
An Analysis of Abstracted Model-Based Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Many methods for Model-based Reinforcement Learning (MBRL) provide guaran-
2 tees for the accuracy of the Markov decision process (MDP) model they can deliver.
3 At the same time, state abstraction techniques allow for a reduction of the size of
4 an MDP while maintaining a bounded loss with respect to the original problem.
5 It may come as a surprise, therefore, that no such guarantees are available when
6 combining both techniques, i.e., where MBRL merely observes abstract states. Our
7 theoretical analysis shows that abstraction can introduce a dependence between
8 samples collected online (i.e., in the real world), which invalidates most results
9 for MBRLs in this setting. Collecting samples using a simulator can avoid this
10 problem. We conclude that we should be careful when applying MBRL methods
11 to abstracted real-world data.

12 1 Introduction

13 When trying to find good solutions to MDPs using Reinforcement Learning (RL) a fundamental
14 problem is the exploration-exploitation dilemma: when to take actions to obtain more information,
15 and when to take actions that maximize reward based on the current knowledge. Tabular MBRL
16 methods have found good ways to deal with this dilemma [7, 28, 14].

17 However, MDPs can be very large, which can be problematic for these methods. One way to deal
18 with this is to reduce the size of the MDP. State abstractions are one way to do this [17, 1]. We
19 are interested in approximate state abstractions since they allow for greater reductions of the MDP,
20 though there is a trade-off with solution quality [1]. Specifically, we assume we have an *approximate*
21 *model similarity abstraction function* ϕ [1] that maps states to abstract states. The environment
22 returns states s , but the agent receives $\phi(s)$, see Figure 1. This setting, which was considered before
23 [22, 2], is what we call *Abstracted RL*, and is the topic of this paper.

24 Abstracted RL corresponds to RL in a Partially
25 Observable MDP (POMDP), as previously de-
26 scribed [5]. It is well known that policies for
27 POMDPs that only base their action on the last
28 observation $\phi(s)$ could be arbitrarily bad [26].
29 However, when we assume that ϕ is an approx-
30 imate model similarity abstraction [1] this worst
31 case may not apply: Based on the observed ab-
32 stract states the agent learns an (empirical) abstract model. If we could show that this learned model
33 is close to an ‘abstract MDP’ (details in Section 2.2), we could give finite-sample guarantees on the
34 performance in the original MDP by combining results from MBRL and abstraction.

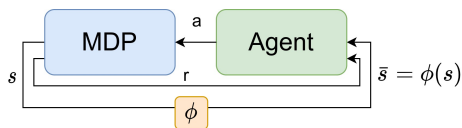


Figure 1: Abstracted RL, the agent observes $\bar{s} = \phi(s)$ instead of s . Image based on Abel et al. [2]

35 However, in MBRL, to guarantee (with high probability) that the learned model is close to the actual
 36 environment model, it is typical (e.g., [28, 14]) to use concentration inequalities such as Theorem
 37 2 from Weissman et al. [30]. But this theorem relies on independent and identically distributed
 38 (i.i.d.) samples for each state-action pair. In this paper, we analyze online collection of such samples
 39 in Abstracted RL and show that they are not independent¹, which means that *most guarantees for*
 40 *existing MBRL methods do not hold in the online Abstracted RL setting.*²

41 On the positive side, when we have access to a simulator, we show how this can be used to collect the
 42 data such that the typical MBRL guarantees hold and we can learn an accurate model. We discuss
 43 that emulating this in the real world is possible, but extremely sample inefficient, thus highlighting
 44 the difficulty of assuming that we would have access to an i.i.d. dataset, as in some earlier works.

45 2 Preliminaries

46 We assume the environment the agent is acting in can be represented by an infinite horizon MDP
 47 $M := \langle S, A, T, R, \gamma \rangle$. Where S is a finite set of states $s \in S$, A a finite set of actions $a \in A$, T
 48 a transition function $T(s'|s, a) = \Pr(s'|s, a)$, R a reward function $R(s, a)$ which gives the reward
 49 received when the agent executes action a in state s , and γ the discount factor ($0 \leq \gamma < 1$).

50 In RL the goal of the agent is to find an optimal policy $\pi^* : S \rightarrow A$ which maximizes the expectation
 51 of the discounted cumulative reward. $V^\pi(s)$ denotes the expected value of the discounted cumulative
 52 reward under policy π starting from state s . Similarly, $Q^\pi(s, a)$ denotes the expected value of the
 53 discounted cumulative reward when first taking action a from state s and then following policy π .

54 2.1 Model-Based RL

55 MBRL methods learn a model from the experiences, these are obtained by the agent acting in the
 56 MDP. The learned model is usually the empirical model, directly based on the experience the agent
 57 obtains [7, 28, 14]. Per state-action pair the agent stores the next-states reached after taking action a
 58 from state s in sequence $Y_{s,a} : Y_{s,a} : \{s^{(1)}, s^{(2)}, \dots, s^{(m)}\}$. We use Y to refer to the collection of
 59 all $Y_{s,a}$. From this the empirical, or learned, model T_Y is constructed, that just counts how often we
 60 have seen the transition to a next-state, and normalizes this:

$$\forall_{s' \in S} T_Y(s'|s, a) \triangleq \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{Y_{s,a}^{(i)} = s'\}, \quad (1)$$

61 where $\mathbb{1}\{\cdot\}$ denotes the indicator function of the specified event, i.e., it is 1 if $Y_{s,a}^{(i)} = s'$ and 0
 62 otherwise.

63 To give finite-sample guarantees on the accuracy of the estimate T_Y ,³ concentration bounds such as
 64 Theorem 2.1 from Weissman et al. [30] are often used, e.g. in Strehl and Littman [28], Jaksch et al.
 65 [14]. However, these typically make use of the fact that samples are i.i.d. In most MBRL settings this
 66 is not a problem under some assumptions, e.g. when the MDP is communicating [25]. In this case
 67 due to the Markov property the obtained samples are i.i.d.

