23], and entropic VaR (EVaR) [2], all of which have been applied to a wide-range of scenarios. CVaR appears to be the most popular one used to model risk-sensitivity in MDPs [15, 10, 12, 7], mainly due to a delicate control it offers for the undesirable tail of return. Despite its popularity and rich interpretation, solving and learning MDPs with CVaR-defined objectives has rendered technically challenging [7]. This has been a key motivation of adopting weaker notions such as nested CVaR [5], at the expense of sacrificing the interpretability. Entropic risk, as another popular notion, has been considered for risk control in MDPs [11, 38, 22, 24, 20]. In the RL literature, it has been mainly considered for the undiscounted settings, despite the popularity of discounted MDPs. A notable exception is [22] that studies, among other things, planning in discounted MDPs under the entropic risk.

In this paper, we study risk-sensitive discounted RL where the agent's objective is formulated using the entropic risk measure. In discounted RL, where future rewards are discounted by a factor of  $\gamma$ , one may identify two main approaches to apply the entropic risk to sequence of rewards,  $(r_t)_{t>0}$ , collected by an RL agent. The first and most intuitive one, which we may call the *non-recursive* approach, consists in directly applying the entropic risk functional to return  $\sum_{t=0}^{\infty} \gamma^t r_t$ . The other approach, which we may call the *recursive* approach, works by applying the risk functional at every step t (see Section 2 for details). The non-recursive approach (e.g., [22]), while being most intuitive, has several drawbacks; e.g., the optimal policy might not be time-consistent (see [25]). In contrast, the recursive approach yields a form of Bellman optimality equation, which is key to developing learning algorithms with provable sample complexity. Therefore, we restrict attention to the RL problem defined using recursive risk-preferences.

#### 1.1 **Main Contributions**

We study risk-sensitive RL in finite discounted MDPs under the recursively applied entropic risk measure, assuming that the agent is given access to a generative model of the environment. The agent's learning performance is assessed via sample complexity defined as the total number T of samples needed to learn, for input  $(\varepsilon, \delta)$ , either an  $\varepsilon$ -optimal policy (which we call policy learning), or an  $\varepsilon$ -close approximation (in max-norm) to the optimal Q-value (which we call Q-value learning), with probability exceeding  $1 - \delta$ .

Specifically, we make the following contributions: We present an algorithm called Model-Based Risk-Sensitive Q-Value Iteration (MB-RS-QVI), a model-based RL algorithms for the RL problem considered, which is derived using a simple plug-in estimator. It is provably sample efficient, despite its simple design. Notably, we report sample complexity bounds on the performance of (MB-RS-

QVI) under Q-value learning (Theorem 1) scaling as 
$$O\left(\frac{SA}{\varepsilon^2(1-\gamma)^2}\left(\frac{e^{i\gamma+1-\gamma-1}}{|\beta|}\right)^2\right)$$
 and policy learning

QVI) under Q-value learning (Theorem 1) scaling as  $\widetilde{O}\left(\frac{SA}{\varepsilon^2(1-\gamma)^2}(\frac{e^{|\beta|\frac{1}{1-\gamma}}-1}{|\beta|})^2\right)$  and policy learning (Theorems 2 and 3), scaling as  $\widetilde{O}\left(\frac{SA}{\varepsilon^2(1-\gamma)^4}(\frac{e^{|\beta|\frac{1}{1-\gamma}}-1}{|\beta|})^2\right)$ , and  $\widetilde{O}\left(\frac{S^2A}{\varepsilon^2(1-\gamma)^2}\frac{e^{|\beta|\frac{1}{1-\gamma}}}{|\beta|^2}\right)$  in any discounted MDP with S states, A action, and discount factor  $\gamma$ , with  $\widetilde{O}$  hiding logarithmic factors. Here,  $\beta$  denotes the risk parameter (see Section 2 for precise definition), where  $\beta > 0$  (resp.  $\beta < 0$ ) corresponds to a risk-averse (resp. risk-seeking) agent. An interesting property of these bounds is the exponential dependence on the effective horizon  $\frac{1}{1-\gamma}$ , which is absent in the conventional (riskneutral) RL, wherein  $\beta = 0$ . Another key contribution of this paper is to derive the first, to our knowledge, sample complexity lower bounds for the discounted RL under entropic risk measure. Our lower bounds establish that for Q-value learning (Theorem 4) and policy learning (Theorem 5), one needs at least  $\widetilde{\Omega}\left(\frac{(S-2)A}{\varepsilon^2}\frac{e^{|\beta|}\frac{1}{1-\gamma}-3}{|\beta|^2}\right)$  and  $\widetilde{\Omega}\left(\frac{(S-2)(A-1)}{\varepsilon^2}\frac{e^{|\beta|}\frac{1}{1-\gamma}-3}{|\beta|^2}\right)$  samples, respectively, to come  $\varepsilon$ -close to optimality. Interestingly, the derived lower bounds assert that exponential dependence on  $\frac{|\beta|}{1-\gamma}$  in both sample complexities are unavoidable thus making the risk-sensitive setting

fundamentally harder than the classical risk-neutral setting.

#### 1.2 Related Work

Finite-sample guarantees for risk-neutral RL. There is a large body of papers studying provably-sample efficient learning algorithms in discounted MDPs. These papers consider a variety of settings, including the generative setting [26, 21, 1], the offline (or batch) setting [41, 33], and the online setting [48, 31]. In particular, in the generative setting – which we consider in this paper - some notable developments include, but not limited to [27, 21, 1, 32, 46, 52, 14]. Among these, [21, 1, 46, 32] present algorithms attaining minimax-optimal sample complexity bounds, although some of these result do not cover the full range of  $\varepsilon$ . We also mention a line of work, comprising e.g. [39, 45], that investigate discounted RL in the generative setting but under distributional robustness. Let us also remark that some recent works – notably [56, 57] – study sample complexity of average-reward MDPs in the generative setting.

Finally, we mention that some studies consider adaptive sampling in the generative setting to account for the heterogeneity across the various state-action pairs in the MDP; see, e.g., [3, 54]. This line of works stand in contrast to the papers cited above that strive for optimizing the performance, in the worst-case sense, via uniformly sampling various state-action pairs.

**Risk-sensitive RL.** There exists a rich literature on decision making under a risk measure. We refer to [42, 35, 29] for some developments in bandits, where the performance is assessed via regret. Extensions to episodic MDPs were pursued in a recent line of work, including [20, 34, 24], which establish near-optimal guarantees on regret. The two lines of work (on bandits and episodic MDPs) constitute the majority of work on risk-sensitive RL, thoroughly studying a variety of risk measures including CVaR, entropic risk, and entropic value-at-risk [2, 49].

Among the various risk measures, CVaR has arguably received a great attention; see, for instance, [17, 18, 13]. For instance, [18, 13] study episodic RL in the regret setting and present algorithms with sub-linear regret in tabular MDPs ([18]) and under function approximation ([13]), while [17] investigates the sample complexity analysis in the generative setting (similar to ours), We also provide a pointer to some work which deal with a class of measures called coherent risk measures that include CVaR as a special case. In this category, we refer to, for instance, [51], which studies policy gradient algorithms, and to [30] that investigates regret minimization in the episodic setting and with function approximation. In the case of infinite-horizon setting, the entropic risk measure is mostly considered in the undiscounted (i.e., the average-reward) setting; examples include [11, 38, 36, 37]. In contrast, little attention is paid to the discounted setting, which is mainly due to technical difficulties caused by discounting. The work [22] is among the few papers investigating solving discounted MDPs under the entropic risk.

**Notations.** For  $n \in \mathbb{N}$ , let  $[n] := \{1, \dots, n\}$ .  $1_A$  denotes the indicator function of an event A. Given a set  $\mathcal{X}$ ,  $\Delta(\mathcal{X})$  denotes the probability simplex over  $\mathcal{X}$ . We use the convention that  $\|\cdot\| := \|\cdot\|_{\infty}$  and explicitly use the subscript  $\|\cdot\|_p$  when using p-norms for which  $1 \le p < \infty$ . We use  $Z = \mathcal{S} \times \mathcal{A}$  to denote the set of all state-action pairs.

## 2 Background

## 2.1 Markov Decision Processes

We write the 6-tuple  $M=(\mathcal{S},\mathcal{A},P,R,\gamma,\beta)$  as a finite, discounted infinite-horizon Markov decision process (MDP), where  $\mathcal{S}=\{1,2,...,S\}$  is the finite state space of size  $S:=|\mathcal{S}|,\,\mathcal{A}=\{1,2,...,A\}$  is the finite action space of size  $A:=|\mathcal{A}|,\,P:\mathcal{S}\times\mathcal{A}\to\Delta(\mathcal{S})$  is the transition probability function,  $R:\mathcal{S}\times\mathcal{A}\to[0,1]$  is the deterministic reward function,  $\gamma\in(0,1)$  is the discount factor, and  $\beta\neq 0$  is the risk-parameter. The agent interacts with the MDP M as follows. At the beginning of the process, M is in some initial state  $s_0\in\mathcal{S}$ . At each time  $t\geq 0$ , the agent is in state  $s_t\in\mathcal{S}$  and decides on an action  $a_t\in\mathcal{A}$  according to some rule. The MDP generates a reward  $r_t:=R(s_t,a_t)$  and a next-state  $s_{t+1}\sim P(\cdot|s_t,a_t)$ . The MDP moves to  $s_{t+1}$  when the next time slot begins, and this process continues ad infinitum. This process yields a growing sequence  $(s_t,a_t,r_t)_{t\geq 0}$ . The agent's goal is to maximize an objective function, as a function of the reward values  $(r_t)_{t\geq 0}$  which involves the two parameters  $\gamma$  and  $\beta$ ; in addition to the discount factor  $\gamma$  that makes future rewards less valuable than the present ones, the risk-parameter  $\beta$  quantifies to what degree the agent seeks or avoids strategies that have more variability in the rewards obtained over time. To concretely define the objective function of the agent, we shall introduce some necessary concepts.

## 2.2 Entropic Risk Preferences

The entropic risk preferences is rooted in expected utility theory. Consider for  $\beta \neq 0$  the class of utility functions  $u(t) = \frac{1}{\beta}(1-e^{-\beta t})$  defined for  $t \in \mathbb{R}$ . The utility u is supposed to describe the preferences of some economic agent in the form of how much utility u(t) she derives from some monetary quantity  $t \in \mathbb{R}$ . For any bounded random variable  $X \in L^{\infty}(\Omega, \mathcal{F}, \mathbb{P})$ , the certainty equivalent  $u^{-1}(\mathbb{E}[u(X)]) = \frac{-1}{\beta}\log(\mathbb{E}[e^{-\beta X}])$  expresses the amount of money that would give the same utility as that of entering in the bet given by the random variable X. We thus define the functional  $\rho: L^{\infty}(\Omega, \mathcal{F}, \mathbb{P}) \to \mathbb{R}$  by

$$\rho(X) = -\frac{1}{\beta} \log(\mathbb{E}[e^{-\beta X}]).$$

We note that there does not seem to be consensus on wether the functional is parametrized with  $\beta$  of  $-\beta$ . We follow this convention considering its widespread use in the actuarial literature [4], but

we remark that some other lines of work appear to prefer the other parametrization. For the case  $\beta \to 0$ , we recover the risk-neutral case that is simply the expectation:  $\lim_{\beta \to 0} \rho(X) = \mathbb{E}[X]$ . It is straightforward to see that  $\rho$  admits the following properties:

$$\rho(X) \le \rho(Y), \quad \text{for any } X \le Y,$$
(1)

$$\rho(c) = c, \quad \text{for any } c \in \mathbb{R}, \tag{2}$$

$$\rho(X) \le \mathbb{E}[X], \quad \text{for } \beta > 0,$$
(3)

$$\rho(X) \ge \mathbb{E}[X] \quad \text{for } \beta < 0, \tag{4}$$

where properties (3)-(4) follow from Jensen's inequality. Using  $\rho$  as a measure of the preference for different random variables, it follows directly from (2)-(4) that  $\rho(X) \leq \rho(\mathbb{E}[X])$  for  $\beta > 0$  and that  $\rho(X) \geq \rho(\mathbb{E}[X])$  for  $\beta < 0$ . If further shows that  $\beta > 0$  is associated with risk-aversion, while  $\beta < 0$  is associated with risk-seeking behavior.

Applying risk in a stochastic dynamic process can be done in several ways and is thus more complicated than for a single-period problem. Two approaches to this end exist in the literature. The first and most intuitive one, which is often called the *static* or *non-recursive* approach, is to apply the functional to the total discounted sum of rewards  $\rho(\sum_{t=0}^{\infty} \gamma^t r_t)$ , which is well-defined under the bounded rewards assumption, i.e.,  $r_t \in [0,1]$  for all t. The other approach, which we may call the *recursive* approach, works by applying the risk at every step t where we give the details below. The non-recursive approach is probably most intuitive but has several drawbacks. Even though there is no obvious optimality equation for the non-recursive case, a solution to the planning problem has been proposed in [22]. In comparison, the planning problem is more straightforward with recursive risk-preferences due to the availability of an optimality equation that allows for simple value-iteration type algorithms. The recursive approach also guarantees the existence of an optimal stationary deterministic policy whereas with non-recursive risk preferences the optimal policy might not be time-consistent (see [25]). In this paper, we study the problem with recursive risk-preferences.

## 2.3 Value Function and Q-function

A bit of notation is required in order to define the state value function V (henceforth V-function) and state-action value function Q (henceforth Q-function) of a policy.

We follow the approach of [4] and [6] but since none of their cases include our  $\beta>0$  case and also only cover V-functions, we give in Section C the full setup with history-dependent policies as well as a full definition of the V and Q-functions, and prove existence of a stationary optimal policy and show that the value functions of this policy satisfy a Bellman optimality equation and similarly that value functions of any policy satisfy a Bellman recursion relation. We give an outline here that only deals with stationary policies, which is justified by the results of Section C.

Let  $v \in \mathbb{R}^S$  and  $\pi : \mathcal{S} \to \mathcal{A}$  be a stationary deterministic policy. We then define  $\rho_{s,a} : \mathbb{R}^S \to \mathbb{R}$  by

$$\rho_{s,a}(v) = -\frac{1}{\beta} \log \left( \mathbb{E}_{s' \sim P(\cdot|s,a)} [e^{-\beta v(s')}] \right)$$
 (5)

and slighltly abusing the notation, we write  $\rho_{s,\pi}$  when  $a=\pi(s)$ , that is  $\rho_{s,\pi}:=\rho_{s,\pi(s)}$ . We then define the operator  $J_\pi:\mathbb{R}^S\to\mathbb{R}^S$  given by  $J_\pi(v)(s)=r(s,\pi(s))+\gamma\rho_{s,\pi}(v)$ . The N-step total discounted utility  $J_N(s,\pi)$  is defined as applying  $J_\pi$  recursively N-times to the 0-map, that is  $J_N(s,\pi):=J_\pi^N(\mathbf{0})(s)$ . Note that the outer-most iteration corresponds to the immediate time-step. By properties (1)-(2) of  $\rho$ , it follows that  $J_N(s,\pi)$  is increasing in N and that  $J_N(s,\pi)\leq \frac{1}{1-\gamma}$ , so that the limit  $N\to\infty$  exists. This limit is considered the value of state s under the policy  $\pi$  of the agent:  $V^\pi(s)=\lim_{N\to\infty}J_N(s,\pi)$ .

