

---

# Optimistic planning in Markov decision processes using a generative model

---

**Bal6zsz Sz6r6nyi**  
INRIA Lille - Nord Europe,  
SequeL project, France /  
MTA-SZTE Research Group on  
Artificial Intelligence, Hungary  
`balazs.szorenyi@inria.fr`

**Gunnar Kedenburg**  
INRIA Lille - Nord Europe,  
SequeL project, France  
`gunnar.kedenburg@inria.fr`

**Remi Munos\***  
INRIA Lille - Nord Europe,  
SequeL project, France  
`remi.munos@inria.fr`

## Abstract

We consider the problem of online planning in a Markov decision process with discounted rewards for any given initial state. We consider the PAC sample complexity problem of computing, with probability  $1 - \delta$ , an  $\epsilon$ -optimal action using the smallest possible number of calls to the generative model (which provides reward and next-state samples). We design an algorithm, called StOP (for Stochastic-Optimistic Planning), based on the “optimism in the face of uncertainty” principle. StOP can be used in the general setting, requires only a generative model, and enjoys a complexity bound that only depends on the local structure of the MDP.

## 1 Introduction

### 1.1 Problem formulation

In a *Markov decision process* (MDP), an agent navigates in a state space  $X$  by making decisions from some action set  $U$ . The dynamics of the system are determined by transition probabilities  $P : X \times U \times X \rightarrow [0, 1]$  and reward probabilities  $R : X \times U \times [0, 1] \rightarrow [0, 1]$ , as follows: when the agent chooses action  $u$  in state  $x$ , then, with probability  $R(x, u, r)$ , it receives reward  $r$ , and with probability  $P(x, u, x')$  it makes a transition to a next state  $x'$ . This happens independently of all previous actions, states and rewards—that is, the system possesses the *Markov property*. See [20, 2] for a general introduction to MDPs. We do not assume that the transition or reward probabilities are fully known. Instead, we assume access to the MDP via a *generative model* (e.g. simulation software), which, for a state-action  $(x, u)$ , returns a reward sample  $r \sim R(x, u, \cdot)$  and a next-state sample  $x' \sim P(x, u, \cdot)$ . We also assume the number of possible next-states to be bounded by  $N \in \mathbb{N}$ .

We would like to find an agent that implements a policy which maximizes the expected cumulative discounted reward  $\mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_t]$ , which we will also refer to as the *return*. Here,  $r_t$  is the reward received at time  $t$  and  $\gamma \in (0, 1)$  is the *discount factor*. Further, we take an *online planning* approach, where at each time step, the agent uses the generative model to perform a simulated search (planning) in the set of policies, starting from the current state. As a result of this search, the agent takes a single action. An expensive global search for the optimal policy in the whole MDP is avoided.

---

\*Current affiliation: Google DeepMind

To quantify the performance of our algorithm, we consider a PAC (Probably Approximately Correct) setting, where, given  $\epsilon > 0$  and  $\delta \in (0, 1)$ , our algorithm returns, with probability  $1 - \delta$ , an  $\epsilon$ -optimal action (i.e. such that the loss of performing this action and then following an optimal policy instead of following an optimal policy from the beginning is at most  $\epsilon$ ). The number of calls to the generative model required by the planning algorithm is referred to as its *sample complexity*. The sample and computational complexities of the planning algorithm introduced here depend on local properties of the MDP, such as the quantity of near-optimal policies starting from the initial state, rather than global features like the MDP’s size.

## 1.2 Related work

The online planning approach and, in particular, its ability to get rid of the dependency on the global features of the MDP in the complexity bounds (mentioned above, and detailed further below) is the driving force behind the Monte Carlo Tree Search algorithms [16, 8, 11, 18].<sup>1</sup> The theoretical analysis of this approach is still far from complete. Some of the earlier algorithms use strong assumptions, others are applicable only in restricted cases, or don’t adapt to the complexity of the problem. In this paper we build on ideas used in previous works, and aim at fixing these issues.

A first related work is the sparse sampling algorithm of [14]. It builds a uniform look-ahead tree of a given depth (which depends on the precision  $\epsilon$ ), using for each transition a finite number of samples obtained from a generative model. An estimate of the value function is then built using empirical averaging instead of expectations in the dynamic programming back-up scheme. This results in an

algorithm with (problem-independent) sample complexity of order  $\left(\frac{1}{(1-\gamma)^3\epsilon}\right)^{\frac{\log K + \log[1/(\epsilon(1-\gamma)^2)]}{\log(1/\gamma)}}$  (neglecting some poly-logarithmic dependence), where  $K$  is the number of actions. In terms of  $\epsilon$ , this bound scales as  $\exp(O([\log(1/\epsilon)]^2))$ , which is non-polynomial in  $1/\epsilon$ .<sup>2</sup> Another disadvantage of the algorithm is that the expansion of the look-ahead tree is uniform; it does not adapt to the MDP.