68 In general, of course the hope is that with enough samples the learned model T_Y becomes accurate.
 69 With accurate we mean that the distance between $T_Y(\cdot|s, a)$ and $T(\cdot|s, a)$ will be small, where the
 70 distance is measured using the L_1 norm, defined as:

$$\|T_Y(\cdot|s, a) - T(\cdot|s, a)\|_1 \triangleq \sum_{s' \in S} |T_Y(s'|s, a) - T(s'|s, a)|. \quad (2)$$

71 Part of theorem 2.1 from Weissman et al. [30], slightly reworded, then gives a guarantee of accuracy:

¹We also show that samples are not identically distributed, but demonstrate that that problem would be resolvable.

²The reader might be puzzled by this statement, since certain guarantees on the combination of abstraction and RL are known. This can be explained by the generality of Abstracted RL: in this setting there is a non-stationarity caused by the clustering of states with different dynamics. There is a lot of related work in other abstraction settings (e.g., state aggregation) where this complication does not occur due to the particularities of their setting [24, 11, 19, 20, 23, 10]. In section 4 we give details to back up our claim for individual papers.

³This is a crucial element in being able to guarantee good performance, where performance can be measured in different ways, e.g. in PAC-MDP terms [28] or in terms of regret [14]. We focus on the model quality.

72 **Lemma 1** (L_1 inequality [30]). *Let $Y_{s,a} = Y^{(1)}, Y^{(2)}, \dots, Y^{(m)}$ be i.i.d. random variables*
 73 *distributed according to $T(\cdot|s, a)$. Then, for all $\epsilon > 0$,*

$$\Pr(\|T_Y(\cdot|s, a) - T(\cdot|s, a)\|_1 \geq \epsilon) \leq (2^{|S|} - 2)e^{-\frac{1}{2}m\epsilon^2}. \quad (3)$$

74 In this way, MBRL can upper bound the probability that the learned model, based on m samples, for
 75 a state-action pair (s, a) $T_Y(\cdot|s, a)$ will be far away ($\geq \epsilon$) from the true model $T(\cdot|s, a)$.

76 The situation is more subtle if the MDP is not communicating, i.e., if there exists $s_1, s_2 \in S$ for which
 77 there is no deterministic policy that eventually leads from s_1 to s_2 . This can create a dependence
 78 between the samples [28]. Intuitively this happens because, if we look at one particular state-action
 79 pair (s, a) , there might be a transition to state s' such that the probability to return to s is 0. Thus if
 80 we would have n outcomes of (s, a) we would immediately know that at least $n - 1$ outcomes were
 81 not state s' . Since as soon as we observe s' , we know the agent would not be able to return to state s .
 82 Strehl and Littman [28] show that in this setting it is still possible to use Lemma 3 as an upper bound.

83 2.2 State abstraction for given models

84 In the planning setting, where the model is known a priori, a state abstraction can be formulated as
 85 a grouping or mapping from ground states to abstract states [18]. This is done with an abstraction
 86 function ϕ , a surjective function that maps from ground states, $s \in S$, to abstract states $\bar{s} \in \bar{S}$:
 87 $\phi(s) : S \rightarrow \bar{S}$. Here \bar{S} is defined as $\bar{S} = \{\phi(s) | s \in S\}$. We use the $\bar{\cdot}$ notation to refer to the abstract
 88 space. We slightly overload notation and let \bar{s} both denote an abstract state as well as the set of
 89 ground states that map to the abstract state \bar{s} , i.e., $\bar{s} = \{g \in S | \phi(g) = \bar{s}\}$, if $\bar{s} \in \bar{S}$. The use should
 90 be clear from the context. We define the probability to transition to an abstract state $\Pr(\bar{s}'|s, a)$ as
 91 follows:

$$\Pr(\bar{s}'|s, a) \triangleq \sum_{s' \in \bar{s}'} T(s'|s, a). \quad (4)$$

92 This is a very general form of state abstraction, that clusters together states with different dynamics
 93 into abstract states. Note that we do assume that the given state abstraction deterministically maps
 94 states to an abstract state. This in contrast to some related work on problems with block structure
 95 [10], where a Markov state can lead to multiple observations (abstract states in our terminology) that
 96 need to be aggregated appropriately to result in a small MDP [4, 10].

97 **Approximate model similarity abstraction** Many different abstraction criteria exist [17], here we
 98 focus on approximate model similarity abstraction [1]. In this abstraction two states can map to the
 99 same abstract state if their behavior is similar, i.e., when the reward function and the transitions to
 100 abstract states are close. Approximate model-similarity is defined as follows:

101 **Definition 1.** *An approximate model-similarity abstraction, $\phi_{model, \eta}$, for fixed η , satisfies:*

$$\begin{aligned} \phi_{model, \eta}(s_1) = \phi_{model, \eta}(s_2) &\implies \forall_a |R(s_1, a) - R(s_2, a)| \leq \eta, \\ &\forall_{\bar{s}' \in \bar{S}, a} |\Pr(\bar{s}'|s_1, a) - \Pr(\bar{s}'|s_2, a)| \leq \eta. \end{aligned} \quad (5)$$

102 From now on we will just refer to $\phi_{model, \eta}$ as ϕ .

103 We note that the abstraction we consider, approximate model-similarity abstraction, is still quite
 104 generic. It can cluster together states that have different transition and reward functions. However, in
 105 the online Abstracted RL setting, the differences in dynamics can cause a dependence between the
 106 samples, as we will show in detail in section 3. E.g. looking at (\bar{s}, a) , the probability that we reach a
 107 state s' depends both on the probability that we reach a particular state $s \in \bar{s}$ and then state s' from s .

108 Returning to abstraction of a given model, it is possible to construct an abstract MDP \bar{M}_ω from
 109 the model of an MDP M and an abstraction function ϕ , where ω is an action-specific⁴ weighting
 110 function, defined as follows:

⁴The action-specific weighting function is a more general weighting function than is typically used, e.g. by Li et al. [18], which is not action-specific, i.e., it only depends on the state s . More formally it is the case where $\forall_{a, a' \in A, s \in S} \omega(s, a) = \omega(s, a')$.