The problem of the agent is then for all initial states  $s \in \mathcal{S}$  to find a policy  $\pi^*$  that solves  $J(s,\pi^*) = \sup_{\pi} J(s,\pi)$ . In [4], the authors consider a more general framework that is not restricted to finite MDPs or stationary policies; they prove that under some conditions —that are trivially fulfilled in the case of finite MDPs— there exists a stationary policy  $\pi^*$  that maximizes the state-value function for all states  $s \in \mathcal{S}$  simultaneously. We show that any optimal policy  $\pi^*$  that solves the agent satisfies the optimality equation:

$$V^*(s) = \max_{a \in \mathcal{A}} \left( R(s, a) - \frac{\gamma}{\beta} \log \left( \mathbb{E}_{s' \sim P(\cdot | s, a)} [e^{-\beta V^*(s')}] \right) \right), \tag{6}$$

where  $V^*$  is the optimal V-function. Also for any stationary deterministic policy  $\pi$ , we have the Bellman recursion:

$$V^{\pi}(s) = R(s, \pi(s)) - \frac{\gamma}{\beta} \log \left( \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} [e^{-\beta V^{\pi}(s')}] \right). \tag{7}$$

Since we are not interested only in the planning problem but also in learning, we also introduce the state-action value function Q. The approach is very similar to that of the value-function. Given  $\pi$ , we define the operator  $L_\pi:\mathbb{R}^S\to\mathbb{R}^{S\times A}$  as follows: for all  $v:\mathcal{S}\to\mathbb{R}$ ,  $L_\pi(v)(s,a)=R(s,\pi(s))+\gamma\rho_{s,\pi}(v)$ . We define the operator  $L:\mathbb{R}^S\to\mathbb{R}^{S\times A}$  for all  $v:\mathcal{S}\to\mathbb{R}$  as:  $L(v)(s,a)=R(s,a)+\gamma\rho_{s,a}(v)$ . We define the N-step total discounted utility of the state-action pair (s,a) under  $\pi$  as  $L_N(s,a,\pi):=(L\circ J_\pi^{N-1}(\mathbf{0}))(s,a)$  and the limit is denoted  $Q^\pi(s,a)\colon Q^\pi(s,a)=\lim_{N\to\infty}L_N(s,a,\pi)$ . Although the authors do not consider state-action value functions in their paper, repeating the arguments of [4] it suffices to consider stationary policies when wanting to solve  $\max_\pi Q^\pi(s,a)$  for all (s,a) and that the solution  $Q^*$  solves the optimality equation:

$$Q^*(s,a) = R(s,a) - \frac{\gamma}{\beta} \log \left( \mathbb{E}_{s' \sim P(\cdot|s,a)} \left[ e^{-\beta \max_{a'} Q^*(s',a')} \right] \right). \tag{8}$$

Similarly, it is clear that the Q-functions satisfy the Bellman recursion relations:

$$Q^{\pi}(s,a) = R(s,a) - \frac{\gamma}{\beta} \log \left( \mathbb{E}_{s' \sim P(\cdot|s,a)} [e^{-\beta V^{\pi}(s')}] \right). \tag{9}$$

#### 2.4 Learning Performance

We consider two types of RL algorithms  $\mathcal{U}$ , namely those that output a Q-function  $Q_T^{\mathcal{U}}: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$  and those that output a policy  $\pi_T^{\mathcal{U}}: \mathcal{S} \to \mathcal{A}$  using T transition samples. Note that any algorithm that outputs a Q-function also outputs a policy, namely the one obtained by acting greedily with respect to the Q-function. There are also ways to obtain a Q-function from a policy. There is however no canonical way to do this as the algorithm cannot simply output  $Q^{\pi^{\mathcal{U}}}$  since the algorithm does not have access to the true MDP. The way we evaluate the quality of an algorithm that outputs a Q-function is by  $\|Q^* - Q_T^{\mathcal{U}}\|$ . For an algorithm that instead outputs a policy, we evaluate the policy in terms of  $\|V^* - V^{\pi_T^{\mathcal{U}}}\|$ . Often we will suppress T from the notation.

**Definition 1**  $((\varepsilon, \delta)$ -correctness). We say that an algorithm  $\mathcal U$  that outputs a Q-function  $Q^{\mathcal U}$  is  $(\varepsilon, \delta)$ -correct on a set of MDPs  $\mathbb M$  if  $\mathbb P(\|Q^*-Q^{\mathcal U}\|\leq \varepsilon)\geq 1-\delta$  for all  $M\in\mathbb M$ . Similarly, we say that an algorithm  $\mathcal U$  that outputs a policy  $\pi^{\mathcal U}$  is  $(\varepsilon, \delta)$ -correct on a set of MDPs  $\mathbb M$  if  $\mathbb P(\|V^*-V^{\pi^{\mathcal U}}\|\leq \varepsilon)\geq 1-\delta$  for all  $M\in\mathbb M$ .

## 3 Model-Based Risk-Sensitive Q-Value Iteration

In this section we describe the model-based value iteration algorithm which aims at finding the optimal Q-function  $Q^*$ . We then give an upper bound on the total number of calls to the generative model needed in order for this algorithm to be  $(\varepsilon, \delta)$ —correct. The model based approach is based on working on an MDP, which may disagree with the true MDP because it does not use the true transition probabilities but an estimate of the transition functions obtained from n calls to each of the state-action pairs in  $\mathcal{S} \times \mathcal{A}$  as described in the algorithm below.

The model-based approach we describe is general in the sense that any oracle that for any  $\varepsilon > 0$  can find an  $\varepsilon$ -optimal policy can be used. We prove the existence of one such oracle in the form of a Q-value iteration very like the one from the classical risk-neutral setting. The proof is very similar to Part (a) in Theorem 3.1 in [4] but is nevertheless provided in Appendix F.

**Lemma 1** (Q-value iteration). Fix a map  $\pi: \mathcal{A} \to \mathcal{S}$ . We then define the operators  $\mathcal{T}^{\pi}, \mathcal{T}: \mathbb{R}^{S \times A} \to \mathbb{R}^{S \times A}$  which for  $f: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$  is given by

$$(\mathcal{T}f)(s,a) = R(s,a) + \frac{-\gamma}{\beta} \log \left( \sum_{s'} P(s'|s,a) e^{-\beta \max_a f(s',a)} \right)$$
$$(\mathcal{T}^{\pi}f)(s,a) = R(s,a) + \frac{-\gamma}{\beta} \log \left( \sum_{s'} P(s'|s,a) e^{-\beta f(s',\pi(s'))} \right)$$

#### **Algorithm 1:** Model estimation

```
Input: Generative model P
   Output: Model estimate \hat{P}
1 Function EstimateModel(n):
        \forall (s,z) \in \mathcal{S} \times Z : m(s,z) = 0
2
        for each z \in Z do
3
             for i = 1, 2, ..., n do
 4
                  s \sim P(\cdot|z)
 5
                 m(s,z) := m(s,z) + 1
 6
 7
             \forall s \in \mathcal{S} : \widehat{P}(s, z) = \frac{m(s, z)}{n}
 8
        end
9
        return \widehat{P}
10
```

The operators  $\mathcal{T}$  and  $\mathcal{T}^{\pi}$  are  $\gamma$ -contractions with respect to the max-norm, i.e., for value-functions  $f_1$  and  $f_2$ , it holds that  $\|\mathcal{T}f_1 - \mathcal{T}f_2\| \leq \gamma \|f_1 - f_2\|$  and  $\|\mathcal{T}^{\pi}f_1 - \mathcal{T}^{\pi}f_2\| \leq \gamma \|f_1 - f_2\|$ .

The above lemma is the basis for the Q-value iteration algorithm:

```
Algorithm 2: RS-QVI(M, k)
```

```
Input: MDP M = (S, A, P, R, \gamma, \beta) and number of iterations k.
    Output: Estimate Q_k of optimal Q-function Q^*
 1 Initialization: \forall (s, a) \text{ set } Q(s, a) = 0
 2 for j = 0, 1, ..., k - 1 do
          for all s \in \mathcal{S} do
 3
               \pi_j(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q_j(s, a)
 4
               for all a \in \mathcal{A} do
 5
                    \mathcal{T}Q_j(s,a) = R(s,a) - \frac{\gamma}{\beta} \log(\mathbb{E}_{s' \sim P(\cdot|s,a)}[e^{-\beta Q_j(s,\pi_j(s))}])
 6
                    Q_{j+1}(s,a) = \mathcal{T}Q_i(s,a)
               end
 8
          end
10 end
11 \forall s \in \mathcal{S} : \pi_k(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q_k(s, a)
12 return Q_k and \pi_k
```

The next lemma shows that if we choose k large enough in the RS-QVI algorithm, we can obtain  $Q_k$  and  $V^{\pi_k}$  that are as close to  $Q^*$  and  $V^*$  as we desire. The proof is postponed to Appendix F.

**Lemma 2.** Fix  $\varepsilon > 0$ . Then there exists some  $k(\varepsilon)$  such that if the number of iterations in RS-QVI exceeds  $k(\varepsilon)$ , then the output of Algorithm 2 (RS-QVI) satisfies  $\|Q_k - Q^*\| < \varepsilon$  and  $\|V^{\pi_k} - V^*\| < \varepsilon$ .

Using Algorithm 2, we introduce the MB-RS-QVI algorithm, which consists in building an empirical model  $\widehat{M}=(\mathcal{S},\mathcal{A},\widehat{P},R,\gamma,\beta)$  via calling the generative model n times – namely,  $\widehat{P}=\texttt{EstimateModel}(n)$  – and then solving it via RS-QVI.

## 3.1 Analysis of MB-RS-QVI

With RS-QVI in place, we need a set of lemmas for the analysis of the sample-complexity of the model based RS-QVI algorithm.

An important result for the analysis is a risk-sensitive version of the simulation lemma [28, 48], which describes how different two Q-functions for the same policy are in two different MDPs that differ only slightly in their rewards and transition functions. The proof is postponed to Appendix F.

**Lemma 3** (Simulation Lemma with Entropic Risk). Consider two MDPs  $M_1 = (S, A, P_1, R, \gamma, \beta)$  and  $M_2 = (S, A, P_2, R, \gamma, \beta)$  differing only in their transition functions. Fix a stationary policy  $\pi$ , and let  $Q_1^{\pi}$  and  $Q_2^{\pi}$  be respective Q-functions of  $\pi$  in  $M_1$  and  $M_2$ . If  $\max_{(s,a)\in\mathcal{S}\times\mathcal{A}}\|P_1(\cdot|s,a)-P_2(\cdot|s,a)\|_1 \leq \tau$ . Then, it holds that

$$\|Q_1^{\pi} - Q_2^{\pi}\| \leq \begin{cases} \frac{\gamma}{1-\gamma} \frac{e^{|\beta|\frac{1}{1-\gamma}}}{|\beta|} \max_{s,a} \left| \sum_{s' \in \mathcal{S}} [P_1(s'|s,a) - P_2(s'|s,a)] e^{-|\beta|(\frac{1}{1-\gamma} - V_1(s'))} \right| & \beta < 0 \\ \frac{\gamma}{1-\gamma} \frac{e^{|\beta|\frac{1}{1-\gamma}}}{|\beta|} \max_{s,a} \left| \sum_{s' \in \mathcal{S}} [P_1(s'|s,a) - P_2(s'|s,a)] e^{-|\beta|V_1(s'))} \right| & \beta > 0 \\ \frac{\gamma}{1-\gamma} \frac{e^{|\beta|\frac{1}{1-\gamma}}}{|\beta|} \max_{s,a} \sum_{s' \in \mathcal{S}} |P_1(s'|s,a) - P_2(s'|s,a)| & \beta \neq 0. \end{cases}$$

It then follows that if for some  $\varepsilon>0$  it holds that any of the  $\max_{s,a}$ -expressions is smaller than  $\varepsilon^{\frac{1-\gamma}{\gamma}}|\beta|e^{-|\beta|\frac{1}{1-\gamma}}$ , then  $\|Q_1^\pi-Q_2^\pi\|\leq\varepsilon$  and ensuring this with high probability is a matter of invoking an appropriate concentration inequality. Using the decompositions from Lemma 4, we get for any  $s\in\mathcal{S}$  that

$$V^{\pi_k}(s) \ge V^*(s) - \|V^{\pi_k} - \widehat{V}^{\pi_k}\| - \|\widehat{V}^{\pi^*} - V^*\| - \|\widehat{V}^{\pi_k} - \widehat{V}^*\|, \tag{10}$$

and similarly we get for any state-action pair  $(s, a) \in \mathcal{S} \times \mathcal{A}$  that

$$Q_k(s,a) \ge Q^*(s,a) - \|\widehat{Q}^{\pi^*} - Q^*\| - \|Q_k - \widehat{Q}^*\|.$$
(11)

The last of the distances on the right-hand side of (10) and (11) can be made arbitrarily small by any optimization oracle for the problem. One such is value-iteration using the model MDP as demonstrated by Lemma 2. Making these terms small enough is thus purely a computational matter and not a statistical one. We thus focus on bounding the remaining distances on the right-hand sides.

**Theorem 1.** There exists a universal constant c such that for any  $\varepsilon > 0$ ,  $\delta \in (0,1)$  and any MDP M with S states and A actions, if the learner is allowed to make

$$T = c \frac{SA}{\varepsilon^2 (1 - \gamma)^2} \left( \frac{e^{|\beta| \frac{1}{1 - \gamma}} - 1}{|\beta|} \right)^2 \log \left( \frac{SA}{\delta} \right)$$
 (12)

calls to the generative model, then  $\mathbb{P}(\|Q^* - Q_k\| \le \varepsilon) \ge 1 - \delta$ .

*Proof.* For any  $\varepsilon>0$ , we can get  $\|Q_k-\widehat{Q}^*\|<\varepsilon/2$  using enough iterations of the optimization oracle by Lemma 2. The term  $\|\widehat{Q}^{\pi^*}-Q^*\|$  can also be made smaller than  $\varepsilon/2$  by the simulation lemma if either  $\max_{s,a} \left|\sum_{s'\in\mathcal{S}} [P_1(s'|s,a)-P_2(s'|s,a)]e^{-|\beta|(\frac{1}{1-\gamma}-V_1(s'))}\right|<\tau$  or

$$\max_{s,a} \left| \sum_{s' \in \mathcal{S}} [P_1(s'|s,a) - P_2(s'|s,a)] e^{-|\beta|V_1(s'))} \right| < \tau$$
, where  $\tau = \frac{\varepsilon}{2} \left[ \frac{\gamma}{1-\gamma} \frac{e^{|\beta|\frac{1}{1-\gamma}}}{|\beta|} \right]^{-1}$  which

can be ensured with probability larger than  $1-\delta$  by picking  $N=\frac{(1-e^{-\beta\frac{1}{1-\gamma}})^2}{2\tau^2}\log(2SA/\delta)$ . Using that the total calls to the generative model is T=SAN and substituting in the value for  $\tau$ , we can ensure for all (s,a) that  $Q_k(s,a)>Q^*(s,a)-\varepsilon$  with probability larger than  $1-\delta$  by using a total number of samples

$$T = 3 \frac{SA}{\varepsilon^2 (1 - \gamma)^2} \left( \frac{e^{|\beta| \frac{1}{1 - \gamma}} - 1}{|\beta|} \right)^2 \log \left( \frac{SA}{\delta} \right). \tag{13}$$

**Theorem 2.** There exists a universal constant c such that for any  $\varepsilon > 0, \delta \in (0,1)$ , and any MDP M with S states and A actions if the learner is allowed to make

$$T = c \frac{SA}{\varepsilon^2 (1 - \gamma)^4} \left( \frac{e^{|\beta| \frac{1}{1 - \gamma}} - 1}{|\beta|} \right)^2 \log \left( \frac{SA}{\delta} \right)$$
 (14)

calls to the generative model then  $\mathbb{P}(\|V^* - V^{\pi_k}\| \leq \varepsilon) \geq 1 - \delta$ 

*Proof.* The proof follows immediately from Theorem 1 and Theorem 7 with the choice of c=6.  $\Box$ 

By taking the limit  $\beta \to 0$  leads to PAC-bounds in the risk-neutral case, but the resulting Q-value learning PAC bound is worse by a factor of  $\frac{1}{1-\gamma}$  [21] and the policy learning PAC-bound is worse by a factor of  $\frac{1}{(1-\gamma)^3}$  [1].

Since the derivations so far lead to a PAC-bound on policy learning that is worse by a factor of  $\frac{1}{(1-\gamma)^2}$  compared to the PAC bound on Q-value learning, we now provide an alternative derivation that leads to a PAC-bound with the same dependence on horizon. The bound achieved based on the Weissman bound however comes with a worse dependence on S, but in cases with long effective horizons and small state-space, this bound might be tighter. Another deficiency of this alternative bound is that it is also less interpretable in the sense that it explodes for  $\beta \to 0$ .

**Theorem 3.** There exists a universal constant c such that for any  $\varepsilon > 0$ ,  $\delta \in (0,1)$ , and any MDP M with S states and A actions, if the learner is allowed to make

$$T = c \frac{SA(S + \log(SA/\delta))}{\varepsilon^2 (1 - \gamma)^2} \frac{e^{2|\beta| \frac{1}{1 - \gamma}}}{|\beta|^2}$$
(15)

calls to the generative model, then  $\mathbb{P}(\|V^* - V^{\pi_k}\| \leq \varepsilon) \geq 1 - \delta$ .

*Proof.* By Lemma 2, we can pick k large enough so that  $\|\widehat{V}^{\pi_k} - \widehat{V}^*\| \le \varepsilon/3$ . By the union bound we have that the two terms  $\|V^{\pi_k} - \widehat{V}^{\pi_k}\|$  and  $\|\widehat{V}^{\pi^*} - V^*\|$  can simultaneously be made smaller than  $\varepsilon/3$  by sampling each state-action pair  $N = 8 \frac{S + \log(\frac{2SA}{\delta})}{\tau^2}$  times where  $\tau = \frac{\varepsilon}{3} \left[ \frac{\gamma}{1 - \gamma} \frac{e^{|\beta|} \frac{1}{1 - \gamma}}{|\beta|} \right]^{-1}$  such that if we sample

$$T = 73 \left[ S + \log \left( \frac{SA}{\delta} \right) \right] \frac{SA}{\varepsilon^2 (1 - \gamma)^2} \frac{e^{2|\beta| \frac{1}{1 - \gamma}}}{|\beta|^2}$$
 (16)

times in total, then we can ensure that for all s simultaneously it holds that  $V^{\pi_k}(s) > V^* - \varepsilon$ .  $\square$ 

**Remark 1.** State-of-the-art techniques for proving optimal upper bounds in the risk-neutral case are critically exploiting linearity of the expectation operator and are thus not readily available in the risk-sensitive case due to the non-linearity of the entropic risk measure.

## 4 Sample Complexity Lower Bounds

In this section, we provide two sample complexity lower bounds. The first one, presented in Theorem 4, concerns the sample complexity of learning the optimal Q-value function  $Q^*$ , whereas the second one, in Theorem 5, is on learning an optimal policy  $\pi^*$ . The proofs are both postponed to Appendix H.

**Theorem 4** (Lower bound for learning  $Q^*$ ). There exist constants  $c_1, c_2 > 0$  such that for any RL algorithm  $\mathcal{U}$  that outputs a Q-function  $Q^{\mathcal{U}}$  and any  $\delta \in (0, \frac{1}{4})$  and  $\varepsilon \in (0, \frac{1}{40} \frac{\gamma}{|\beta|} (1 - e^{-|\beta| \frac{1}{1-\gamma}}))$ , the following holds: if the total number T of transitions satisfies

$$T \le \frac{(S-2)A\gamma^2}{c_1\varepsilon^2} \frac{\left(e^{|\beta|\frac{1}{1-\gamma}} - 3\right)}{|\beta|^2} \log\left(\frac{(S-2)A}{c_2\delta}\right),\,$$

then there is some MDP M with S states and A actions for which  $\mathbb{P}(\|Q_M^* - Q_T^{\mathcal{U}}\| > \varepsilon) \geq \delta$ .

**Theorem 5** (Lower bound for learning  $\pi^*$ ). There exist constants  $c_1, c_2 > 0$  such that for any RL algorithm  $\mathcal U$  that outputs a policy  $Q^{\mathcal U}$  and any  $\delta \in (0, \frac{1}{4})$  and  $\varepsilon \in (0, \frac{1}{50} \frac{\gamma}{|\beta|} (1 - e^{-|\beta| \frac{1}{1-\gamma}}))$ , it holds that if the total number T of transitions satisfies

$$T \le \frac{(S-2)(A-1)\gamma^2}{c_1\varepsilon^2} \frac{\left(e^{|\beta|\frac{1}{1-\gamma}} - 3\right)}{|\beta|^2} \log\left(\frac{(S-2)}{c_2\delta}\right),$$

then there is some MDP M with S states and A actions for which  $\mathbb{P}(\|V_M^* - V^{\pi_T^{\mathcal{U}}}\| > \varepsilon) \geq \delta$ .

The lower bounds in Theorems 4-5 establish that an exponential dependence of the sample complexity on the effective horizon  $\frac{1}{1-\gamma}$  is unavoidable; more precisely, they assert that a scaling of  $e^{|\beta|\frac{1}{1-\gamma}}$  is unavoidable in both Q-learning and policy-learning. Recalling the sample complexity bounds of MB-RS-QVI (Theorems 1-2), we observe a similar exponential dependence. However, there remains a gap of order  $\frac{1}{(1-\gamma)^2}e^{|\beta|\frac{1}{1-\gamma}}$ . This remains as an interesting open question as to whether closing this gap can be done via a more elegant analysis of MB-RS-QVI or it calls for more novel algorithmic ideas, but in any case the lower bounds show that risk-sensitive agents face a fundamentally harder problem than risk-neutral agents where the sample complexity is polynomial in  $\frac{1}{1-\gamma}$ . The proofs of lower bounds are given in Section H in the appendix, but we briefly sketch the ideas below for the case where the algorithm outputs a Q-function; the other case is proven using quite similar ideas.

We consider a class of hard-to-learn MDPs. In the class there are two absorbing states  $s^G$  and  $s^B$  where in  $s^G$  the agent always receives a reward of R=1 and in  $s^B$  the agent always receives a reward of R=0 irrespective of the action taken. For all other state-action pairs z the reward is zero and the only possible transition is to either of the states  $s^G$  and  $s^B$ .  $P(s^G|z)=q$  and  $P(s^B|z)=1-q$ , for some q>0. This construction critically allows us to calculate explicitly  $Q^*(z)$  for a given parameter q and for two different MDPs  $M_0$ ,  $M_1$  in the class where  $q_0=p$  and  $q_1=p+\alpha$  for appropriately chosen values of p and  $\alpha$  we are able to ensure that  $Q^*_{M_1}(z)-Q^*_{M_0}(z)>2\varepsilon$  which means that any specific algorithmic output  $Q^U(z)$  cannot be  $\varepsilon$ -close to both  $Q^*_{M_1}(z)$  and  $Q^*_{M_0}(z)$ . We then show by a likelihood ratio argument that any algorithm  $\mathcal U$  that is  $(\varepsilon,\delta)$ -correct on  $M_0$ , i.e. that  $\mathbb P_0(|Q^*_{M_0}(z)-Q^U(z)|\leq \varepsilon)>\delta$ , will also satisfy that  $\mathbb P_1(|Q^*_{M_0}(z)-Q^U(z)|\leq \varepsilon)>\delta$  provided that the algorithm does not try out z enough times on  $M_0$  and exactly because  $Q^*_{M_1}(z)-Q^*_{M_0}(z)>2\varepsilon$ , the event  $\{|Q^*_{M_0}(z)-Q^U(z)|\leq \varepsilon\}$  is disjoint from the event on being  $\varepsilon$ -close to  $Q^*_{M_1}$ . The final part of the proof is to exploit that the different state-action pairs contain no information about each other which allows for an independence argument for the estimation of  $Q^U(z)$  and  $Q^U(z')$  for  $z\neq z'$ . While doing this analysis, we fix an inaccuracy in the proof of Lemma 17 in [21] that arises where they lower-bound the likelihood ratio of two Bernoulli random variables with biases  $p\geq \frac{1}{2}$  and  $p+\alpha$  on a high probability event. We also mention that we extend the result to hold for  $p<\frac{1}{2}$ .

For policy learning we consider almost the same class of MDPs but augment with a known state  $a_0$  that is used in the analysis.

**Remark 2.** It is worth remarking that the best lower bound in the risk-neutral setting is derived in [21] using a richer construction than above. However, with a risk-sensitive learning objective, the optimal state-action value function in the construction of [21] does not admit an analytical solution, which is needed for the delicate tuning of the transition probabilities.

## 5 Conclusion and Future Works

We have studied the sample-complexity of learning the optimal Q-function and that of learning an optimal policy in finite discounted MDPs, where the agent has recursive risk-preferences given by the entropic risk measure and has access to a generative model. We introduced an algorithm, called MB-RS-QVI, and derived PAC-type bounds on its sample complexity for both learning which have derived bounds from analyzing the MB-RS-QVI algorithm that uses the model given by the plugin estimator from samples generated by a simulator. We also derive lower bounds. The upper bounds show that PAC-learning is possible but the lower bounds show that dependence on  $e^{|\beta|} \frac{1}{1-\gamma}$  is unavoidable and thus that learning is fundamentally harder for risk-senstive agents relative to risk-neutral agents. The bounds that we derive on the sample complexity of learning the optimal Q-value are of order

$$\mathcal{O}\bigg(\log(SA/\delta)\bigg)\frac{SA}{\varepsilon^2(1-\gamma)^2}\bigg(\frac{e^{|\beta|\frac{1}{1-\gamma}}-1}{|\beta|}\bigg)^2\bigg), \qquad \Omega\bigg(\frac{(S-2)A}{\varepsilon^2}\frac{e^{|\beta|\frac{1}{1-\gamma}}-3}{|\beta|^2}\log((S-2)A/\delta)\bigg)$$
(17)

while the bounds we derive on the sample complexity of learning an optimal policy are of order

$$\mathcal{O}\bigg(\log(SA/\delta)\bigg)\frac{SA}{\varepsilon^2(1-\gamma)^4}\bigg(\frac{e^{|\beta|\frac{1}{1-\gamma}}-1}{|\beta|}\bigg)^2\bigg), \qquad \Omega\bigg(\frac{(S-2)(A-1)}{\varepsilon^2}\frac{e^{|\beta|\frac{1}{1-\gamma}}-3}{|\beta|^2}\log((S-2)/\delta)\bigg)$$
(18)

where we also give the alternative bound of order

$$\mathcal{O}\left(\left(S + \log(SA/\delta)\right) \frac{SA}{\varepsilon^2 (1-\gamma)^2} \frac{e^{2|\beta|\frac{1}{1-\gamma}}}{|\beta|^2}\right) \tag{19}$$

which might be tighter in cases of long horizon and small state space. These constitute the first bounds, to our knowledge, on the sample complexities of entropic risk-sensitive agents in the discounted MDP setting. The upper and lower bounds derived in this paper leave open gaps in  $\frac{1}{1-\gamma}$ . Since the constructions in the lower bounds are not the ones used in the tightest lower bounds derived in the risk-neutral setting, one possibility is that the lower bounds can be improved by considering a more carefully chosen set of hard-to-learn MDPs with the challenge being to control the gap in V-values or Q-values under different parameters. Also since the plugin-estimator model-based QVI algorithm is provably optimal in the risk-neutral setting, we believe that this might also be the case for risk-sensitive agents but that more tools are needed to develop a more careful analysis. Another future direction is that of developing model-free algorithms for this setting and analyzing their statistical efficiency. Another interesting research direction is to consider function approximation, as in [55]. As other future directions, one may consider more complicated RL settings such as offline RL [41], where data is collected under a fixed (but unknown) behavior policy, and online RL [48, 31], where the agent's learned policy impacts the data collection process. And finally one may also consider the problem where the learner have non-recursive risk-preferences instead.

## Acknowledgments

The authors would like to acknowledge the support from Independent Research Fund Denmark, grant number 1026-00397B.

## References

- [1] Alekh Agarwal, Sham Kakade, and Lin F Yang. "Model-based reinforcement learning with a generative model is minimax optimal". In: *Conference on Learning Theory*. PMLR. 2020, pp. 67–83.
- [2] Amir Ahmadi-Javid. "Entropic value-at-risk: A new coherent risk measure". In: *Journal of Optimization Theory and Applications* 155 (2012), pp. 1105–1123.
- [3] Aymen Al Marjani and Alexandre Proutiere. "Adaptive sampling for best policy identification in Markov decision processes". In: *International Conference on Machine Learning*. PMLR. 2021, pp. 7459–7468.
- [4] Hubert Asienkiewicz and Anna Jaśkiewicz. "A note on a new class of recursive utilities in Markov decision processes". In: *Applicationes Mathematicae* 44 (2017), pp. 149–161.
- [5] Nicole Bäuerle and Alexander Glauner. "Markov decision processes with recursive risk measures". In: *European Journal of Operational Research* 296.3 (2022), pp. 953–966.
- [6] Nicole Bäuerle and Anna Jaśkiewicz. "Markov decision processes with risk-sensitive criteria: An overview". In: Mathematical Methods of Operations Research 99.1 (2024), pp. 141–178.
- [7] Nicole Bäuerle and Jonathan Ott. "Markov decision processes with average-value-at-risk criteria". In: *Mathematical Methods of Operations Research* 74 (2011), pp. 361–379.
- [8] Marc G Bellemare, Will Dabney, and Mark Rowland. Distributional reinforcement learning. MIT Press, 2023.
- [9] Tomasz R Bielecki and Stanley R Pliska. "Risk-sensitive dynamic asset management". In: *Applied Mathematics and Optimization* 39 (1999), pp. 337–360.
- [10] Lorenzo Bisi et al. "Risk-averse policy optimization via risk-neutral policy optimization". In: *Artificial Intelligence* 311 (2022), p. 103765.

- [11] Vivek S Borkar and Sean P Meyn. "Risk-sensitive optimal control for Markov decision processes with monotone cost". In: *Mathematics of Operations Research* 27.1 (2002), pp. 192–209.
- [12] Daniel Brown, Scott Niekum, and Marek Petrik. "Bayesian robust optimization for imitation learning". In: Advances in Neural Information Processing Systems 33 (2020), pp. 2479–2491.
- [13] Yu Chen et al. "Provably Efficient Iterated CVaR Reinforcement Learning with Function Approximation and Human Feedback". In: *The Twelfth International Conference on Learning Representations*. 2024.
- [14] Zaiwei Chen et al. "Finite-sample analysis of contractive stochastic approximation using smooth convex envelopes". In: *Advances in Neural Information Processing Systems* 33 (2020), pp. 8223–8234.
- [15] Yinlam Chow and Mohammad Ghavamzadeh. "Algorithms for CVaR optimization in MDPs". In: *Advances in neural information processing systems* 27 (2014).
- [16] Erick Delage and Shie Mannor. "Percentile optimization for Markov decision processes with parameter uncertainty". In: *Operations research* 58.1 (2010), pp. 203–213.
- [17] Zilong Deng, Simon Khan, and Shaofeng Zou. "Near-Optimal Sample Complexity for Iterated CVaR Reinforcement Learning with a Generative Model". In: *AISTATS*. 2025.
- [18] Yihan Du, Siwei Wang, and Longbo Huang. "Provably Efficient Risk-Sensitive Reinforcement Learning: Iterated CVaR and Worst Path". In: *The Eleventh International Conference on Learning Representations*. 2023.
- [19] Damien Ernst et al. "Clinical data based optimal STI strategies for HIV: A reinforcement learning approach". In: *Proceedings of the 45th IEEE Conference on Decision and Control*. IEEE. 2006, pp. 667–672.
- [20] Yingjie Fei et al. "Risk-sensitive reinforcement learning: Near-optimal risk-sample tradeoff in regret". In: *Advances in Neural Information Processing Systems* 33 (2020), pp. 22384–22395.
- [21] Mohammad Gheshlaghi Azar, Rémi Munos, and Hilbert J Kappen. "Minimax PAC bounds on the sample complexity of reinforcement learning with a generative model". In: *Machine learning* 91 (2013), pp. 325–349.
- [22] Jia Lin Hau, Marek Petrik, and Mohammad Ghavamzadeh. "Entropic risk optimization in discounted MDPs". In: *International Conference on Artificial Intelligence and Statistics*. PMLR. 2023, pp. 47–76.
- [23] Ronald A Howard and James E Matheson. "Risk-sensitive Markov decision processes". In: *Management science* 18.7 (1972), pp. 356–369.
- [24] Xiaoyan Hu and Ho-fung Leung. "A tighter problem-dependent regret bound for risk-sensitive reinforcement learning". In: *International Conference on Artificial Intelligence and Statistics*. PMLR. 2023, pp. 5411–5437.
- [25] Stratton C Jaquette. "A utility criterion for Markov decision processes". In: *Management Science* 23.1 (1976), pp. 43–49.
- [26] Sham Machandranath Kakade. *On the sample complexity of reinforcement learning*. University of London, University College London (United Kingdom), 2003.
- [27] Michael Kearns and Satinder Singh. "Finite-sample convergence rates for Q-learning and indirect algorithms". In: *Advances in neural information processing systems* 11 (1998).
- [28] Michael Kearns and Satinder Singh. "Near-optimal reinforcement learning in polynomial time". In: *Machine learning* 49 (2002), pp. 209–232.
- [29] Najakorn Khajonchotpanya, Yilin Xue, and Napat Rujeerapaiboon. "A revised approach for risk-averse multi-armed bandits under CVaR criterion". In: *Operations Research Letters* 49.4 (2021), pp. 465–472.
- [30] Thanh Lam et al. "Risk-aware reinforcement learning with coherent risk measures and non-linear function approximation". In: *The Eleventh International Conference on Learning Representations*, 2022.
- [31] Tor Lattimore and Marcus Hutter. "Near-optimal PAC bounds for discounted MDPs". In: *Theoretical Computer Science* 558 (2014), pp. 125–143.
- [32] Gen Li et al. "Breaking the sample size barrier in model-based reinforcement learning with a generative model". In: *Advances in neural information processing systems* 33 (2020), pp. 12861–12872.