111 **Definition 2.** We refer to the weight associated with a ground state, $s \in S$, and action, $a \in A$, by
 112 $\omega(s, a)$. We have: $\forall_{s \in S, a \in A} 0 \leq \omega(s, a) \leq 1$ and $\sum_{s' \in \phi(s)} \omega(s', a) = 1$.

113 The weighting function can be used to create abstract transition and reward functions, which are a
 114 weighted average of the ground function. In this way, from M , ϕ and any ω we can *construct* an
 115 abstract MDP \bar{M}_ω :

116 **Definition 3.** Given an MDP M , ϕ , and ω , $\bar{M}_\omega = \langle \bar{S}, A, \bar{T}_\omega, \bar{R}_\omega, \gamma \rangle$ is constructed as:

$$\bar{S} = \{\phi(s) \mid s \in S\}, A = A, \gamma = \gamma, \quad (6)$$

$$\forall_{\bar{s} \in \bar{S}, a \in A} \bar{R}_\omega(\bar{s}, a) = \sum_{s \in \bar{s}} \omega(s, a) R(s, a), \quad (7)$$

$$\forall_{\bar{s}, \bar{s}' \in \bar{S}, a \in A} \bar{T}_\omega(\bar{s}' | \bar{s}, a) = \sum_{s \in \bar{s}} \sum_{s' \in \bar{s}'} \omega(s, a) T(s' | s, a). \quad (8)$$

117 An abstract MDP \bar{M}_ω is just an MDP. This means we can use any planning method we like to find
 118 an optimal policy $\bar{\pi}^*$ for \bar{M}_ω .

119 What we are interested in is the performance of a policy on the abstract space, when applied on the
 120 original problem M . Any policy on the abstract space $\bar{\pi}$ can be used in M as follows $\bar{\pi}(s) := \bar{\pi}(\phi(s))$,
 121 leading to $V^{\bar{\pi}^*}$. It has been shown that we can upper bound the loss in performance due to using an
 122 optimal policy for \bar{M}_ω , $\bar{\pi}^*$ in M instead of using the optimal solution for M [8, 1, 29]:

123 **Lemma 2** (Lemma 4 [29]). *An approximate model similarity abstraction (Definition 1), has sub-*
 124 *optimality bounded in η : $\forall_{s \in S} V^*(s) - V^{\bar{\pi}^*}(s) \leq \frac{2\eta + 2\gamma(|\bar{S}| - 1)\eta}{(1 - \gamma)^2}$.*

125 3 Abstracted MBRL and the problem of online data collection

126 As explained, we are interested in Abstracted RL, where we have an approximate model similarity
 127 abstraction function ϕ and an MDP M . The agent acts in M but only observes $\phi(s)$ using abstraction
 128 function ϕ , as in Figure 1. This setting can also be seen as a POMDP, where the states are hidden and
 129 there is a deterministic observation function, $o = \phi(s)$. However, in contrast to the usual POMDP
 130 settings, we look for a myopic (memoryless) policy. While we know that in general this can lead to
 131 arbitrarily bad results [26], in this case the value loss would be bounded in the *planning* setting by
 132 Lemma 2. However, now we assume we are in the Abstracted RL setting, and the result for planning
 133 may not hold for the learned model.

134 We assume we know S, A, R, γ , and ϕ (and thus \bar{S}), but do not know the transition function.⁵ Since
 135 we do not know the transition function we can neither simply do planning on M nor can we construct
 136 an abstract MDP, using Definition 3, and solve that. Instead, we let the agent interact with M but
 137 use ϕ to let the agent observe $\phi(s)$, instead of s , and build a learned (abstract) model from the
 138 observations it obtains. We show the general Abstracted MBRL procedure in Algorithm 1.

139 The agent collects data for every abstract state-action pair (\bar{s}, a) , which is stored as sequences $\bar{Y}_{\bar{s}, a}$:

$$\bar{Y}_{\bar{s}, a} : \{\bar{s}'^{(1)}, \bar{s}'^{(2)}, \dots, \bar{s}'^{(m)}\}. \quad (9)$$

140 Similar to before in (1), we construct a learned model \bar{T}_Y , now looking at the abstract next-states that
 141 were reached:

$$\bar{T}_Y(\bar{s}' | \bar{s}, a) \triangleq \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{\bar{Y}_{\bar{s}, a}^{(i)} = \bar{s}'\}. \quad (10)$$

142 If this model would be equal, or close, to the transition function \bar{T}_ω of an abstract MDP \bar{M}_ω , for
 143 some valid ω , we could upper bound the loss in performance due to applying learned policy $\bar{\pi}^*$ to M
 144 instead of the optimal policy π^* [1, 29].

⁵The assumption that the reward function is known simplifies our arguments but can be relaxed.

Algorithm 1 Procedure: Abstracted MBRL

Input: $M, \phi, \delta, \epsilon, \pi$
 $\bar{Y} = \text{COLLECTSAMPLES}(M, \phi, \delta, \epsilon, \pi)$
The sampling results in sequences $\bar{Y}_{\bar{s},a}$, one for every pair (\bar{s}, a) :
 $\bar{Y}_{\bar{s},a} = \phi(s^{(1)}), \dots, \phi(s^{(m)})$
 $= \bar{s}^{(1)}, \dots, \bar{s}^{(m)}$
for all $(\bar{s}, a, \bar{s}') \in \bar{S} \times A \times \bar{S}$ **do**
 $\bar{T}_Y(\bar{s}'|\bar{s}, a) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{\bar{Y}_{\bar{s},a}^{(i)} = \bar{s}'\}$
end for
 $\bar{M}_Y := \langle \bar{S}, A, \bar{T}_Y, \bar{R}, \gamma \rangle$
 $\bar{\pi}_Y^* = \text{Value Iteration}(\bar{M}_Y)$
Apply $\bar{\pi}_Y^*$ to M