- [33] Gen Li et al. "Settling the sample complexity of model-based offline reinforcement learning". In: The Annals of Statistics 52.1 (2024), pp. 233–260.
- [34] Hao Liang and Zhiquan Luo. "Regret bounds for risk-sensitive reinforcement learning with lipschitz dynamic risk measures". In: *International Conference on Artificial Intelligence and Statistics*. PMLR. 2024, pp. 1774–1782.
- [35] Odalric-Ambrym Maillard. "Robust risk-averse stochastic multi-armed bandits". In: *Algorithmic Learning Theory: 24th International Conference, ALT 2013, Singapore, October 6-9, 2013. Proceedings 24.* Springer. 2013, pp. 218–233.
- [36] Alexandre Marthe, Aurélien Garivier, and Claire Vernade. "Beyond average return in markov decision processes". In: Advances in Neural Information Processing Systems 36 (2023), pp. 56488–56507.
- [37] Alexandre Marthe et al. "Efficient Risk-sensitive Planning via Entropic Risk Measures". In: *arXiv preprint arXiv:2502.20423* (2025).
- [38] Gergely Neu, Anders Jonsson, and Vicenç Gómez. "A unified view of entropy-regularized markov decision processes". In: *arXiv preprint arXiv:1705.07798* (2017).
- [39] Kishan Panaganti and Dileep Kalathil. "Sample complexity of robust reinforcement learning with a generative model". In: *International Conference on Artificial Intelligence and Statistics*. PMLR. 2022, pp. 9582–9602.
- [40] Martin L. Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- [41] Paria Rashidinejad et al. "Bridging offline reinforcement learning and imitation learning: A tale of pessimism". In: *Advances in Neural Information Processing Systems* 34 (2021), pp. 11702–11716.
- [42] Amir Sani, Alessandro Lazaric, and Rémi Munos. "Risk-aversion in multi-armed bandits". In: *Advances in neural information processing systems* 25 (2012).
- [43] Maria Grazia Scutella and Raffaella Recchia. "Robust portfolio asset allocation and risk measures". In: *Annals of Operations Research* 204.1 (2013), pp. 145–169.
- [44] Alexander Shapiro, Darinka Dentcheva, and Andrzej Ruszczynski. *Lectures on stochastic programming: Modeling and theory.* SIAM, 2021.
- [45] Laixi Shi et al. "The curious price of distributional robustness in reinforcement learning with a generative model". In: *Advances in Neural Information Processing Systems* 36 (2023), pp. 79903–79917.
- [46] Aaron Sidford et al. "Variance reduced value iteration and faster algorithms for solving Markov decision processes". In: *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*. 2018, pp. 770–787.
- [47] Satinder P Singh and Richard C Yee. "An upper bound on the loss from approximate optimal-value functions". In: *Machine Learning* 16.3 (1994), pp. 227–233.
- [48] Alexander L Strehl and Michael L Littman. "An analysis of model-based interval estimation for Markov decision processes". In: *Journal of Computer and System Sciences* 74.8 (2008), pp. 1309–1331.
- [49] Xihong Su, Marek Petrik, and Julien Grand-Clément. "Risk-averse Total-reward MDPs with ERM and EVaR". In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 39. 19. 2025, pp. 20646–20654.
- [50] Richard S Sutton, Andrew G Barto, et al. *Reinforcement learning: An introduction*. Vol. 1. 1. MIT press Cambridge, 1998.
- [51] Aviv Tamar et al. "Policy gradient for coherent risk measures". In: *Advances in neural information processing systems* 28 (2015).
- [52] Mengdi Wang. "Randomized linear programming solves the Markov decision problem in nearly linear (sometimes sublinear) time". In: *Mathematics of Operations Research* 45.2 (2020), pp. 517–546.
- [53] Tsachy Weissman et al. "Inequalities for the L1 deviation of the empirical distribution". In: *Hewlett-Packard Labs, Tech. Rep* (2003), p. 125.
- [54] Andrea Zanette, Mykel J Kochenderfer, and Emma Brunskill. "Almost horizon-free structure-aware best policy identification with a generative model". In: *Advances in Neural Information Processing Systems* 32 (2019).

- [55] Dongruo Zhou, Jiafan He, and Quanquan Gu. "Provably efficient reinforcement learning for discounted mdps with feature mapping". In: *International Conference on Machine Learning*. PMLR. 2021, pp. 12793–12802.
- [56] Matthew Zurek and Yudong Chen. "Span-Based Optimal Sample Complexity for Weakly Communicating and General Average Reward MDPs". In: *The Thirty-eighth Annual Conference on Neural Information Processing Systems*. 2024.
- [57] Matthew Zurek and Yudong Chen. "The Plugin Approach for Average-Reward and Discounted MDPs: Optimal Sample Complexity Analysis". In: 36th International Conference on Algorithmic Learning Theory. 2025.

## **A** Technical Lemmas

Recall that  $Q_k$  is the Q-function output by the algorithm after k iterations,  $\pi_k$  is the greedy policy with respect to  $Q_k$ , and that  $\pi^*$  is an optimal policy of the true MDP M.

The first lemma establishes a decomposition result for MB-RS-QVI, whose proof follows very similar lines to the proof of Lemma 3 in [1].

**Lemma 4.** For any state-action pair  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,

$$Q_k(s, a) \ge Q^*(s, a) - \|Q_k - \widehat{Q}^*\| - \|\widehat{Q}^{\pi^*} - Q^*\|.$$

Further, for any state  $s \in S$ ,

$$V^{\pi_k}(s) \ge Q^*(s) - \|V^{\pi_k} - \widehat{V}^{\pi_k}\| - \|\widehat{V}^{\pi_k} - \widehat{V}^*\| - \|\widehat{V}^{\pi^*} - V^*\|.$$

*Proof.* For any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , we have

$$Q_k(s,a) - Q^*(s,a) = Q_k(s,a) - \hat{Q}^*(s,a) + \hat{Q}^*(s,a) - Q^*(s,a)$$

$$\geq Q_k(s,a) - \hat{Q}^*(s,a) + \hat{Q}^{\pi^*}(s,a) - Q^*(s,a)$$

$$\geq -\|Q_k - \hat{Q}^*\| - \|\hat{Q}^{\pi^*} - Q^*\|.$$

Similarly, for any  $s \in \mathcal{S}$ , we have

$$V^{\pi_{k}}(s) - V^{*}(s) = V^{\pi_{k}}(s) - \widehat{V}^{\pi_{k}}(s) + \widehat{V}^{\pi_{k}}(s) - \widehat{V}^{*}(s) + \widehat{V}^{*}(s) - V^{*}(s)$$

$$\geq V^{\pi_{k}}(s) - \widehat{V}^{\pi_{k}}(s) + \widehat{V}^{\pi_{k}}(s) - \widehat{V}^{*}(s) + \widehat{V}^{\pi^{*}}(s) - V^{*}(s)$$

$$\geq -\|V^{\pi_{k}} - \widehat{V}^{\pi_{k}}\| - \|\widehat{V}^{\pi_{k}} - \widehat{V}^{*}\| - \|\widehat{V}^{\pi^{*}} - V^{*}\|,$$

and the lemma follows.

Next, we present two lemmas that collect a few useful inequalities. Some of these may be standard results, but for concreteness, we collect them here.

Lemma 5. It holds that

$$\log(1-x) \ge -x - x^2 + x^3 \qquad \forall x \in [0, \frac{1}{5}]$$

$$\log(1-x) \ge -x - 2x \qquad \forall x \in [0, \frac{1}{2}]$$

$$\log(1+x) \ge x - x^2 \qquad \forall x \in [0, \infty)$$

$$\log(1+x) \ge \frac{x}{2} \qquad \forall x \in [0, 1].$$

*Proof.* We only prove the first claim, as the rest could be proven using the technique after some elementary calculations.

Let f(x)=(1-x) and  $g(x)=-x-x^2+x^3$ . It holds that f(0)=g(0), and since we have  $f'(x)=\frac{1}{1-x}$  and  $g'(x)=-1-2x+3x^2$ , it follows easily that

$$f'(x) \ge g'(x) \Leftrightarrow 0 \le x(1 - 5x + 3x^2)$$

where the inequality is satisfied for all  $x \in [0, \frac{5-\sqrt{13}}{6}] \subseteq [0, \frac{1}{5}]$ . The result then follows from the fundamental theorem of calculus.

**Lemma 6.** Let  $\alpha > 1$ . For any  $x \in [0, \frac{1}{\alpha}]$ , it holds that

$$1 - (1 - x)^{\alpha} \ge \frac{x\alpha}{2} \,.$$

*Proof.* Define  $f(x) = 1 - (1-x)^{\alpha} - \frac{x\alpha}{2}$ . Since  $f''(x) = -\alpha(\alpha-1)(1-x)^{\alpha-2} < 0$ , f is strictly concave. Further, since f(0) = 0 and  $f(\frac{1}{\alpha}) = \frac{1}{2}(1-\frac{1}{\alpha})^{\alpha} > \frac{1}{2}-\frac{1}{e} > 0$ , f is positive on the interval  $[0,\frac{1}{\alpha}]$  and the result follows.

## **B** Risk Measures

In this section we give a very brief introduction to risk measures with proper definitions and key examples to the extend needed. The reason for this is that the literature is varied and that some inequalities and vocabulary change slightly, in particular over the choice of working with rewards or losses.

**Definition 2.** Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a background probability space and  $\mathcal{M}$  some convex cone of random variables defined on the background space. That is for any  $X, Y \in \mathcal{M}$  and  $\lambda > 0$  it holds that  $X + Y \in \mathcal{M}$  and  $\lambda X \in \mathcal{M}$ . A functional  $\psi : \mathcal{M} \to \mathbb{R}$  is said to be a risk measure if it satisfies the following properties:

$$\psi(0) = 0 (Normalization) (20)$$

if 
$$X \le Y$$
 then  $\psi(X) \ge \psi(Y)$  (Monotonicity) (21)

$$\psi(X+c) = \psi(X) - c \quad \forall c \in \mathbb{R}$$
 (Translation invariance) (22)

Any risk measure that in addition satisfies the properties

$$\psi(cX) = c\psi(X) \quad \forall a > 0$$
 (Positive homogeneity) (23)

$$\psi(X+Y) \le \psi(X) + \psi(Y) \tag{Subadditivity}$$

is called a coherent risk measure. A weaker notion is convex risk measure which is a risk measure that satisfies

$$\psi(\lambda X + (1 - \lambda)Y) \le \lambda \psi(X) + (1 - \lambda)\psi(Y) \quad \forall \lambda \in [0, 1]$$
 (Convexity) (25)

and finally a risk-measure  $\psi$  is called law-invariant if  $\psi(X)$  only depends on the distribution of X under  $\mathbb{P}$ .

We now mention some example of risk measures: The risk measure given by

$$ERM_{\beta}(X) = \frac{1}{\beta} \log \left( \mathbb{E}[e^{-\beta X}] \right)$$
 (26)

is known as the entropic risk measure with parameter  $\beta \neq 0$ . It is not positive homogeneous and hence not coherent. For  $\beta < 0$  it is not even convex. Letting  $\beta \to 0$  one recovers the expectation and letting  $\beta \to \infty$  one recovers the essential infimum risk measure.

The risk measure given by

$$VaR_{\alpha}(X) := q_{\alpha}(X) := \inf\{x \in \mathbb{R} : F_X(x) \ge \alpha\}$$
(27)

is called the value-at-risk at level  $\alpha \in (0,1)$  and is in general not sub-additive hence also not coherent.

The risk measure given by

$$CVaR_{\alpha}(X) := \frac{1}{1-\alpha} \int_{\alpha}^{1} VaR_{u}(X) du$$
 (28)

is known as the conditional Value-at-risk (CVaR) or sometimes as the expected shortfall (ES) and is known to be a coherent risk-measure.

All the examples so far are evidently law-invariant.

The actual functional that will be used to rank random variables is the negative of the ERM-risk measure given in the example above with the interpretation being that a lower quantity of risk is preferable. It follows directly then that the functional  $\rho: \mathcal{M} \to \mathbb{R}$  given by  $\rho(X) := -\text{ERM}_{\beta}(X)$  has the following properties:

$$\rho(0) = 0 \tag{Normalization}$$

$$\text{if } X \leq Y \text{ then } \rho(X) \leq \rho(Y) \tag{Monotonicity} \tag{30}$$

$$\rho(X+c) = \rho(X) + c$$
 (Translation invariance) (31)

It is common in the literature to overload notation and also refer to  $\rho$  as the ERM and we will do so and henceforth we will no longer care about risk measures, but only about this specific functional  $\rho$ . It follows immediately from the normalization and translation invariance that for any real number  $c \in \mathbb{R}$  it holds that  $\rho(c) = c$ .

We will often use the short-hand notation  $\rho_{s,a}(V(s'))$  as  $\rho$  applied to the random variable X that takes on the values  $\{V(s')\}_{s'\in S}$  with probabilities  $\mathbb{P}(X=V(s'))=P(s'|s,a)$ .

## C Bellman Optimality and Bellman Recursions

In this section we properly define the state-value functions and state-action value functions of any possibly history-dependent policy  $\pi$  and show that the problem of finding an optimal policy can be achieved by a stationary policy and that the value functions satisfy Bellman recursions when the value functions are defined iteratively with respect to the ERM. Several similar results exist in the literature e.g. [4] and [6] which also covers the  $\beta>0$  case. These results are derived under more general assumptions on  $\mathcal S$  and  $\mathcal A$ . These general assumptions are trivially satisfied when  $\mathcal S$  and  $\mathcal A$  are finite but their proofs require assumptions on the functionals to ensure the existence of a stationary optimal policy usually by invoking a measurable selection theorem. We avoid this complication by only considering finite  $\mathcal S$  and  $\mathcal A$  and we in turn also give the first proof for state-action value functions and not just for value-functions which is needed as we consider the problem of learning.

Let  $M=(\mathcal{S},\mathcal{A},P,R,\gamma,\rho)$  be a finite MDP with  $\rho$  being the ERM, and  $R(s,a)\in[0,1]$  a deterministic reward function. Let  $D=\mathcal{S}\times\mathcal{A},\,H_1=\mathcal{S}$  and  $H_k=D^{k-1}\times\mathcal{S}$  for  $k\geq 2$  the set of all possible histories up to stage k. A policy  $\pi=(\pi_k)_{k\in\mathbb{N}}$  a sequence of maps  $\pi_k:H_k\to\mathcal{A}$ . We denote the set of all policies  $\Pi$  and identify the set of all stationary policies with the set of measurable maps F from  $\mathcal{S}$  to  $\mathcal{A}$  which is simply the set of all maps from  $\mathcal{S}$  to  $\mathcal{A}$  since all maps between finite sets that are both equipped with the discrete topology are measurable with respect to the induced Borel  $\sigma$ -algebras. Let  $B(H_k)$  be the set of maps  $V_k:H_k\to\mathbb{R}$  equipped with the supremum normal let  $\pi=(\pi_k)_{k\in\mathbb{N}}$  be any policy. For any  $V_{k+1}\in B(H_{k+1})$  and  $h_k\in H_k$  we denote by  $\rho_{h_k,\pi_k}(V_{k+1})$  the functional  $\rho$  applied to the random variable concentrated on the set  $\{V_{k+1}(h_k,\pi_k(h_k),s')\}_{s'\in S}$  with  $\mathbb{P}(s_{k+1}=s')=\mathbb{P}(s'|s_k,\pi_k(h_k))$ . By monotonicity of  $\rho$  we get that  $\rho_{h_k,\pi_k}(V_{k+1})\leq \|V_{k+1}\|$ 

Next, we define the operators  $L_{\pi_k}: B(H_{k+1}) \to B(H_k)$  by

$$(L_{\pi_k}V_{k+1})(h_k) = L_{\pi_k, V_{k+1}}(h_k) := R(s_k, \pi_k(h_k)) + \gamma \rho_{h_k, \pi_k}(V_{k+1})$$
(32)

and similarly we define  $L_a: B(H_{k+1}) \to B(H_k)$  by

$$(L_a V_{k+1})(h_k) = L_{a,V_{k+1}}(h_k) := R(s_k, a) + \gamma \rho_{s_k, a}(V_{k+1})$$

with  $\rho_{s_k,a}$  defined analogously as  $\rho_{h_k,\pi_k}$  as above. By the basic properties of risk-measures it follows directly that  $0 \le L_{\pi_k V_{k+1}}(h_k) \le 1 + \gamma \|V_{k+1}\|$  and similarly for  $L_a$ .

For any initial state  $s_1 = s$  we define the N-step discounted utility as

$$J_N(s, a, \pi) := (L_a \circ L_{\pi_2} \circ \dots \circ L_{\pi_N}) \mathbf{0}(s) \tag{33}$$

where  $\mathbf{0}(h_k) = 0$  for all  $h_k \in H_k$  and all  $k \in \mathbb{N}$ .

By monotonicity of  $\rho$ , it holds that the sequence  $(J_N(s,a,\pi))_{N\in\mathbb{N}}$  is non decreasing and bounded in the interval  $[0,\frac{1}{1-\gamma}]$  for any  $s,a,\pi\in\mathcal{S}\times\mathcal{A}\times\Pi$  and so the limit

$$J(s, a, \pi) := \lim_{N \to \infty} J_N(s, a, \pi)$$

exists for any state s, any action a and any policy  $\pi$ .

The problem of the agent is to find  $J^*(s,a) = \sup_{\pi \in \Pi} J(s,a,\pi)$  and an optimal policy  $\pi^*$  that solves  $J(s,a,\pi^*) = J^*(s,a)$ .

**Theorem 6.** There exist a unique non-negative function  $Q \in B(S \times A)$  (non-negative maps from  $S \times A \to \infty$  equipped with sup-norm) and a stationary decision rule  $f^* : S \to A$  such that

$$Q(s,a) = R(s,a) + \gamma \rho_{s,a} (\max_{a'} Q(s',a')),$$
(34)

$$= R(s, a) + \gamma \rho_{s, a}(Q(s', f^*(s'))). \tag{35}$$

Moreover,  $Q(s,a) = J^*(s,a) = J(s,a,f^*)$  meaning that  $f^*$  is an optimal stationary policy.