Algorithm 2 COLLECTSAMPLES Online

Input: $M, \phi, \delta, \epsilon, \pi$
 $s = \text{initial state}$
// The number of samples m is based on the simulator analysis, Theorem 1.
 $\kappa = \delta / (|\bar{S}| |A|)$
 $m = \lceil \frac{2[\ln(2^{|\bar{S}|} - 2) - \ln(\kappa)]}{\epsilon^2} \rceil$
for all $\bar{s} \in \bar{S}$ **do**
 $\bar{Y}_{\bar{s},a} = []$
end for
while $\min_{(\bar{s},a)} |\bar{Y}_{\bar{s},a}| < m$ **do**
 $\bar{s} = \phi(s)$
 $a = \pi(\bar{s})$
 $s' = \text{Step}(s, a)$
 $\bar{Y}_{\bar{s},a}.\text{append}(\phi(s'))$
 $s = s'$
end while
Return: Return all $\bar{Y}_{\bar{s},a}$

146 Our main question is: do the finite-sample model learning guarantees of MBRL algorithms still hold
147 in the Abstracted RL setting?

148 **3.1 Online data collection**

149 In this section we follow the MBRL method from Algorithm 1, collecting samples online using
150 Algorithm 2.⁶ Starting from an initial state the agent follows a policy π . Instead of observing the
151 states s , the agent observes abstract states $\bar{s} = \phi(s)$, see Figure 1.

152 We make two important assumptions in order to make analysis possible. We assume that the MDP
153 is ergodic [25]⁷ and that the policy assigns a positive probability to every action in every abstract
154 state. Together this can guarantee that Algorithm 2 can obtain any finite number of samples for every
155 abstract state-action pair within finite time.

156 Our question is, can we still use Lemma 1 to guarantee that we learn an accurate model?

157 Since we learn an abstract transition model \bar{T}_Y , we want to be able to guarantee that this learned
158 model will be close to the transition model of some abstract MDP. To define this transition model,
159 we first look at how the data is collected.

160 In the online data collection, a sample in $\bar{Y}_{\bar{s},a}$ is drawn when the agent takes action a when it is
161 in a ground state $s \in \bar{s}$. Specifically the i -th abstract $\bar{Y}_{\bar{s},a}^{(i)} = \bar{s}'$ is drawn from (ground) state
162 $X_{\bar{s},a}^{(i)} = s \in \bar{s}$:

$$\bar{Y}_{\bar{s},a}^{(i)} \sim \Pr(\cdot | X_{\bar{s},a}^{(i)} = s, a). \quad (11)$$

163 Let $X_{\bar{s},a} = (X_{\bar{s},a}^{(i)})_{i=1}^m$ denote the sequence of ground states $s \in \bar{s}$ from which the agent took action
164 a . Each ground state gets a weight according to how often it was sampled from, which we formalize
165 with the weighting function $\omega_X: \forall_{(\bar{s},a), s \in \bar{s}} \omega_X(s, a) \triangleq \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{X_{\bar{s},a}^{(i)} = s\}$. We use ω_X to define
166 \bar{T}_{ω_X} analogous to (8):

$$\forall_{(\bar{s},a), \bar{s}'} \bar{T}_{\omega_X}(\bar{s}'|\bar{s}, a) = \sum_{s \in \bar{s}} \omega_X(s, a) \sum_{s' \in \bar{s}'} T(s'|s, a). \quad (12)$$

⁶The order of m in Algorithm 2, the number of samples we want to collect, is based on the analysis of Model-based Interval Estimation (MBIE) [28].

⁷An ergodic, or recurrent, MDP is an MDP where, under every stationary policy, every state is recurrent (i.e., asymptotically every state will be visited infinitely often) [25].

167 We want to have a concentration inequality to provide bounds on the deviation of the learned model
 168 \bar{T}_Y from \bar{T}_{ω_X} , we refer to this inequality as the abstract L1 inequality, similar in form to (3):

$$P(|\bar{T}_Y(\cdot|\bar{s}, a) - \bar{T}_{\omega_X}(\cdot|\bar{s}, a)|_1 \geq \epsilon) \leq \delta, \quad (13)$$

169 where $\bar{T}_Y(\cdot|\bar{s}, a)$ is defined according to (10) and \bar{T}_{ω_X} according to (12).

170 If we could directly obtain i.i.d. samples from \bar{T}_{ω_X} and base our learned model \bar{T}_Y on the obtained
 171 samples, then we would be able to show that the abstract L1 inequality holds by applying Lemma 1.
 172 Since in this case, we would have m i.i.d. samples per abstract state-action pair, distributed according
 173 to $\bar{T}_{\omega_X}(\cdot|\bar{s}, a)$.

174 However, the samples are not guaranteed to be i.i.d. when the agent follows Algorithm 2 to collect
 175 the samples. Since every sample $\bar{Y}^{(i)}$ was obtained after taking action a from state $X_{\bar{s},a}^{(i)} = s \in \bar{s}$, as
 176 in (11). These can have different distributions if $X_{\bar{s},a}^{(i)} \neq X_{\bar{s},a}^{(j)}$.

177 **Non Identically Distributed** While Lemma 1 assumes i.i.d. random variables, we show that it also
 178 holds when the random variables are independent but not (necessarily) identically distributed.

179 **Lemma 3.** *Let $X_{\bar{s},a} = s_1, \dots, s_m$ be a sequence of states $s \in \bar{s}$ and let*
 180 *$\bar{Y}_{\bar{s},a} = \bar{Y}^{(1)}, \bar{Y}^{(2)}, \dots, \bar{Y}^{(m)}$ be independent random variables distributed according to*
 181 *$\Pr(\cdot|s_1, a), \dots, \Pr(\cdot|s_m, a)$ (Eqn. 4). Then, for all $\epsilon > 0$,*

$$\Pr(|\bar{T}_Y(\cdot|\bar{s}, a) - \bar{T}_{\omega_X}(\cdot|\bar{s}, a)|_1 \geq \epsilon) \leq (2^{|\bar{S}|} - 2)e^{-\frac{1}{2}m\epsilon^2}. \quad (14)$$

182 The proof can be found in Appendix B. It mostly follows the proof by Weissman et al. [30], which uses
 183 Hoeffding's inequality [12] and the union bound [6].⁸ Lemma 3 shows that the fact that Hoeffding's
 184 inequality does not need identically distributed data can be carried over to the setting from Lemma 1.