*Proof.* We start by proving existence of Q. Let  $L: B(\mathcal{S} \times \mathcal{A}) \to B(\mathcal{S} \times \mathcal{A})$  denote the operator given by

$$LQ(s,a) := R(s,a) + \gamma \rho_{s,a} (\max_{a'} Q(s',a')).$$
(36)

Let  $Q_1, Q_2 \in B(\mathcal{S} \times \mathcal{A})$ . We then have for all (s, a) that

$$LQ_1(s,a) - LQ_2(s,a) = \gamma \left[ \rho_{s,a} (\max_{a'} Q_1(s',a')) - \rho_{s,a} (\max_{a'} Q_2(s',a')) \right]$$
(37)

$$= \gamma \left[ \rho_{s,a} \left( \max_{a'} Q_1(s',a') \right) - \rho_{s,a} \left( \max_{a'} Q_1(s',a') - \max_{a'} Q_1(s',a') + \max_{a'} Q_2(s',a') \right) \right]$$
(38)

$$\leq \gamma \left[\rho_{s,a}(\max_{a'} Q_1(s',a')) - \rho_{s,a}(\max_{a'} Q_1(s',a') + \max_{a'} \{Q_2(s',a') - Q_1(s',a'))\}\right)]$$
(39)

$$\leq \gamma \left[\rho_{s,a}(\max_{a'} Q_1(s',a')) - \rho_{s,a}(\max_{a'} Q_1(s',a') + \|Q_1 - Q_2\|)\right] \tag{40}$$

$$= \gamma \|Q_1 - Q_2\|. \tag{41}$$

We start by showing that  $L: B(\mathcal{S} \times \mathcal{A}) \to B(\mathcal{S} \times \mathcal{A})$ , that is it takes non-negative functions and returns non-negative functions. By normalization and monotonicity of  $\rho$  we have for any non-negative  $Q \in B(\mathcal{S} \times \mathcal{A})$  that

$$LQ(s,a) = R(s,a) + \gamma \rho_{s,a}(\max_{a'} Q(s,a)) \ge 0 + \gamma \rho_{s,a}(0) = 0$$
(42)

and since we by a completely similar argument have  $LQ_2(s,a) - LQ_1(s,a) \le \gamma \|Q_1 - Q_2\|$ , we have that L is a contraction, and since  $\mathcal{S} \times \mathcal{A}$  we can identify  $B(\mathcal{S} \times \mathcal{A})$  with the closed subset of the complete metric space  $(\mathbb{R}^{S \times A}, \|\cdot\|)$  that consists of vectors with non-negative coordinates. Since this subspace is closed, it is also a complete metric space and the existence of Q then follows from the Banach fixed point theorem.

Since there are only finitely many states and actions we can pick a stationary decision rule where  $f^*(s)$  is an arbitrary element of  $\operatorname{argmax}_a Q(s, a)$ .

Let V be the function given  $V(s) := Q(s, f^*(s))$  for all s. We then see that

$$V(s) \ge R(s, a) + \gamma \rho_{s, a}(V(s')) \tag{43}$$

for every  $s \in \mathcal{S}$ . Let  $\pi = (\pi_k)_k \in \mathbb{N}$  be any policy in  $\Pi$ . The above inequality then shows that for any history  $h_k, k \in \mathbb{N}$  we have that  $V(s_k) \geq L_{\pi_k} V(h_k)$  and furthermore we note that  $Q(s_1, a) = L_a V(h_1)$ . This implies for any  $N \in \mathbb{N}$  that

$$Q(s,a) \ge (L_a \circ L_{\pi_2} \circ \cdots \circ L_{\pi_N})V(s) \ge (L_a \circ L_{\pi_2} \circ \cdots \circ L_{\pi_N})\mathbf{0}(s) = J_N(s,a,\pi), \tag{44}$$

where we have used that  $Q(s, a) \ge 0$ . Finally taking the limit we find that  $Q(s, a) \ge J(s, a, \pi)$ .

Finally we aim to that  $Q(s,a) \leq J(s,a,f^*)$ . By induction, we wish to show that  $V(s) \leq J_N(s,f^*(s),f^*) + \gamma^N \|V\|$  for all  $N \in \mathbb{N}$ . For the induction step, we start by noting that  $J_1(s,f^*(s),f^*) = R(s,f^*(s))$  and so

$$V(s) = R(s, f^*(s)) + \gamma \rho_{s, f^*(s)}(V(s'))$$
(45)

$$\leq R(s, f^*(s)) + \gamma \rho_{s, f^*(s)}(\|V\|)$$
 (46)

$$= R(s, f^*(s)) + \gamma ||V|| \tag{47}$$

$$= J_1(s, f^*(s), f^*) + \gamma ||V||, \tag{48}$$

for all  $(s,a) \in \mathcal{S} \times \mathcal{A}$ . For the induction step, we assume that  $V(s) \leq J_N(s,f^*(s),f^*) + \gamma^N ||V||$ . By using that  $V(s) = L_{f^*}V(s)$  and that L is monotone, we see that

$$V(s) = L_{f^*}V(s) \tag{49}$$

$$\leq L_{f^*}(J_N(s, f^*(s), f^*) + \gamma^N ||V||)$$
 (50)

$$= \left( R(\cdot, f^*(\cdot)) + \gamma \rho_{\cdot, f^*(\cdot)} (J_N(\cdot, f^*(\cdot), f^*) + \gamma^N ||V|| \right) (s)$$

$$(51)$$

$$= J_{N+1}(s, f^*(s), f^*) + \gamma^{N+1} ||V||, \tag{52}$$

from which taking the limit  $N \to \infty$ , we get that  $V(s) \le J(s, f^*(s), f^*)$ .

Finally, since

$$Q(s,a) = L_a V(s) < L_a J(s, f^*(s), f^*) = J(s, a, f^*),$$
(53)

the conclusion holds.

Since this shows that an optimal stationary policy exists, it will suffice to consider only stationary policies and one can by completely analogous arguments show that for any stationary policy  $\pi$ , there exists a non-negative map  $Q^{\pi} \in B(\mathcal{S} \times \mathcal{A})$  such that  $Q^{\pi}(s,a) = J(s,a,\pi)$  such that  $Q^{\pi}$  satisfies the Bellman recursion:

$$Q^{\pi}(s,a) = R(s,a) + \gamma \rho_{s,a}(Q^{\pi}(s',\pi(s'))), \tag{54}$$

and similarly for state-value functions  $V^{\pi}(s) := Q^{\pi}(s, \pi(s))$ .

We also remark that in the proof we see directly that  $Q(s,a) \in [0, \frac{1}{1-\gamma}]$  for all (s,a).

## **D** Deterioration of Greedy Policy

Next we show a result that bounds the quality of a greedy policy with respect to the quality of the value-function for which the policy is greedy. The result is a generalization of [47] from the expectation to that of the ERM. Throughout, we use the notation  $\rho_{s,a}(V(s'))$  as shorthand notation for  $\rho$  applied to the categorical random variable X with support  $\{V(s')\}_{s'\in\mathcal{S}}$  where  $\mathbb{P}(X=V(s'))=P(s'|s,a)$ .

**Theorem 7.** Let  $\widehat{V}$  be a function for which  $\|V^* - \widehat{V}\| < \varepsilon$  and let  $\pi^G := \operatorname{argmax}_a[R(s,a) + \gamma \rho_{s,a}(\widehat{V}(s'))]$  be a greedy policy with respect to  $\widehat{V}$ . Let the value function of this greedy policy be denoted  $V^G := V^{\pi^G}$ . It then holds that

$$||V^* - V^G|| \le \frac{2\gamma}{1 - \gamma} \varepsilon. \tag{55}$$

*Proof.* Let  $\bar{s}$  be a state such that  $\|V^* - V^G\| = V^*(\bar{s}) - V^G(\bar{s})$ . We then consider the two actions  $a^* := \pi^*(\bar{s})$  (pick any if more) and  $a^G := \pi^G(\bar{s})$ . Since  $\pi^G$  is greedy w.r.t.  $V^G$ , we have that

$$R(\bar{s}, a^*) + \gamma \rho_{\bar{s}, a^*}(\widehat{V}(s')) \le R(\bar{s}, a^G) + \gamma \rho_{\widehat{s}, a^G}(\widehat{V}(s')).$$

By assumption, it holds for any  $s \in \mathcal{S}$  that

$$V^*(s) - \varepsilon \le \widehat{V}(s) \le V^*(s) + \varepsilon.$$

By monotonicity and translation invariance of  $\rho$ , we thus get

$$R(\bar{s}, a^*) + \gamma \rho_{\bar{s}, a^*}(\hat{V}(s')) \ge R(\bar{s}, a^*) + \gamma \rho_{\bar{s}, a^*}(V^*(s') - \varepsilon)$$
 (56)

$$= R(\bar{s}, a^*) + \gamma \rho_{\bar{s}, a^*}(V^*(s')) - \gamma \varepsilon, \tag{57}$$

and similarly we have

$$R(\bar{s}, a^G) + \gamma \rho_{\bar{s}, a^G}(\hat{V}(s')) \le R(\bar{s}, a^G) + \gamma \rho_{\bar{s}, a^G}(V^*(s')) + \gamma \varepsilon, \tag{58}$$

which collectively implie

$$R(\bar{s}, a^*) - R(\bar{s}, a^G) \le 2\gamma \varepsilon + \gamma \left(\rho_{\bar{s}, a^G}(V^*(s')) - \rho_{\bar{s}, a^*}(V^*(s'))\right).$$
 (59)

Finally, we obtain

$$\begin{split} V^*(\bar{s}) - V^G(\bar{s}) &= R(\bar{s}, a^*) - R(\bar{s}, a^G) + \gamma \rho_{\bar{s}, a^*}(V^*(s')) - \gamma \rho_{\bar{s}, a^G}(V^G(s')) \\ &\leq 2\gamma \varepsilon + \gamma \rho_{\bar{s}, a^G}(V^*(s')) - \gamma \rho_{\bar{s}, a^*}(V^*(s') + \gamma \rho_{\bar{s}, a^*}(V^*(s')) - \gamma \rho_{\bar{s}, a^G}(V^G(s')) \\ &= 2\gamma \varepsilon + \gamma \left(\rho_{\bar{s}, a^G}(V^*(s') - \rho_{\bar{s}, a^G}(V^G(s'))\right) \\ &= 2\gamma \varepsilon + \gamma \|V^* - V^G\|, \end{split}$$

from which the result follows.

## E Concentration Results

In this section we collect some lemmas on concentration, they directly complement lemma 3.

Let N denote the number of calls to the generative model on each state-action pair such that the total number of calls is SAN. Let  $\widehat{P}(s'|s,a)$  denote the plug-in estimator obtained from N samples of  $s' \sim P(\cdot|s,a)$ , that is  $\widehat{P}(s'|s,a) = \frac{\sum_{n=1}^{N} \mathbb{1}_{\{s_n=s'\}}}{N}$  where  $s_n$  is the outcome of the random variable  $X_n$  taking values in S according to  $P(\cdot|s,a)$ .

**Lemma 7.** Fix  $\varepsilon > 0$ . If every state-action pair has been tried at least  $N = 8 \frac{S + \log(\frac{SA}{\delta})}{\tau^2}$  times, then it holds that  $\max_{(s,a)} \|P(\cdot|s,a) - \widehat{P}(\cdot|s,a)\|_1 \le \varepsilon$  with probability at least  $1 - \delta$ .

*Proof.* Using the Weissman inequality [53], any confidence interval for a state-action pair that have been tried m times have size  $2\sqrt{2[\log(2^S-2)-\log(\delta_P)]/N}$ . Setting this size to be smaller than  $\tau$  and solving for N, we find that

$$N \ge \frac{8}{\varepsilon^2} \left( \log(2^S - 2) - \log(\delta_P) \right). \tag{60}$$

Noting that  $8(S + \log(1/\delta_P))/\varepsilon^2 \ge 8[\log(2^S - 2) - \log(\delta_P)]/\varepsilon^2$  and substituting  $\delta_P = \frac{\delta}{SA}$ , the result follows by the union bound since each of the SA confidence balls contain  $P(\cdot|s,a)$  with probability at least  $1 - \delta$ .

**Theorem 8.** Let  $\pi$  be any fixed policy. For  $\beta > 0$  we have that if  $N > \frac{(1-e^{-\beta\frac{1}{1-\gamma}})^2}{2\varepsilon^2}\log(2SA/\delta)$  then it holds that with probability at least  $1-\delta$  that

$$\max_{s,a} \left| \sum_{s'} [P(s'|s,a) - \widehat{P}(s'|s,a)] e^{-\beta V^{\pi}(s')} \right| < \varepsilon \tag{61}$$

and if  $\beta < 0$  it holds that if  $N > \frac{(1 - e^{-\beta \frac{1}{1 - \gamma}})^2}{2\varepsilon^2} \log(2SA/\delta)$ , then it holds that with probability at least  $1 - \delta$  that

$$\max_{s,a} \Big| \sum_{s'} [P(s'|s,a) - \widehat{P}(s'|s,a)] e^{|\beta|(V^{\pi}(s') - \frac{1}{1-\gamma})} \Big| < \varepsilon.$$
 (62)

*Proof.* We only prove the first claim:  $\beta > 0$  as the other case is completely similar. We note that for the random variable  $\sum_{s'} 1_{\{X_n = s'\}} (s'|s,a) e^{-\beta V^{\pi}(s')}$ , we have that

$$\mathbb{E}\left[\sum_{s'} 1_{\{X_n = s'\}}(s'|s, a)e^{-\beta V^{\pi}(s')}\right] = \sum_{s'} \mathbb{E}\left[1_{\{X_n = s'\}}\right]e^{-\beta V^{\pi}(s')}$$
(63)

$$= \sum_{s'} P(s'|s,a)e^{-\beta V^{\pi}(s')}$$
 (64)

and that it is bounded in  $[e^{-\beta \frac{1}{1-\gamma}}, 1]$ . Also, since

$$\sum_{s'} \widehat{P}(s'|s,a)e^{-\beta V^{\pi}(s')} = \frac{1}{N} \sum_{n=1}^{N} \sum_{s'} 1_{\{X_n = s'\}} (s'|s,a)e^{-\beta V^{\pi}(s')}, \tag{65}$$

it follows directly from Hoeffding's inequality that

$$\mathbb{P}\left(\left|\sum_{s'} [P(s'|s,a) - \widehat{P}(s'|s,a)]e^{-\beta V^{\pi}(s')}\right| \ge \varepsilon\right) \le 2\exp\left(-\frac{2N\varepsilon^2}{\left(1 - e^{-\beta\frac{1}{1-\gamma}}\right)^2}\right). \tag{66}$$

Thus, by picking  $N=\frac{(1-e^{-\beta\frac{1}{1-\gamma}})^2}{2\varepsilon^2}\log(2SA/\delta)$  and a union bound,

$$\mathbb{P}\bigg(\max_{s,a} |\sum_{s'} [P(s'|s,a) - \widehat{P}(s'|s,a)]e^{-\beta V^{\pi}(s')}| \ge \varepsilon\bigg) \le \delta.$$
 (67)

## F Analysis of MB-RS-OVI: Missing Proofs

#### F.1 Proof of Lemma 1

*Proof.* We only give the proof for  $\mathcal{T}$  as the claim for  $\mathcal{T}^{\pi}$  could be proven using extremely similar lines.

Consider two maps  $Q: \mathbb{R}^{S \times A} \to \mathbb{R}^{S \times A}$  and  $W: \mathbb{R}^{S \times A} \to \mathbb{R}^{S \times A}$ , and let  $Q' = \mathcal{T}Q$  and  $W' = \mathcal{T}W$  be their respective  $\mathcal{T}$ -transforms. Let (s,a) be any pair such that  $|Q'(s,a) - W'(s,a)| = \|Q' - W'\|_{\infty}$ , and assume without loss of generality that  $Q'(s,a) \geq W'(s,a)$ . Further, define

$$V(s) := \max_a Q(s,a)\,, \qquad X(s) := \max_a W(s,a)\,.$$

Assuming that  $\beta > 0$  (the case  $\beta < 0$  is completely similar), we then have

$$\begin{split} \|Q' - W'\| &= Q'(s,a) - W'(s,a) \\ &= -\frac{\gamma}{\beta} \log \left( \sum_{s'} P(s'|s,a) e^{-\beta V(s')} \right) + \frac{\gamma}{\beta} \log \left( \sum_{s'} P(s'|s,a) e^{-\beta X(s')} \right) \\ &= -\frac{\gamma}{\beta} \log \left( \sum_{s'} P(s'|s,a) e^{-\beta X(s') - \beta (V(s') - X(s'))} \right) + \frac{\gamma}{\beta} \log \left( \sum_{s'} P(s'|s,a) e^{-\beta X(s')} \right) \\ &\leq -\frac{\gamma}{\beta} \log \left( \sum_{s'} P(s'|s,a) e^{-\beta X(s') - \beta \|V - X\|} \right) + \frac{\gamma}{\beta} \log \left( \sum_{s'} P(s'|s,a) e^{-\beta X(s')} \right) \\ &= \gamma \|V - X\| \\ &\leq \gamma \|Q - W\| \,, \end{split}$$

and the lemma follows.

#### F.2 Proof of Lemma 2

*Proof.* By Lemma 1, we have that  $\mathcal{T}$  is a  $\gamma$ -contraction and that  $Q^*$  is its unique fixed point. We thus have  $\|Q_k - Q^*\| = \|\mathcal{T}Q_{k-1} - \mathcal{T}Q^*\| \le \gamma \|Q_{k-1} - Q^*\|$ . Applying this inequality k times yields

$$||Q_k - Q^*|| \le \gamma^k ||Q_0 - Q^*|| \le \frac{\gamma^k}{1 - \gamma}.$$

Solving  $\frac{\gamma^k}{1-\gamma}$  for k, we get that if  $k>\frac{\log(\frac{1}{(1-\gamma)\varepsilon})}{\log(1/\gamma)}$ , then  $\|Q_k-Q^*\|<\varepsilon$ , thus proving the first claim.

To show the other claim, we start by noting that  $\|V^{\pi_k} - V^*\| \le \|Q^{\pi_k} - Q^*\|$ . Note also that by design we have that  $\mathcal{T}^{\pi_k}Q^{\pi_k} = Q^{\pi_k}$  and that  $\mathcal{T}Q_k = \mathcal{T}^{\pi_k}Q_k$ . Thus,

$$||Q^{\pi_k} - Q^*|| \le ||Q^{\pi_k} - Q_k|| + ||Q_k - Q^*||.$$
 (68)

The first term in the right-hand side is bounded as follows:

$$\begin{aligned} \|Q^{\pi_k} - Q_k\| &= \|\mathcal{T}^{\pi_k} Q^{\pi_k} - Q_k\| \\ &\leq \|\mathcal{T}^{\pi_k} Q^{\pi_k} - \mathcal{T} Q_k\| + \|\mathcal{T} Q_k - Q_k\| \\ &= \|\mathcal{T}^{\pi_k} Q^{\pi_k} - \mathcal{T}^{\pi_k} Q_k\| + \|\mathcal{T} Q_k - \mathcal{T} Q_{k-1}\| \\ &\leq \gamma \|Q^{\pi_k} - Q_k\| + \gamma \|Q_k - Q_{k-1}\| \,, \end{aligned}$$

which means that

$$\|Q^{\pi_k} - Q_k\| \le \frac{\gamma}{1 - \gamma} \|Q_k - Q_{k-1}\| \le \frac{\gamma^k}{1 - \gamma} \|Q_1 - Q_0\| \le \frac{\gamma^k}{(1 - \gamma)^2}.$$
 (69)

The proof is completed by observing that picking  $k > \frac{\log(\frac{2}{(1-\gamma)^2\varepsilon})}{\log(1/\gamma)}$  implies  $||V^{\pi_k} - V^*|| < \varepsilon$ .

#### F.3 Proof of Lemma 3

*Proof.* There are four different cases to consider, namely the combinations arrising from  $\beta>0$  vs  $\beta<0$  and wether on the state-action pair (s,a) that realizes the maximum it holds that  $Q_1(s,a)>Q_2(s,a)$  or  $Q_2(s,a)>Q_1(s,a)$ .

Case 1:  $\beta > 0, Q_1(s, a) > Q_2(s, a)$ .

$$\begin{split} \|Q_1 - Q_2\| &= \frac{\gamma}{\beta} \ln \left( \frac{\sum_{s'} P_2(s'|s,a) e^{-\beta V_2(s')}}{\sum_{s'} P_1(s'|s,a) e^{-\beta V_1(s')}} \right) \\ &= \frac{\gamma}{\beta} \ln \left( \frac{\sum_{s'} P_2(s'|s,a) e^{-\beta V_1(s')} + \beta (V_2(s') - V_1(s'))}{\sum_{s'} P_1(s'|s,a) e^{-\beta V_1(s')}} \right) \\ &\leq \frac{\gamma}{\beta} \ln \left( e^{\beta \|V_1 - V_2\|} \frac{\sum_{s'} P_2(s'|s,a) e^{-\beta V_1(s')}}{\sum_{s'} P_1(s'|s,a) e^{-\beta V_1(s')}} \right) \\ &= \gamma \|V_1 - V_2\| + \frac{\gamma}{\beta} \ln \left( \frac{\sum_{s'} P_2(s'|s,a) e^{-\beta V_1(s')}}{\sum_{s'} P_1(s'|s,a) e^{-\beta V_1(s')}} \right) \\ &\leq \gamma \|Q_1 - Q_2\| + \frac{\gamma}{\beta} \ln \left( 1 + \frac{\sum_{s'} P_2(s'|s,a) e^{-\beta V_1(s')}}{\sum_{s'} P_1(s'|s,a) e^{-\beta V_1(s')}} \right) \\ &\leq \gamma \|Q_1 - Q_2\| + \frac{\gamma}{\beta} \frac{\sum_{s'} P_2(s'|s,a) e^{-\beta V_1(s')} - \sum_{s'} P_1(s'|s,a) e^{-\beta V_1(s')}}{\sum_{s'} P_1(s'|s,a) e^{-\beta V_1(s')}} \\ &\leq \gamma \|Q_1 - Q_2\| + \frac{\gamma}{\beta} \frac{|\sum_{s'} [P_2(s'|s,a) - P_1(s'|s,a)] e^{-\beta V_1(s')}}{e^{-\beta \frac{1}{1-\gamma}}} \,. \end{split}$$

Rearranging the terms yields the asserted result

$$||Q_1 - Q_2|| \le \frac{\gamma}{1 - \gamma} \frac{e^{\beta \frac{1}{1 - \gamma}}}{\beta} \Big| \sum_{s'} [P_2(s'|s, a) - P_1(s'|s, a)] e^{-\beta V_1(s')} \Big|.$$

Case 2:  $\beta > 0$  and  $Q_1(s,a) < Q_2(s,a)$ . The proof is very similar but the extension  $V_2(s) = V_1(s) + V_2(s) - V_1(s)$  is now done in the numerator instead:

$$\begin{split} \|Q_1 - Q_2\| &= \frac{\gamma}{\beta} \ln \left( \frac{\sum_{s'} P_1(s'|s,a) e^{-\beta V_1(s')}}{\sum_{s'} P_2(s'|s,a) e^{-\beta V_2(s')}} \right) \\ &= \frac{\gamma}{\beta} \ln \left( \frac{\sum_{s'} P_1(s'|s,a) e^{-\beta V_1(s')}}{\sum_{s'} P_1(s'|s,a) e^{-\beta V_1(s')}} \right) \\ &\leq \frac{\gamma}{\beta} \ln \left( e^{\beta \|V_1 - V_2\|} \frac{\sum_{s'} P_1(s'|s,a) e^{-\beta V_1(s')}}{\sum_{s'} P_2(s'|s,a) e^{-\beta V_1(s')}} \right) \\ &= \gamma \|V_1 - V_2\| + \frac{\gamma}{\beta} \ln \left( \frac{\sum_{s'} P_1(s'|s,a) e^{-\beta V_1(s')}}{\sum_{s'} P_2(s'|s,a) e^{-\beta V_1(s')}} \right) \\ &\leq \gamma \|Q_1 - Q_2\| + \frac{\gamma}{\beta} \ln \left( 1 + \frac{\sum_{s'} P_2(s'|s,a) e^{-\beta V_1(s')}}{\sum_{s'} P_1(s'|s,a) e^{-\beta V_1(s')}} \right) \\ &\leq \gamma \|Q_1 - Q_2\| + \frac{\gamma}{\beta} \frac{\sum_{s'} P_1(s'|s,a) e^{-\beta V_1(s')} - \sum_{s'} P_2(s'|s,a) e^{-\beta V_1(s')}}{\sum_{s'} P_2(s'|s,a) e^{-\beta V_1(s')}} \\ &\leq \gamma \|Q_1 - Q_2\| + \frac{\gamma}{\beta} \frac{\sum_{s'} P_1(s'|s,a) - P_2(s'|s,a) e^{-\beta V_1(s')}}{\sum_{s'} P_2(s'|s,a) e^{-\beta V_1(s')}} \\ &\leq \gamma \|Q_1 - Q_2\| + \frac{\gamma}{\beta} \frac{\sum_{s'} [P_1(s'|s,a) - P_2(s'|s,a)] e^{-\beta V_1(s')}}{\sum_{s'} P_2(s'|s,a)] e^{-\beta V_1(s')}}, \end{split}$$

which again yields

$$||Q_1 - Q_2|| \le \frac{\gamma}{1 - \gamma} \frac{e^{\beta \frac{1}{1 - \gamma}}}{\beta} |\sum_{s'} [P_2(s'|s, a) - P_1(s'|s, a)] e^{-\beta V_1(s')}|.$$

Case 3:  $\beta < 0$  and  $Q_1(s, a) > Q_2(s, a)$ .

$$\begin{split} \|Q_1 - Q_2\| &= \frac{\gamma}{|\beta|} \ln \left( \frac{\sum_{s'} P_1(s'|s,a) e^{|\beta|V_1(s')}}{\sum_{s'} P_2(s'|s,a) e^{|\beta|V_2(s')}} \right) \\ &= \frac{\gamma}{|\beta|} \ln \left( \frac{\sum_{s'} P_1(s'|s,a) e^{|\beta|V_1(s')}}{\sum_{s'} P_1(s'|s,a) e^{|\beta|V_1(s')} - |\beta|(V_2(s') - V_1(s'))} \right) \\ &\leq \frac{\gamma}{|\beta|} \ln \left( \frac{\sum_{s'} P_1(s'|s,a) e^{|\beta|V_1(s')}}{\sum_{s'} P_1(s'|s,a) e^{|\beta|V_1(s')}} \right) \\ &= \gamma \|V_1 - V_2\| + \frac{\gamma}{|\beta|} \ln \left( \frac{\sum_{s'} P_1(s'|s,a) e^{|\beta|V_1(s')}}{\sum_{s'} P_2(s'|s,a) e^{|\beta|V_1(s')}} \right) \\ &\leq \gamma \|Q_1 - Q_2\| + \frac{\gamma}{\beta} \ln \left( 1 + \frac{\sum_{s'} P_1(s'|s,a) e^{|\beta|V_1(s')}}{\sum_{s'} P_2(s'|s,a) e^{|\beta|V_1(s')}} \right) \\ &\leq \gamma \|Q_1 - Q_2\| + \frac{\gamma}{|\beta|} \left| \sum_{s'} [P_1(s'|s,a) - P_2(s'|s,a)] e^{|\beta|V_1(s')} \right| \\ &\gamma \|Q_1 - Q_2\| + \frac{\gamma}{|\beta|} e^{|\beta| \frac{1}{1-\gamma}} \left| \sum_{s'} [P_2(s'|s,a) - P_1(s'|s,a)] e^{|\beta| \left[V_1(s') - \frac{1}{1-\gamma}\right]} \right|, \end{split}$$

which implies

$$||Q_1 - Q_2|| \le \frac{\gamma}{1 - \gamma} \frac{e^{|\beta| \frac{1}{1 - \gamma}}}{\beta} \Big| \sum_{s'} [P_2(s'|s, a) - P_1(s'|s, a)] e^{|\beta| [V_1(s') - \frac{1}{1 - \gamma}]} \Big|.$$

Case 4:  $\beta < 0$  and  $Q_2(s,a) > Q_1(s,a)$ . The proof of this case is similar to the other three cases and is omitted and the final part of the lemma follows by the triangle inequality and the fact that  $e^{-|\beta|V_1(s)} < 1$  and  $e^{-|\beta|(\frac{1}{1-\gamma}-V_1(s))} < 1$ .

## G Lower Bound on Bernoulli Likelihood Ratio

In this section, we revisit and develop a technical result that bounds the likelihood ratio of two samples under different hypotheses on a high probability event. Parts of the proof closely resembles parts of Lemma 17 in [21]; however, we stress that our treatment fixes an error in the proof, which however requires slightly stronger assumptions than those imposed in [21]. In addition, while the result in [21] only considers  $p \ge \frac{1}{2}$ , ours deal with both cases of  $p \ge \frac{1}{2}$  and  $p < \frac{1}{2}$ .

Let  $p \in (0,1)$  and  $\tilde{p} = \max\{p, 1-p\}$ . Let  $\alpha \in (0, \frac{1-\tilde{p}}{5}]$ . Consider two coins (Bernoulli random variables), one with bias q=p and one with bias  $q=p+\alpha$ . We name the two statistical hypotheses  $H_0: q=p$  and  $H_1: q=p+\alpha$ .

Let W be the outcome of flipping one of the coins t times and the associated likelihood function under hypothesis m as

$$L_m(w) := \mathbb{P}_m(W = w) \tag{70}$$

for hypothesis  $H_m$  with  $m \in \{0,1\}$  and for every possible history of outcomes w, and where  $\mathbb{P}_m(W=w)$  denotes the probability of observing the history w under the hypothesis  $H_m$ . The likelihood function defines a random variable  $L_m(W)$ , where W is the stochastic process of realized coin tosses.

Let  $t \in \mathbb{N}$  and  $\theta = \exp\left(-\frac{c_1\alpha^2t}{p(1-p)}\right)$ . Let k be the number of successes in the t trials and

$$\tilde{k} = \begin{cases} k & \text{if } p \ge \frac{1}{2} \\ t - k & \text{if } p < \frac{1}{2} \end{cases}.$$

Finally, we define the event  $\mathcal{E}$  as

$$\mathcal{E} = \left\{ \tilde{p}t - \tilde{k} \le \sqrt{2p(1-p)\log(\frac{c_2}{2\theta})} \right\},\,$$

where  $c_2 \ge 2$  is any constant.

**Theorem 9.** For  $c_1=32$ , it holds that  $\frac{L_1(W)}{L_0(W)}1_{\mathcal{E}}\geq \frac{2\theta}{c_2}1_{\mathcal{E}}$ .

*Proof.* We distinguish two cases depending on the value of p.

**Case 1:**  $p \ge \frac{1}{2}$ . The likelihood ratio can be written as

$$\frac{L_1(W)}{L_2(W)} = \frac{(p+\alpha)^k (1-p-\alpha)^{t-k}}{p^k (1-p)^{t-k}} = \left(1 + \frac{\alpha}{p}\right)^k \left(1 - \frac{\alpha}{1-p}\right)^{t-k}$$
$$= \left(1 + \frac{\alpha}{p}\right)^k \left(1 - \frac{\alpha}{1-p}\right)^{k \cdot \frac{1-p}{p}} \left(1 - \frac{\alpha}{1-p}\right)^{t - \frac{k}{p}}.$$

We start by bounding the second factor using that  $\log(1-x) \ge -x - x^2 + x^3$  for  $x \in [0, \frac{1}{5}]$  (Lemma 5) and that  $\exp(x) \ge 1 + x$  for all x along with our assumption that  $\alpha \le \frac{1-p}{5}$ :

$$\left(1 - \frac{\alpha}{1 - p}\right)^{\frac{1 - p}{p}} \ge \exp\left(\frac{1 - p}{p}\left[-\frac{\alpha}{1 - p} - \frac{\alpha^2}{(1 - p)^2} + \frac{\alpha^3}{(1 - p)^3}\right]\right) 
\ge 1 - \frac{1 - p}{p}\left[\frac{\alpha}{1 - p} + \frac{\alpha^2}{(1 - p)^2} - \frac{\alpha^3}{(1 - p)^3}\right] 
= 1 - \frac{\alpha}{p} - \frac{\alpha^2}{p(1 - p)} + \frac{\alpha^3}{p(1 - p)^2} 
\ge 1 - \frac{\alpha}{p} - \frac{\alpha^2}{p(1 - p)} + \frac{\alpha^3}{p^2(1 - p)} 
= \left(1 - \frac{\alpha}{p}\right)\left(1 - \frac{\alpha^2}{p(1 - p)}\right),$$

where we have used that  $p \ge 1 - p$ .

Using this along with the fact that  $k \le t$  and  $p \ge 1 - p$ , it follows that

$$\frac{L_1(W)}{L_0(W)} \ge \left(1 - \frac{\alpha^2}{p^2}\right)^k \left(1 - \frac{\alpha^2}{p(1-p)}\right)^k \left(1 - \frac{\alpha}{1-p}\right)^{t-\frac{k}{p}}$$

$$\ge \left(1 - \frac{\alpha^2}{p(1-p)}\right)^{2k} \left(1 - \frac{\alpha}{1-p}\right)^{t-\frac{k}{p}}$$

$$\ge \left(1 - \frac{\alpha^2}{p(1-p)}\right)^{2t} \left(1 - \frac{\alpha}{1-p}\right)^{t-\frac{k}{p}}.$$

Note that we have  $\alpha^2 \leq \frac{(1-p)^2}{25} \leq \frac{p(1-p)}{25} \leq \frac{p(1-p)}{2}$ . Using this and the fact that  $\log(1-x) \geq -2x$  for  $x \in [0, \frac{1}{2}]$ , we obtain

$$\left(1 - \frac{\alpha^2}{p(1-p)}\right)^{2t} \ge \exp\left(-4t\frac{\alpha^2}{p(1-p)}\right)$$
$$= \theta^{\frac{4}{c_1}}$$
$$\ge \left(\frac{2\theta}{c_2}\right)^{\frac{4}{c_1}},$$

where we have used that  $\frac{2}{c_2} \ge 1$ .