185 **Independence** We may be tempted to assume the samples are independent, i.e.,

$$\forall_{\bar{s}_1, \dots, \bar{s}_m \in (\bar{S})^m} \Pr(\bar{Y}_{\bar{s},a}^{(1)} = \bar{s}_1, \dots, \bar{Y}_{\bar{s},a}^{(m)} = \bar{s}_m) = \Pr(\bar{Y}_{\bar{s},a}^{(1)} = \bar{s}_1) \cdots P(\bar{Y}_{\bar{s},a}^{(m)} = \bar{s}_m) \quad (15)$$

186 however, this may not be the case:

187 **Observation 1.** *When collecting samples online, i.e., based on Algorithm 2, the samples cannot be*
 188 *assumed to be independent.*

189 The following counterexample illustrates this.

190 **Counterexample** To show that the samples may not be independent,
 191 we will give a counterexample. We use the example MDP and
 192 abstraction in Figure 2, where we have 4 (ground) states, 3 abstract
 193 states and only 1 action. We look at the transition probability from
 194 abstract state A, $\bar{T}_Y(\cdot|A)$.

195 We will consider two samples and show that for at least one combi-
 196 nation of \bar{s}_1 and \bar{s}_2 the samples are not independent. Consider
 197 $\bar{s}_1 = \bar{s}_2 = B$. That is, the first two times that we experience a
 198 transition from the abstract state A, we end up in B.

199 Let state 1 be the starting state. Then we have $\Pr(\bar{Y}_A^{(1)} = B) =$
 200 $\Pr(B|1) = 0.6$ and

$$\Pr(\bar{Y}_A^{(2)} = B) = \sum_{\bar{s} \in \bar{S}} \Pr(\bar{Y}_A^{(2)} = B | \bar{Y}_A^{(1)} = \bar{s}) \Pr(\bar{Y}_A^{(1)} = \bar{s}) \quad (16)$$

$$= 0 + 0.6 \cdot 0.6 + 0.4 \cdot 0.4 = 0.52. \quad (17)$$

201 So then we end up with: $\Pr(\bar{Y}_A^{(1)} = B) \Pr(\bar{Y}_A^{(2)} = B) = 0.6 \cdot 0.52 = 0.321$. And for the joint
 202 probability: $\Pr(\bar{Y}_A^{(1)} = B, \bar{Y}_A^{(2)} = B) = \Pr(\bar{Y}_A^{(1)} = B) \Pr(\bar{Y}_A^{(2)} = B | \bar{Y}_A^{(1)} = B) = 0.6 \cdot 0.6 =$
 203 0.36 .

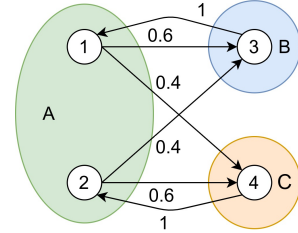


Figure 2: Simple MDP, with only 1 action, and abstraction. The small circles are ground states (1,2,3,4). A, B and C are the abstract states. The numbers along the arrows show the transition probabilities, e.g. $P(3|1) = 0.6$.

⁸Hoeffding's inequality and the union bound can be found in Appendix A.

204 Thus we have that $\Pr(\bar{Y}_A^{(1)} = B, \bar{Y}_A^{(2)} = B) \neq \Pr(\bar{Y}_A^{(1)} = B) \Pr(\bar{Y}_A^{(2)} = B)$, the samples are not
 205 independent. Leading us to the second observation:

206 **Observation 2.** *As independence cannot be guaranteed, Lemmas 1 and 3 cannot be readily applied*
 207 *to show that the abstract L1 inequality holds.*

208 3.2 Simulator data collection

209 Here we also want to give a guarantee in the form of the abstract L1 inequality from (13). While in
 210 the previous section we found this was not possible because the samples were dependent, here we
 211 assume that we have access to a simulator. To some extent this is not surprising, but to the best of our
 212 knowledge, this is the first work that explicitly shows how to combine MBRL and abstraction, using
 213 a simulator. We assume that this allows us to select (or move to) any state and draw a sample from its
 214 transition function. This we call the independent samples assumption:

215 **Assumption 1** (Independent samples). *We assume we can obtain independent samples, e.g. for any*
 216 *state-action pair (s, a) we can draw samples directly from its transition function $T(\cdot|s, a)$.*

217 In case a simulator of the MDP is available this is a reasonable assumption. For every (\bar{s}, a) the
 218 simulator sampling procedure (Algorithm 3 in Appendix B) selects a prototype $x_{\bar{s}, a} \in \bar{s}$ to sample
 219 from. We define a weighting function $\omega_x(s, a)$ that has weight 1 if s is the prototype $x_{\bar{s}, a}$ and 0
 220 otherwise:

$$\forall_{(\bar{s}, a), s \in \bar{s}} \omega_x(s, a) \triangleq \mathbb{1}\{s = x_{\bar{s}, a}\}. \quad (18)$$

221 Then we use this ω_x to define the abstract transition function \bar{T}_{ω_x} according to (8). $\bar{T}_{\omega_x}(\bar{s}'|\bar{s}, a) =$
 222 $\sum_{s' \in \bar{s}'} T(s'|s = x_{\bar{s}, a}, a)$. This way the samples that we collect for one pair (\bar{s}, a) are i.i.d., they are
 223 independent because of our assumption of independent samples and identically distributed because
 224 we sample from the prototype. This means we can use Lemma 1. We show that with the simulator
 225 we can combine MBRL and abstraction, and still learn an accurate model, that is, we can guarantee
 226 that \bar{T}_Y will be close to \bar{T}_{ω_x} , with high probability:

227 **Theorem 1.** *Under assumption 1, and following the procedure in Algorithm 1, with the data*
 228 *collection from Algorithm 3 (Appendix B), with inputs $|\bar{S}|, A, \epsilon$ and δ . For \bar{T}_Y constructed by the*
 229 *algorithm we have that with probability $1 - \delta$, the following holds:*

$$\forall_{(\bar{s}, a)} \|\bar{T}_Y(\cdot|\bar{s}, a) - \bar{T}_{\omega_x}(\cdot|\bar{s}, a)\|_1 \leq \epsilon. \quad (19)$$

230 By Assumption 1 we can obtain any number of independent samples for each abstract state action
 231 pair (\bar{s}, a) . Using Lemma 1 we can then derive the number of samples m that is required for each
 232 pair (\bar{s}, a) such that, after applying a union bound, we obtain the bounds in (19). The full proof can
 233 be found in Appendix B.

234 4 Related work

235 There is a lot of work that considers the combination of abstraction with either planning or (online)
 236 RL. In a lot of these works the dependence of samples that arises in Abstracted RL is not an issue
 237 due to various assumptions, similarly to how in MBRL dependence of samples is often not an issue
 238 because of the Markov property and the assumption that the MDP is communicating [25]. Often
 239 this is either due the assumption that data has been obtained i.i.d., the specific type of abstraction, or
 240 because access to an MDP model is assumed.

241 One paper that does give a result for the Abstracted RL setting is the work by Abel et al. [2].
 242 They show that in this setting R-MAX [7] no longer maintains its guarantees when paired with any
 243 type of state-abstraction function, though their example is specifically for approximate Q-function
 244 abstractions. They also show that the expected trajectory of a learning agent in a constructed
 245 abstract MDP (Definition 3) is not the same as in Abstracted RL. Their work makes clear there is a
 246 complication when combining MBRL and abstraction, here we further investigated the cause of this
 247 complication, the dependence between samples.

248 For planning in constructed abstract MDPs, some main results for exact state-abstractions come
 249 from Li et al. [18] and for approximate state-abstractions from Abel et al. [1]. The results from Abel
 250 et al. [1] allow for quantifying an upper bound on performance for policies found in a constructed

251 abstract MDP, as in section 2.2. Taïga et al. [29] build on this by giving a result for performing
252 RL on top of the constructed abstract MDP. They provide upper bounds for this setting when using
253 MBIE with exploratory bonus (MBIE-EB) [28]. In addition, they give an example to show that in
254 this combination you cannot guarantee optimal performance in the original MDP. Still, they show
255 that an upper bound on the loss in value can be given.

256 Both Paduraru et al. [24] and Jiang et al. [15] deal with the issue of dependence by making the
257 explicit assumption that samples are obtained i.i.d. Paduraru et al. [24] consider the setting where
258 we are given a dataset for a continuous domain and then use discretization to aggregate states into
259 abstract states. They then give PAC-style guarantees on the learned abstract model and the value that
260 a policy based on this model can achieve in the real MDP. Instead of using the L1 deviation bound
261 from Weissman et al. [30], Paduraru et al. [24] use a similar bound for i.i.d. samples by Devroye and
262 Györfi [9], which requires a minimum amount of samples. Another difference is that their results
263 calculate the probability that the model will be ϵ -accurate given a fixed dataset. They assume that the
264 data has been gathered i.i.d., but our Lemma 3 shows that merely independent data would be enough.
265 At the same time, our results show that when we collect data online in the Abstracted RL setting,
266 their guarantees will not hold.

267 Jiang et al. [15] operate in the abstraction selection setting, where the agent is provided with a set of
268 abstraction functions (state representations). They do not assume that any of the abstraction functions
269 results in a Markov model, but they do assume a given dataset, with data that was collected i.i.d. They
270 give a bound directly on how accurate the Q-values based on the (implicitly) learned model will be,
271 rather than on the accuracy of the model itself. As we showed, the assumption that the data is i.i.d.
272 is not a trivial assumption, since it means the data cannot just have been collected online. They do
273 mention that samples will not be strictly independent if a fixed exploration policy is used to collect
274 data but do not mention what the implications are.

275 There are quite a few other papers in the abstraction selection setting, several of these assume that
276 the given set of state representations contains a Markov model [11, 19, 23]. Hallak et al. [11] give
277 asymptotic guarantees for selecting the correct model and on building an exact MDP model. The
278 assumption that there is an MDP model in the given set of representations is crucial in their analysis
279 since for this ‘true model’ the samples are i.i.d. Similarly, both Maillard et al. [19] and Ortner
280 et al. [23] also assume that the given set of state representations contains a Markov model. They
281 create an algorithm for which they obtain regret bounds, their analysis also makes use of the Markov
282 representation.

283 Other work in the abstraction selection setting does not assume that the set of abstraction functions
284 contains a Markov model [16, 22]. However, Ortner et al. [22] use Theorem 2.1 from Weissman
285 et al. [30] that requires i.i.d. samples, which we have shown here cannot be guaranteed in this setting.
286 Lattimore et al. [16] operate in a setting more general than MDPs, where the dynamics of the true
287 environment depend arbitrarily on a history of actions, rewards, and observations. The agent gets
288 as input a finite set of environments, one of which is the true environment. Since the input includes
289 the full model of each environment, the agent does not have to learn a transition model. Instead, to
290 obtain regret bounds, they directly compare the rewards the agent obtains to the expected rewards of
291 the given environments and eliminate environments that are implausible given the observed rewards.