Now on the event  $\mathcal{E}$ , we have that  $t-\frac{k}{p} \leq \sqrt{2\frac{1-p}{p}t\log(\frac{c_2}{2\theta})}$ . Using this along with the fact that  $\frac{1}{c_1}\log(\frac{c_2}{2\theta}) \leq \frac{\alpha^2t}{p(1-p)}$ , which follows since

$$\log(\frac{c_2}{2\theta}) = \log\left(\frac{c_2}{2}\exp\left[\frac{c_1\alpha^2t}{p(1-p)}\right]\right) \le \log\left(\exp\left[\frac{c_1\alpha^2t}{p(1-p)}\right]\right) = \frac{c_1\alpha^2t}{p(1-p)},$$

we obtain that

$$\begin{split} \left(1 - \frac{\alpha}{1 - p}\right)^{t - \frac{k}{p}} &\geq \left(1 - \frac{\alpha}{1 - p}\right)^{\sqrt{2\frac{1 - p}{p}t\log(c_2/(2\theta))}} \\ &\geq \exp\left(-2\frac{\alpha}{1 - p}\sqrt{2\frac{1 - p}{p}t\log(c_2/(2\theta))}\right) \\ &= \exp\left(-2\sqrt{2}\sqrt{\frac{\alpha^2 t}{p(1 - p}\log(c_2/(2\theta))}\right) \\ &\geq \exp\left(-\frac{2\sqrt{2}}{\sqrt{c_1}}\log(c_2/(2\theta))\right) \\ &= \left(\frac{2\theta}{c_2}\right)^{\frac{2\sqrt{2}}{\sqrt{c_1}}}. \end{split}$$

Putting these together, we see that

$$\frac{L_1(W)}{L_2(W)} 1_{\mathcal{E}} \ge \left(\frac{2\theta}{c_2}\right)^{\frac{2\sqrt{2}}{\sqrt{c_1}} + \frac{2(1-p)}{p \cdot c_1} + \frac{2}{c_1}} 1_{\mathcal{E}},$$

so that choosing  $c_1 = 32$  yields the claimed result:

$$\frac{L_1(W)}{L_2(W)} 1_{\mathcal{E}} \ge \frac{2\theta}{c_2} 1_{\mathcal{E}} .$$

Case 2:  $p < \frac{1}{2}$ . Define m = t - k, which is now the number of failed coin flips. Hence,

$$\frac{L_1(W)}{L_0(W)} = \frac{(1-p-\alpha)^m (p+\alpha)^{t-m}}{(1-p)^m p^{t-m}} = \left(1 - \frac{\alpha}{1-p}\right)^m \left(1 + \frac{\alpha}{p}\right)^{t-m}$$
$$= \left(1 - \frac{\alpha}{1-p}\right)^m \left(1 + \frac{\alpha}{p}\right)^{m\frac{p}{1-p}} \left(1 + \frac{\alpha}{p}\right)^{t-\frac{m}{1-p}}.$$

Again, using  $\exp(1+x) \ge x$  for all  $x \in \mathbb{R}$  and using that  $\log(1+x) \ge x - x^2$  for all  $x \ge 0$ , we get that

$$\left(1 + \frac{\alpha}{p}\right)^{\frac{p}{1-p}} \ge \exp\left(\frac{p}{1-p}\left[\frac{\alpha}{p} - \frac{\alpha^2}{p^2}\right]\right)$$

$$\ge 1 + \frac{\alpha}{1-p} - \frac{\alpha^2}{p(1-p)}$$

$$\ge 1 + \frac{\alpha}{1-p} - \frac{\alpha^2}{p(1-p)} - \frac{\alpha^3}{p(1-p)^2}$$

$$= \left(1 + \frac{\alpha}{1-p}\right)\left(1 - \frac{\alpha^2}{p(1-p)}\right).$$

Using this along with the fact that (1-p) > p and  $m \le t$ , we have

$$\frac{L_1(W)}{L_2(W)} \ge \left(1 - \frac{\alpha^2}{(1-p)^2}\right)^m \left(1 - \frac{\alpha^2}{p(1-p)}\right)^m \left(1 + \frac{\alpha}{p}\right)^{t - \frac{m}{1-p}} \\
\ge \left(1 - \frac{\alpha^2}{p(1-p)}\right)^{2t} \left(1 - \frac{\alpha}{p}\right)^{t - \frac{m}{1-p}}.$$

Again, using  $\log(1-x) \ge -2x$  for  $x \in [0, \frac{1}{2}]$ , we get that

$$\left(1 - \frac{\alpha^2}{p(1-p)}\right)^{2t} \ge \exp\left(-4t\frac{\alpha^2}{p(1-p)}\right)$$

$$\ge \theta^{\frac{4}{c_1}}$$

$$\ge \left(\frac{2\theta}{c_2}\right)^{\frac{4}{c_1}}.$$

On the event  $\mathcal{E}$ , we have that  $t - \frac{m}{1-p} \leq \sqrt{\frac{2pt\alpha^2}{1-p}\log(\frac{c_2}{2\theta})}$ . Using this along with the fact that  $\frac{1}{c_1}\log(\frac{c_2}{2\theta}) \leq \frac{\alpha^2t}{p(1-p)}$ , we get on the event  $\mathcal{E}$  that

$$\left(1 - \frac{\alpha}{p}\right)^{t - \frac{m}{1 - p}} \ge \left(1 - \frac{\alpha}{p}\right)^{\sqrt{\frac{2p}{1 - p}t\log(\frac{c_2}{2\theta})}}$$

$$\ge \exp\left(-2\sqrt{\frac{2t}{p(1 - p)}\log(\frac{c_2}{2\theta})}\right)$$

$$\ge \exp\left(-\frac{2\sqrt{2}}{\sqrt{c_1}}\log(\frac{c_2}{2\theta})\right)$$

$$= \left(\frac{2\theta}{c_2}\right)^{\frac{2\sqrt{2}}{\sqrt{c_1}}}.$$

We thus get the desired result for  $c_1 = 32$ :

$$\frac{L_1(W)}{L_0(W)} 1_{\mathcal{E}} \ge 1_{\mathcal{E}} \left(\frac{2\theta}{c_2}\right)^{\frac{4}{c_1} + \frac{2\sqrt{2}}{\sqrt{c_1}}} \ge 1_{\mathcal{E}} \left(\frac{2\theta}{c_2}\right).$$

## H Proofs of Lower Bounds

## H.1 Lower Bound for Q-value Learning

For a lower bound we construct the following class of MDPs with S':=S+2 states and A actions where the first states are labelled  $S_1,\ldots,s_S,s^G,s^B$  and the actions are labelled  $a_1,\ldots,a_A$ . The states  $s^G$  and  $s^B$  are absorbing under any actions and  $R(s^G,a)=1$  for all j and  $R(s^B,a)=0$  for all  $a\in A$ . For the states  $s\in\{s_1,\ldots,s_S\}$ , we have that R(s,a)=0 for all  $a\in A$ . We have SA state-action pair combinations from  $\{s_1,\ldots,s_S\}\times A=:Z$  on which we assume some ordering allowing us to write  $z_i,i\in[SA]$ . Finally for all state-action pairs  $z_i\in[SA]$  we have  $P(s^G|z_i)=q_i$  and  $P(s^B|z_i)=1-q_i$  for some  $q_i\in[0,1]$ . The structure of this class of MDPs allows us to get lower bounds on the samples needed to learn the Q-value of each state-action pair  $z_i$  and then use the fact that samples used to learn the Q-values for different state-action pairs bring no information on eachother to get the final bound.

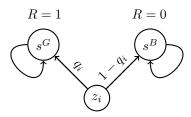


Figure 1: Dynamics and rewards of the hard-to-learn MDP class

For any state-action pair we can explicitly calculate the state-action value-functions

$$\begin{split} Q(z_i) &= \frac{-\gamma}{\beta} \log(q_i e^{-\beta \frac{1}{1-\gamma}} + 1 - q_i), \\ Q(s^G, a) &= \frac{1}{1-\gamma}, \\ Q(s^B, a) &= 0. \end{split}$$

Denote the collection of all such MDPs by M.

Fix any index i and consider the two hypotheses  $H_0^i:q_i=p$  and  $H_1^i:q_i=p+\alpha$  where p and  $\alpha$  are given by

$$p = \begin{cases} 1 - e^{-\beta \frac{1}{1-\gamma}} & \text{for } \beta > 0, \\ e^{-|\beta| \frac{1}{1-\gamma}} & \text{for } \beta < 0, \end{cases}$$

and

$$\alpha = 8\varepsilon \frac{|\beta|}{\gamma} \frac{1}{e^{|\beta| \frac{1}{1-\gamma}} - 1},$$

for any  $\varepsilon$  in the range  $\varepsilon < \frac{1}{40} \frac{\gamma}{|\beta|} (1 - e^{-|\beta|})$  .

We use  $M_0$  to denote an MDP where  $H_0^i$  holds and  $M_1$  to denote an MDP where instead  $H_0^i$  holds and  $\mathbb{E}_0$  and  $\mathbb{P}_0$  as the expectations operator and probability operator under  $H_1^i$  and similarly  $\mathbb{E}_1$  and  $\mathbb{P}_1$  under  $H_0^i$ . Fix any  $(\varepsilon, \delta)$ -correct Q-algorithm  $\mathcal{U}$ . We start by showing that with these parameter we have that  $Q_{M_1}^*(z_i) - Q_{M_0}^*(z_i) > 2\varepsilon$ , which we do by casing on the sign of  $\beta$ :

Case 1:  $\beta < 0$ . In this case  $p = e^{-|\beta| \frac{1}{1-\gamma}}$ . We then have

$$\begin{split} Q_{M_1}^*(z_i) - Q_{M_0}^*(z_i) &= \frac{\gamma}{|\beta|} \log \left( \frac{(p+\alpha)e^{|\beta|\frac{1}{1-\gamma}} + 1 - p - \alpha}{pe^{|\beta|\frac{1}{1-\gamma}} + 1 - p} \right) \\ &= \frac{\gamma}{|\beta|} \log \left( 1 + \frac{\alpha(e^{|\beta|\frac{1}{1-\gamma}} - 1)}{pe^{|\beta|\frac{1}{1-\gamma}} + 1 - p} \right) \\ &\geq \frac{\gamma}{|\beta|} \frac{\alpha}{2} \frac{e^{|\beta|\frac{1}{1-\gamma}} - 1}{pe^{|\beta|\frac{1}{1-\gamma}} + 1 - p} \\ &\geq \frac{\gamma}{|\beta|} \frac{\alpha}{4} (e^{|\beta|\frac{1}{1-\gamma}} - 1) \\ &= 2\varepsilon \,, \end{split}$$

where we have used that  $p = e^{-|\beta| \frac{1}{1-\gamma}}$  and the fact that  $\log(1+x) \ge \frac{x}{2}$  for  $x \in [0,1]$ .

Case 2:  $\beta > 0$ . The case for  $\beta > 0$  is similar, although in this case we have  $p = 1 - e^{-\beta \frac{1}{1 - \gamma}}$  and use the inequality  $\log(1 + x) \le x$  for all x > -1 to get that

$$Q_{M_{1}}^{*}(z_{i}) - Q_{M_{0}}^{*}(z_{i}) = -\frac{\gamma}{\beta} \log \left( \frac{(p+\alpha)e^{-\beta\frac{1}{1-\gamma}} + 1 - p - \alpha}{pe^{-\beta\frac{1}{1-\gamma}} + 1 - p} \right)$$

$$= -\frac{\gamma}{\beta} \log \left( 1 - \frac{\alpha(1 - e^{-\beta\frac{1}{1-\gamma}})}{1 - p + pe^{-\beta\frac{1}{1-\gamma}}} \right)$$

$$= -\frac{\gamma}{\beta} \log \left( 1 - \frac{\alpha(1 - e^{-\beta\frac{1}{1-\gamma}})}{(1 - p)e^{-\beta\frac{1}{1-\gamma}}} \right)$$

$$\geq \frac{\gamma}{\beta} \alpha \frac{1 - e^{-\beta\frac{1}{1-\gamma}}}{(1 + p)e^{-\beta\frac{1}{1-\gamma}}}$$

$$\geq \frac{\gamma}{\beta} \alpha \frac{1 - e^{-\beta\frac{1}{1-\gamma}}}{2e^{-\beta\frac{1}{1-\gamma}}}$$

$$\geq \frac{\gamma}{\beta} \alpha \frac{e^{\beta\frac{1}{1-\gamma}} - 1}{2}$$

$$= 4\varepsilon$$

In particular, this means that the events  $B_0 := \{|Q_{M_0}^*(z_i) - Q_t^{\mathcal{U}}(z_i)| \leq \varepsilon\}$  and  $B_1 := \{|Q_{M_1}^*(z_i) - Q_t^{\mathcal{U}}(z_i)| \leq \varepsilon\}$  are disjoint events Let t be the number of times the algorithm tries  $z_i$ . Since  $\mathcal{U}$  is  $(\varepsilon, \delta)$ -correct it holds that  $\mathbb{P}_0(B_0) \geq 1 - \delta \geq \frac{3}{4}$ .

Let k be the number of transitions from  $z_i$  to  $z_i^G$  in the t trials. We then define  $\tilde{k}, \tilde{p}$  and  $\theta$  by

$$\begin{split} \tilde{k} &:= \begin{cases} k & \text{if } p \geq \frac{1}{2} \\ t - k & \text{if } p < \frac{1}{2} \end{cases} \\ \theta &:= \exp \big( -\frac{32\alpha^2 t}{p(1-p)} \big) \\ \tilde{p} &= \max\{p, 1-p\} \end{split}$$

and the event

$$\mathcal{E} = \left\{ \tilde{p}t - \tilde{k} \le \sqrt{2p(1-p)t\log(\frac{8}{2\theta})} \right\}$$

for which we have  $\mathbb{P}_0(\mathcal{E}) > \frac{3}{4}$  by Lemma 16 in [21] and thus  $\mathbb{P}_0(B_0 \cap \mathcal{E}) > \frac{1}{2}$ . Now by Theorem 9 we get that

$$\mathbb{P}_1(B_0) \geq \mathbb{P}_1(B_0 \cap \mathcal{E}) = \mathbb{E}_1[1_{\mathcal{E}}1_{B_0}] = \mathbb{E}_0\left[\frac{L_1}{L_0}1_{\mathcal{E}}1_{B_0}\right] \geq \frac{\theta}{4}\mathbb{E}_0[1_{\mathcal{E}}1_{B_0}] = \frac{\theta}{4}\mathbb{P}_0(\mathcal{E} \cap B_0) \geq \frac{\theta}{8}.$$

Solving for t in  $\frac{\theta}{8} > \delta$  we find

$$t < \frac{p(1-p)}{32\alpha^2}\log(\frac{1}{8\delta}),$$

and since

$$\frac{p(1-p)}{\alpha^2} = \frac{\gamma^2}{|\beta|^2} \frac{e^{-|\beta|\frac{1}{1-\gamma}}(1-e^{-|\beta|\frac{1}{1-\gamma}})}{64\varepsilon^2} (e^{|\beta|\frac{1}{1-\gamma}} - 1)^2$$
$$\geq \frac{\gamma^2}{64\varepsilon^2} \frac{e^{|\beta|\frac{1}{1-\gamma}} - 3}{|\beta|^2},$$

we conclude that if the algorithm  $\mathcal{U}$  tries the state-action pair  $z_i$  less than

$$\tilde{T}(\varepsilon,\delta) := \frac{\gamma^2}{2048\varepsilon^2} \frac{e^{|\beta|\frac{1}{1-\gamma}} - 3}{|\beta|^2} \log(\frac{1}{8\delta})$$

times under the hypothesis  $H_0^i$ , then  $\mathbb{P}_1(B_0) > \delta$  and  $B_0 \subset B_1^c$ .

Next we use the fact that the structure of the MDPs is such that information on the Q-value of any state-action pair in Z carries no information on the Q-values of any other state-action pair in Z.

Let n:=SA. If the number of total transition samples is less than  $\frac{n}{2}\tilde{T}(\varepsilon,\delta)$  there must be at least n/2 state-action pairs  $z_i$  that has been tried no more than  $\tilde{T}(\varepsilon,\delta)$  times which without loss of generality we might assume are the state-action pairs  $\{z_i\}_{i=1}^{n/2}$ .