292 Another way to deal with the issue of dependence is by looking at convergence in the limit [27, 13, 20].
293 Singh et al. [27] give an asymptotic result for the convergence of Q-learning and TD(0) in MDPs
294 with soft state aggregation. Soft state aggregation means that a state s belongs to a cluster x with
295 some probability $P(x|s)$, this means a state s can belong to several clusters. The state-abstraction
296 functions we consider are a special case of this, where each state is part of exactly one abstract state
297 (or cluster). Their result relies on having a stationary policy that assigns a non-zero probability to
298 every action in every state and the assumption that the MDP is ergodic. Together these imply there is
299 a limiting state distribution, and using this they show convergence asymptotically. Our main interest
300 is in finite-samples guarantees with policies that change due to exploration, whereas this work gives
301 convergence guarantees in the limit using a fixed policy.

302 Hutter [13] gives a variety of results focusing on both approximate and exact abstractions in envi-
303 ronments without MDP assumptions. Several of these are in the planning setting, similar to those of
304 Abel et al. [1]. Most relevant for us is their Theorem 12, which for online RL shows convergence in
305 the limit of the empirical transition function under weak conditions, e.g. if the abstract process itself

306 is an MDP. Under this condition however the problem reduces to RL in an (abstract) MDP, rather
307 than Abstracted RL.

308 Majeed and Hutter [20] build on the work by Hutter [13] and focus on the combination of model-free
309 RL and exact abstraction. They show that, under the condition of state uniformity, q-learning can be
310 shown to converge in the limit to the optimal solution. State uniformity means that histories that are
311 grouped together have the same optimal q-values. In contrast to our setting, they look at an exact
312 abstraction, extending it to approximate aggregation was left as an open question.

313 Other related work is in the area of MDPs with rich observations or block structure [4, 10]. However,
314 in that setting each observation can be generated only from a *single* hidden state, which means that
315 the issue of non-i.i.d. data due to abstraction does not arise. In contrast, in our setting multiple
316 (hidden) states generate the same observation. Azzadenesheli et al. [4] state their setting can be
317 seen as an aggregation problem, where the observations can be aggregated to form a small (latent)
318 MDP. But in our case, we do not try to learn the MDP (as it is not small). Du et al. [10] describe
319 that their setting is similar to exact model similarity (or bisimulation), but we focus on approximate
320 model similarity which is what introduces the problems as described here.

321 5 Discussion

322 When collecting samples online in Abstracted RL, there is a potential dependence between samples,
323 meaning we cannot use the typically used concentration results that assume i.i.d. samples, e.g.
324 Theorem 2.1 from Weissman et al. [30], the empirical Bernstein inequality [3, 21] or the Chernoff
325 bound. In case the samples are only weakly dependent, it may be that concentration inequalities
326 for (weakly) dependent variables are a viable alternative through which we can come to guarantees
327 on the learned model. Alternatively, it may be possible to change the sampling process to ensure
328 independent samples. One way to ensure independent samples is to, as in the simulator setting, select
329 a prototype state and only use the samples collected from this state. Though in this case, we will be
330 discarding information when we reach a state $s \in \bar{s}$ that is not the prototype.

331 Our assumption on the simulator that we can go/reset to any state to draw samples from it can be
332 relaxed, though it may mean that the procedure takes considerably more time. Consider the case
333 where we cannot just reset the simulator to the state s from which we want to sample, and instead, it
334 would behave like the MDP. In this case, we would have to take the right actions to arrive at the state
335 s from which we would like to sample. Since we assume we do not know T , this may take a long
336 time. This also shows the difficulty of assuming that in the MBRL setting somehow have access to an
337 i.i.d. dataset, as has been assumed in some earlier work [24, 15].

338 6 Conclusion

339 We analyzed Abstracted RL: the combination of MBRL and state abstraction when the model of
340 the MDP is not available. We have shown that in Abstracted RL samples obtained online cannot
341 be assumed to be independent. Since many current guarantees from MBRL methods rely on this
342 assumption, their guarantees do not hold in this setting. And in fact, no current methods exist that
343 give (correct) finite-sample quality guarantees for the models learned in this setting. This also means
344 that current results that rely on an i.i.d. assumption cannot be readily transferred to the Abstracted
345 RL setting.

346 In addition, we show that with a simulator, since we can draw independent samples, it is still possible
347 to give guarantees on the accuracy of the model. However, having access to a simulator may often
348 not be possible. An important step is to see if the MBRL guarantees can be adapted to Abstracted RL
349 for online sample collection.

350 References

351 [1] David Abel, David Hershkowitz, and Michael Littman. Near optimal behavior via approximate
352 state abstraction. In *International Conference on Machine Learning*, pages 2915–2923, 2016.