Let  $T_i$  be the number of times the algorithm has tried  $z_i$  for  $i \leq n/2$  Due to the structure of the MDPs in  $\mathbb{M}$  it is sufficient to consider only the algorithms that outputs an estimate of  $Q_{T_i}^{\mathcal{U}}$  based on samples from  $z_i$  since any other samples can yield no information on  $Q^*(z_i)$ 

Thus by defining the events  $\Lambda_i := \{|Q_{M_1}^*(z_i) - Q_{T_i}^{\mathcal{U}}(z_i)| > \varepsilon\}$  we have that  $\Lambda_i$  and  $\Lambda_j$  are conditionally independent given  $T_i$  and  $T_j$ . We then have

$$\begin{split} & \mathbb{P}_{1}(\{\Lambda_{i}^{c}\}_{1 \leq i \leq n/2} \cap \{T_{i} \leq \tilde{T}(\varepsilon, \delta)\}_{1 \leq i \leq n/2}) \\ & = \sum_{t_{1}=0}^{\tilde{T}(\varepsilon, \delta)} \cdots \sum_{t_{n/2}=0}^{\tilde{T}(\varepsilon, \delta)} \mathbb{P}_{1}(\{T_{i} = t_{i}\}_{1 \leq i \leq n/2}) \mathbb{P}_{1}(\{\Lambda_{i}^{c}\}_{1 \leq i \leq n/2} \cap \{T_{i} = t_{i}\}_{1 \leq i \leq n/2}) \\ & = \sum_{t_{1}=0}^{\tilde{T}(\varepsilon, \delta)} \cdots \sum_{t_{n/2}=0}^{\tilde{T}(\varepsilon, \delta)} \mathbb{P}_{1}(\{T_{i} = t_{i}\}_{1 \leq i \leq n/2}) \prod_{1 \leq i \leq n/2} \mathbb{P}_{1}(\Lambda_{i}^{c} \cap \{T_{i} = t_{i}\}) \\ & = \sum_{t_{1}=0}^{\tilde{T}(\varepsilon, \delta)} \cdots \sum_{t_{n/2}=0}^{\tilde{T}(\varepsilon, \delta)} \mathbb{P}_{1}(\{T_{i} = t_{i}\}_{1 \leq i \leq n/2}) (1 - \delta)^{n/2}, \end{split}$$

where we have used the law of total probability from line one to two and from two to three follows from independence. We now have directly that

$$\mathbb{P}_1(\{\Lambda_i^c\}_{1 < i < n/2} | \{T_i \le \tilde{T}(\varepsilon, \delta)\}_{1 < i < n/2}) \le (1 - \delta)^{\frac{n}{2}}.$$

Thus, if the total number of transitions T is less than  $\frac{n}{2}\tilde{T}(\varepsilon,\delta)$ , then

$$\mathbb{P}_{1}(\|Q^{*} - Q_{T}^{\mathcal{U}}\| > \varepsilon) \geq \mathbb{P}_{1}\left(\bigcup_{z \in S \times A} \Lambda(z)\right)$$

$$= 1 - \mathbb{P}_{1}\left(\bigcap_{1 \leq i \leq n/2} \Lambda_{i}^{c}\right)$$

$$\geq 1 - \mathbb{P}_{1}(\{\Lambda_{i}^{c}\}_{1 \leq i \leq n/2} | \{T_{z_{i}} \leq \tilde{T}(\varepsilon, \delta)\}_{1 \leq i \leq n/2})$$

$$\geq 1 - (1 - \delta)^{n/2}$$

$$\geq \frac{\delta n}{4},$$

when  $\delta \frac{n}{2} \le 1$  by Lemma 6. By setting  $\delta' = \delta \frac{n}{4}$  and substituting back S' we obtain the result. This shows that if the number of samples is smaller than

$$T = \frac{(S'-2)A}{4096} \frac{\gamma^2}{\varepsilon^2} \frac{e^{|\beta|\frac{1}{1-\gamma}} - 3}{|\beta|^2} \log(\frac{(S'-2)A}{32\delta})$$
 (71)

on the MDP corresponding to the hypothesis  $H_0:\{H_0^i|1\leq i\leq n\}$  it holds that  $\mathbb{P}_1(\|Q_{M_1}^*-Q_T^{\mathcal{U}}\|>\varepsilon)>\delta'$ .

## **H.2** Lower Bound for Policy Learning

For a lower bound we construct the following class of MDPs with S':=S+2 states and A':=A+1 actions where the first states are labelled  $s_1,\ldots,s_S,s^G,s^B$  and the actions are labelled  $a_0,a_1,\ldots,a_A$ . The states  $s^G$  and  $s^B$  are absorbing under any actions and  $R(s^G,a)=1$  for all j and  $R(s^B,a)=0$  for all  $a\in A$ . For the states  $s\in\{s_1,\ldots,s_S\}$ , we have that R(s,a)=0 for all  $a\in A$ .

From the state  $s_i$  with probabilities that depend on the action taken the agent will then end up in either a good state  $s^G$  which is absorbing and yields the maximal unit reward under all actions or in the bad state  $s^B$  which is also absorbing but which yields no reward under any action. The different MDPs thus differ only in their transition probabilities in the choice states  $s_i$ .

Fix an index  $1 \le i \le S$ . We then consider the following set of possible parameters called hypotheses  $H_l^i, l \in \{0, 1, 2, \dots, A\}$  given by

$$\begin{split} H_0^i : & q(s_i, a_0) = p + \alpha & q(s_i, a) = p \text{ for } a \neq a_0 \\ H_l : & q(s_i, a_0) = p + \alpha & q(s_i, a) = p \text{ for } a \notin \{a_0, l\} \\ \end{split} \qquad \qquad q(s_i, a_l) = p + 2\alpha \,,$$

where p and  $\alpha$  are given by

$$p = \begin{cases} 1 - e^{-\beta \frac{1}{1-\gamma}} & \beta > 0, \\ e^{-|\beta| \frac{1}{1-\gamma}} & \beta < 0, \end{cases}$$
$$\alpha = \frac{5|\beta|}{\gamma} \frac{\varepsilon}{e^{|\beta| \frac{1}{1-\gamma}} - 1},$$

where we allow for

$$0<\varepsilon<\frac{\gamma}{50|\beta|}\big(1-e^{-|\beta|\frac{1}{1-\gamma}}\big)\,,$$

which ensures that  $\alpha \leq \frac{e^{-|\beta|\frac{1}{1-\gamma}}}{10}$ 

Consider a fixed hypothesis  $H_l^i$  for some  $l \neq 0$  and the sub-MDP that only consists of the states  $\{s_i, s^G, s^B\}$ . Here the optimal action is  $a^* = a_l$ , the second best action is  $a_0$  and all other actions are even worse so the value-error over all states in the triplet for any suboptimal choice of actions will be at least as large as  $V^*(s_i) - V^0(s_i)$  where  $V^0$  is the value by choosing a = 0. We now show that any non-optimal action is  $\varepsilon$ -bad on  $s_i$ .

**Case 1:**  $\beta > 0$ .

$$V^{*}(s_{i}) - V^{0}(s_{i}) = -\frac{\gamma}{\beta} \log \left( \frac{(p+2\alpha)e^{-\beta\frac{1}{1-\gamma}} + 1 - p - 2\alpha}{(p+\alpha)e^{-\beta\frac{1}{1-\gamma}} + 1 - p - \alpha} \right)$$

$$= \frac{-\gamma}{\beta} \log \left( 1 - \alpha \frac{1 - e^{-\beta\frac{1}{1-\gamma}}}{pe^{-\beta\frac{1}{1-\gamma}} + 1 - p - \alpha(1 - e^{-\beta\frac{1}{1-\gamma}})} \right)$$

$$> \frac{\gamma}{\beta} \alpha \frac{1 - e^{-\beta\frac{1}{1-\gamma}}}{pe^{-\beta\frac{1}{1-\gamma}} + 1 - p - \alpha(1 - e^{-\beta\frac{1}{1-\gamma}})}$$

$$\geq \frac{\gamma}{\beta} \alpha \frac{1 - e^{-\beta\frac{1}{1-\gamma}}}{pe^{-\beta\frac{1}{1-\gamma}} + 1 - p}$$

$$= \frac{\gamma}{\beta} \alpha \frac{1 - e^{-\beta\frac{1}{1-\gamma}}}{(1+p)e^{-\beta\frac{1}{1-\gamma}}}$$

$$\geq \frac{\gamma}{\beta} \alpha \frac{1 - e^{-\beta\frac{1}{1-\gamma}}}{2e^{-\beta\frac{1}{1-\gamma}}}$$

$$\geq \frac{\gamma}{2\beta} \alpha(1 - e^{-\beta\frac{1}{1-\gamma}})$$

$$\geq \varepsilon,$$

where we have used  $\log(1+x) > x$  for  $x \in (-1, \infty) \setminus \{0\}$ .

**Case 2:**  $\beta$  < 0.

$$V^{*}(s_{i}) - V^{0}(s_{i}) = \frac{\gamma}{|\beta|} \log \left( \frac{(p+2\alpha)e^{|\beta|\frac{1}{1-\gamma}} + 1 - p - 2\alpha}{(p+\alpha)e^{|\beta|\frac{1}{1-\gamma}} + 1 - p - \alpha} \right)$$

$$= \frac{\gamma}{|\beta|} \log \left( 1 + \alpha \frac{e^{|\beta|\frac{1}{1-\gamma}} - 1}{pe^{-\beta\frac{1}{1-\gamma}} + 1 - p + \alpha(e^{|\beta|\frac{1}{1-\gamma}} - 1)} \right)$$

$$> \frac{\gamma}{2|\beta|} \alpha \frac{e^{|\beta|\frac{1}{1-\gamma}} - 1}{pe^{-\beta\frac{1}{1-\gamma}} + 1 - p + \alpha(e^{|\beta|\frac{1}{1-\gamma}} - 1)}$$

$$\geq \frac{\gamma}{2|\beta|} \alpha \frac{e^{|\beta|\frac{1}{1-\gamma}} - 1}{2 + \alpha(e^{|\beta|\frac{1}{1-\gamma}} - 1)}$$

$$\geq \frac{\gamma}{2|\beta|} \alpha \frac{e^{|\beta|\frac{1}{1-\gamma}} - 1}{2 + \frac{1}{10}}$$

$$\geq \frac{\gamma}{2|\beta|} \alpha (e^{|\beta|\frac{1}{1-\gamma}} - 1)$$

$$\geq \varepsilon,$$

where we have used  $\log(1+x) > \frac{x}{2}$  for  $x \in (0,1)$ .

Now haven shown that all non-optimal actions are  $\varepsilon$ -bad, we wish to show that any algorithm that is  $(\varepsilon, \delta)$ -correct on  $H_0^i$ , i.e. choosing the action  $a_0$  with probability at least  $1 - \delta$ , will also have a probability of choosing  $a_0$  on  $H_l^i$  that is larger than  $\delta$  provided that  $a_l$  is not tried sufficiently many times under  $H_0^i$ .

Let  $\mathbb{P}_l$  and  $\mathbb{E}_l$  denote the probability operator and expectation operator under the hypothesis  $H_l^i$ . Let  $t:=t_l^i$  be the number of times the algorithm tries action l in  $s_i$  under  $H_0$ . Assuming that  $\delta \in (0,\frac{1}{4})$  and using that the algorithm is  $(\varepsilon,\delta)$ -correct we have that  $\mathbb{P}_0(B) \geq 1 - \delta \geq \frac{3}{4}$  where  $B = \{\pi^{\mathcal{U}}(s_i) = a_0\}$  is the event that the algorithm outputs the action  $a_0$ . Let  $\theta = \exp\left(-\frac{32\alpha^2t}{p(1-p)}\right)$ . Fix some  $t \in \mathbb{N}$  and let k be the number of transitions to  $s_i^G$  in the t trials and

$$\tilde{k} = \begin{cases} k & \text{if } p \ge \frac{1}{2} \\ t - k & \text{if } p < \frac{1}{2}. \end{cases}$$

Finally, we define the event  $\mathcal{E}$  as

$$\mathcal{E} = \left\{ \tilde{p}t - \tilde{k} \le \sqrt{2p(1-p)\log(\frac{8}{2\theta})} \right\}. \tag{72}$$

Form the Chernoff-Hoeffding bound and as shown in [21] we have that  $\mathbb{P}_0(\mathcal{E}) > \frac{3}{4}$  and so  $\mathbb{P}_0(B \cap \mathcal{E}) > \frac{1}{2}$ . From Theorem 9, we get that

$$\mathbb{P}_{1}(B) \geq \mathbb{P}_{1}(B \cap \mathcal{E}) = \mathbb{E}_{1}[1_{B}1_{\mathcal{E}}] \geq \mathbb{E}_{0}\left[\frac{L_{1}(W)}{L_{0}(W)}1_{\mathcal{E}}1_{B}\right] \geq \mathbb{E}_{0}\left[\frac{\theta}{4}1_{\mathcal{E}}1_{B}\right] = \frac{\theta}{4}\mathbb{P}_{0}(\mathcal{E} \cap B) \geq \frac{\theta}{8}.$$
(73)

Now solving for  $\frac{\theta}{8} > \delta$ , we see that if

$$t < \tilde{t}(\varepsilon, \delta) := \frac{1}{800} \log(\frac{1}{8\delta}) \frac{\gamma^2}{\varepsilon^2} \cdot \frac{e^{|\beta| \frac{1}{1-\gamma}} - 3}{|\beta|^2}$$
 (74)

then  $\mathbb{P}_1(B) > \delta$  and the event B is containing the event that the algorithm does not choose the optimal action  $a_l$ .

Since this holds for all the A hypotheses  $H_l^i, l=1,2,...,A$  it follows that the algorithm needs at least  $\tilde{T}(\varepsilon,\delta):=A\tilde{t}(\varepsilon,\delta)$  samples to be  $(\varepsilon,\delta)$ -correct on the state  $s_i$ .

Next we use the fact that the structure of the MDPs is such that information used to determine  $\pi^*(s_i)$  carries no information to determine  $\pi^*(s_i)$  for  $i \neq j$ .

If the number of total transition samples is less than  $\frac{S}{2}\tilde{T}(\varepsilon,\delta)$  then there must be at least  $\frac{S}{2}$  states in the set  $\{s_i\}_{i=1}^S$  for which some action (apart from  $a_0$ ) has been tried no more than  $\tilde{T}(\varepsilon,\delta)$  times which without loss of generality we might assume are the states  $\{s_i\}_{i=1}^{S/2}$ . and that it is action  $a_1$  that has been tried out at most  $\tilde{T}(\varepsilon,\delta)$  times in each of these states.

Let  $T_i$  be the number of times the algorithm has tried sampled any action on  $s_i$  for  $i \leq S/2$  Due to the structure of the MDPs in  $\mathbb M$  it is sufficient to consider only the algorithms that yields an estimate of  $\pi^{\mathcal U}_{T_i}$  based on samples from  $s_i$  since any other samples can yield no information on  $\pi^*(s_i)$ .

Thus, by defining the events  $\Lambda_i := \{|V_{M_1}^*(s_i) - V^{\pi_{T_i}^{\mathcal{U}}}(s_i)| > \varepsilon\}$  we have that  $\Lambda_i$  and  $\Lambda_j$  are conditionally independent given  $T_i$  and  $T_j$ . We then have that for the MDP  $M_1 \in \mathbb{M}$  (The one corresponding to the hypothesis  $H_1 := \{H_1^i | 1 \le i \le n\}$ ) it holds that

$$\begin{split} & \mathbb{P}(\{\Lambda_i^c\}_{1 \leq i \leq S/2} \cap \{T_i \leq \tilde{T}(\varepsilon, \delta)\}_{1 \leq i \leq S/2}) \\ & = \sum_{t_1 = 0}^{\tilde{T}(\varepsilon, \delta)} \cdots \sum_{t_{S/2} = 0}^{\tilde{T}(\varepsilon, \delta)} \mathbb{P}(\{T_i = t_i\}_{1 \leq i \leq S/2}) \mathbb{P}(\{\Lambda_i^c\}_{1 \leq i \leq S/2} \cap \{T_i = t_i\}_{1 \leq i \leq S/2}) \\ & = \sum_{t_1 = 0}^{\tilde{T}(\varepsilon, \delta)} \cdots \sum_{t_{S/2} = 0}^{\tilde{T}(\varepsilon, \delta)} \mathbb{P}(\{T_i = t_i\}_{1 \leq i \leq S/2}) \prod_{1 \leq i \leq S/2} \mathbb{P}(\Lambda_i^c \cap \{T_i = t_i\}) \\ & = \sum_{t_1 = 0}^{\tilde{T}(\varepsilon, \delta)} \cdots \sum_{t_{S/2} = 0}^{\tilde{T}(\varepsilon, \delta)} \mathbb{P}(\{T_i = t_i\}_{1 \leq i \leq S/2}) (1 - \delta)^{S/2} \,, \end{split}$$

where we have used the law of total probability from line one to two and from two to three follows from independence. We now have directly that

$$\mathbb{P}\left(\left\{\Lambda_i^c\right\}_{1 \le i \le S/2} \middle| \left\{T_i \le \tilde{T}(\varepsilon, \delta)\right\}_{1 \le i \le S/2}\right) \le (1 - \delta)^{\frac{S}{2}}.$$

Thus, if the total number of transitions T is less than  $\frac{S}{2}\tilde{T}(\varepsilon,\delta)$  on the MDP  $M_0$  corresponding to the hypothesis  $H_0:\{H_0^i|1\leq i\leq n\}$ , then on  $M_1$  it holds that

$$\mathbb{P}(\|V^* - V^{\pi_T^{\mathcal{U}}}\| > \varepsilon) \ge \mathbb{P}\left(\bigcup_{1 \le i \le S/2} \Lambda(z)\right)$$

$$= 1 - \mathbb{P}\left(\bigcap_{1 \le i \le S/2} \Lambda_i^c\right)$$

$$\ge 1 - \mathbb{P}(\{\Lambda_i^c\}_{1 \le i \le S/2} | \{T_{z_i} \le \tilde{T}(\varepsilon, \delta)\}_{1 \le i \le S/2})$$

$$\ge 1 - (1 - \delta)^{S/2}$$

$$\ge \frac{\delta S}{4},$$

when  $\frac{\delta S}{2} \leq 1$  by Lemma 6. By setting  $\delta' = \delta \frac{S}{4}$  and substituting back S' and A' we obtain the result. This shows that if the number of samples is smaller than

$$T = \frac{(S'-2)(A'-1)}{1600} \log(\frac{S'-2}{32\delta}) \frac{\gamma^2}{\varepsilon^2} \cdot \frac{e^{|\beta| \frac{1}{1-\gamma}} - 3}{|\beta|^2}$$

on  $M_0$  then on  $M_1$  it holds that  $\mathbb{P}(\|V^* - V^{\pi_T^{\mathcal{U}}}\| > \varepsilon) > \delta$ .