- 353 [2] David Abel, Dilip Arumugam, Lucas Lehnert, and Michael Littman. State abstractions for
354 lifelong reinforcement learning. In *International Conference on Machine Learning*, pages
355 10–19, 2018.
- 356 [3] Jean-Yves Audibert, Rémi Munos, and Csaba Szepesvári. Tuning bandit algorithms in stochastic
357 environments. In *International conference on algorithmic learning theory*, pages 150–165.
358 Springer, 2007.
- 359 [4] Kamyar Azizzadenesheli, Alessandro Lazaric, and Animashree Anandkumar. Reinforcement
360 learning in rich-observation mdps using spectral methods. *arXiv preprint arXiv:1611.03907*,
361 2016.
- 362 [5] Aijun Bai, Siddharth Srivastava, and Stuart J Russell. Markovian state and action abstractions
363 for mdps via hierarchical mcts. In *IJCAI*, pages 3029–3039, 2016.
- 364 [6] George Boole. *An investigation of the laws of thought: on which are founded the mathematical*
365 *theories of logic and probabilities*. Dover Publications, 1854.
- 366 [7] Ronen I Brafman and Moshe Tennenholtz. R-max-a general polynomial time algorithm for
367 near-optimal reinforcement learning. *Journal of Machine Learning Research*, 3(Oct):213–231,
368 2002.
- 369 [8] Richard Dearden and Craig Boutilier. Abstraction and approximate decision-theoretic planning.
370 *Artificial Intelligence*, 89(1-2):219–283, 1997.
- 371 [9] L. Devroye and L. Györfi. *Nonparametric Density Estimation: The LI View*. Wiley Interscience
372 Series in Discrete Mathematics. Wiley, 1985.
- 373 [10] Simon Du, Akshay Krishnamurthy, Nan Jiang, Alekh Agarwal, Miroslav Dudik, and John
374 Langford. Provably efficient rl with rich observations via latent state decoding. In *International*
375 *Conference on Machine Learning*, pages 1665–1674. PMLR, 2019.
- 376 [11] Assaf Hallak, Dotan Di-Castro, and Shie Mannor. Model selection in markovian processes. In
377 *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and*
378 *data mining*, pages 374–382, 2013.
- 379 [12] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of*
380 *the American Statistical Association*, 58(301):13–30, 1963.
- 381 [13] Marcus Hutter. Extreme state aggregation beyond markov decision processes. *Theoretical*
382 *Computer Science*, 650:73–91, 2016.
- 383 [14] Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement
384 learning. *Journal of Machine Learning Research*, 11(Apr):1563–1600, 2010.
- 385 [15] Nan Jiang, Alex Kulesza, and Satinder Singh. Abstraction selection in model-based reinforce-
386 ment learning. In *International Conference on Machine Learning*, pages 179–188, 2015.
- 387 [16] Tor Lattimore, Marcus Hutter, Peter Sunehag, et al. The sample-complexity of general rein-
388 forcement learning. In *Proceedings of the 30th International Conference on Machine Learning*.
389 Journal of Machine Learning Research, 2013.
- 390 [17] Lihong Li. *A unifying framework for computational reinforcement learning theory*. PhD thesis,
391 Rutgers University-Graduate School-New Brunswick, 2009.
- 392 [18] Lihong Li, Thomas J Walsh, and Michael L Littman. Towards a unified theory of state abstraction
393 for mdps. In *ISAIM*, 2006.
- 394 [19] Odalric-Ambrym Maillard, Phuong Nguyen, Ronald Ortner, and Daniil Ryabko. Optimal
395 regret bounds for selecting the state representation in reinforcement learning. In *International*
396 *Conference on Machine Learning*, pages 543–551. PMLR, 2013.
- 397 [20] Sultan Javed Majeed and Marcus Hutter. On q-learning convergence for non-markov decision
398 processes. In *IJCAI*, pages 2546–2552, 2018.

- 399 [21] Andreas Maurer and Massimiliano Pontil. Empirical Bernstein bounds and sample variance
400 penalization. *arXiv preprint arXiv:0907.3740*, 2009.
- 401 [22] Ronald Ortner, Odalric-Ambrym Maillard, and Daniil Ryabko. Selecting near-optimal ap-
402 proximate state representations in reinforcement learning. In *International Conference on*
403 *Algorithmic Learning Theory*, pages 140–154. Springer, 2014.
- 404 [23] Ronald Ortner, Matteo Pirodda, Alessandro Lazaric, Ronan Fruit, and Odalric-Ambrym Maillard.
405 Regret bounds for learning state representations in reinforcement learning. In *Advances in*
406 *Neural Information Processing Systems*, pages 12738–12748, 2019.
- 407 [24] Cosmin Paduraru, Robert Kaplow, Doina Precup, and Joelle Pineau. Model-based reinforcement
408 learning with state aggregation. In *8th European Workshop on Reinforcement Learning*, 2008.
- 409 [25] Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*.
410 John Wiley & Sons, 2014.
- 411 [26] Satinder P Singh, Tommi Jaakkola, and Michael I Jordan. Learning without state-estimation in
412 partially observable Markovian decision processes. In *Machine Learning Proceedings 1994*,
413 pages 284–292. Elsevier, 1994.
- 414 [27] Satinder P Singh, Tommi Jaakkola, and Michael I Jordan. Reinforcement learning with soft
415 state aggregation. In *Advances in neural information processing systems*, pages 361–368, 1995.
- 416 [28] Alexander L Strehl and Michael L Littman. An analysis of model-based interval estimation for
417 Markov decision processes. *Journal of Computer and System Sciences*, 74(8):1309–1331, 2008.
- 418 [29] Adrien Ali Taïga, Aaron Courville, and Marc G Bellemare. Approximate exploration through
419 state abstraction. *arXiv preprint arXiv:1808.09819*, 2018.
- 420 [30] Tsachy Weissman, Erik Ordentlich, Gadiel Seroussi, Sergio Verdu, and Marcelo J Weinberger.
421 Inequalities for the ℓ_1 deviation of the empirical distribution. *Hewlett-Packard Labs, Tech. Rep.*,
422 2003.

423 Checklist

- 424 1. For all authors...
- 425 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
426 contributions and scope? [Yes]
- 427 (b) Did you describe the limitations of your work? [Yes] In Section 6 we describe one the
428 limitations of our work.
- 429 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 430 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
431 them? [Yes]
- 432 2. If you are including theoretical results...
- 433 (a) Did you state the full set of assumptions of all theoretical results? [Yes] We state our
434 general assumptions on the environment in Section 2 and more specific assumptions in
435 Section 3, Section 3.1 and Section 3.2.
- 436 (b) Did you include complete proofs of all theoretical results? [Yes] In the Appendix.
- 437 3. If you ran experiments...
- 438 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
439 mental results (either in the supplemental material or as a URL)? [N/A]
- 440 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
441 were chosen)? [N/A]
- 442 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
443 ments multiple times)? [N/A]
- 444 (d) Did you include the total amount of compute and the type of resources used (e.g., type
445 of GPUs, internal cluster, or cloud provider)? [N/A]

- 446 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 447 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 448 (b) Did you mention the license of the assets? [N/A]
- 449 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 450
- 451 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 452 using/curating? [N/A]
- 453 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 454 information or offensive content? [N/A]
- 455 5. If you used crowdsourcing or conducted research with human subjects...
- 456 (a) Did you include the full text of instructions given to participants and screenshots, if
- 457 applicable? [N/A]
- 458 (b) Did you describe any potential participant risks, with links to Institutional Review
- 459 Board (IRB) approvals, if applicable? [N/A]
- 460 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 461 spent on participant compensation? [N/A]