An algorithm which addresses this appears in [21]. It avoids evaluating some unnecessary branches of the look-ahead tree of the sparse sampling algorithm. However, the provided sample bound does not improve on the one in [14], and it is possible to show that the bound is tight (for both algorithms). In fact, the sample complexity turns out to be super-polynomial even in the pure Monte Carlo setting (i.e., when  $K = 1$ ):  $1/\epsilon^{2+(\log C)/\log(1/\gamma)}$ , with  $C \geq \frac{1}{\epsilon^2(1-\gamma)^4}$ .

Close to our contribution are the planning algorithms [13, 3, 5, 15] (see also the survey [18]) that follow the so-called “optimism in the face of uncertainty” principle for online planning. This principle has been extensively investigated in the multi-armed bandit literature (see e.g. [17, 1, 4]). In the planning problem, this approach translates to prioritizing the most promising part of the policy space during exploration. In [13, 3, 5], the sample complexity depends on a measure of the quantity of near-optimal policies, which gives a better understanding of the real hardness of the problem than the uniform bound in [14].

The case of deterministic dynamics and rewards is considered in [13]. The proposed algorithm has sample complexity of order  $(1/\epsilon)^{\frac{\log \kappa}{\log(1/\gamma)}}$ , where  $\kappa \in [1, K]$  measures (as a branching factor) the quantity of nodes of the planning tree that belong to near-optimal policies. If all policies are very good, many nodes need to be explored in order to distinguish the optimal policies from the rest, and therefore,  $\kappa$  is close to the number of actions  $K$ , resulting in the minimax bound of  $(1/\epsilon)^{\frac{\log K}{\log(1/\gamma)}}$ . Now if there is structure in the rewards (e.g. when sub-optimal policies can be eliminated by observing the first rewards along the sequence), then the proportion of near-optimal policies is low, so  $\kappa$  can be small and the bound is much better. In [3], the case of stochastic rewards have been considered. However, in that work the performance is not compared to the optimal (closed-loop) policy, but to the best open-loop policy (i.e. which does not depends on the state but only on the sequence of actions). In that situation, the sample complexity is of order  $(1/\epsilon)^{\max(2, \frac{\log(\kappa)}{\log(1/\gamma)})}$ .

The deterministic and open-loop settings are relatively simple, since any policy can be identified with a sequence of actions. In the general MDP case however, a policy corresponds to an exponentially

<sup>1</sup>A similar planning approach has been considered in the control literature, such as the model-predictive control [6] or in the AI community, such as the  $A^*$  heuristic search [19] and the  $AO^*$  variant [12].

<sup>2</sup>A problem-independent lower bound for the sample complexity, of order  $(1/\epsilon)^{1/\log(1/\gamma)}$ , is provided too.

wide tree, where several branches need to be explored. The closest work to ours in this respect is [5]. However, it makes the (strong) assumption that a full model of the rewards and transitions is available. The sample complexity achieved is again  $(1/\epsilon)^{\frac{\log(\kappa)}{\log(1/\gamma)}}$ , but where  $\kappa \in (1, NK]$  is defined as the branching factor of the set of nodes that simultaneously (1) belong to near-optimal policies, and (2) whose “contribution” to the value function at the initial state is non-negligible.

### 1.3 The main results of the paper

Our main contribution is a planning algorithm, called StOP (for Stochastic Optimistic Planning) that achieves a polynomial sample complexity in terms of  $\epsilon$  (which can be regarded as the leading parameter in this problem), and which is, in terms of this complexity, competitive to other algorithms that can exploit more specifics of their respective domains. It benefits from possible reward or transition probability structures, and does not require any special restriction or knowledge about the MDP besides having access to a generative model. The sample complexity bound is more involved than in previous works, but can be upper-bounded by:

$$(1/\epsilon)^{2 + \frac{\log \kappa}{\log(1/\gamma)} + o(1)} \quad (1)$$

The important quantity  $\kappa \in [1, KN]$  plays the role of a branching factor of the set of important states  $\mathcal{S}^{\epsilon,*}$  (defined precisely later) that “contribute” in a significant way to near-optimal policies. These states have a non-negligible probability to be reached when following some near-optimal policy. This measure is similar (but with some differences illustrated below) to the  $\kappa$  introduced in the analysis of OP-MDP in [5]. Comparing the two, (1) contains an additional constant of 2 in the exponent. This is a consequence of the fact that the rewards are random and that we do not have access to the true probabilities, only to a generative model generating transition and reward samples.

In order to provide intuition about the bound, let us consider several specific cases (the derivation of these bounds can be found in Section E):

- **Worst-case.** When there is no structure at all, then  $\mathcal{S}^{\epsilon,*}$  may potentially be the set of all possible reachable nodes (up to some depth which depends on  $\epsilon$ ), and its branching factor is  $\kappa = KN$ . The sample complexity is thus of order (neglecting logarithmic factors)  $(1/\epsilon)^{2 + \frac{\log(KN)}{\log(1/\gamma)}}$ . This is the same complexity that uniform planning algorithm would achieve. Indeed, uniform planning would build a tree of depth  $h$  with branching factor  $KN$  where from each state-action one would generate  $m$  rewards and next-state samples. Then, dynamic programming would be used with the empirical Bellman operator built from the samples. Using Chernoff-Hoeffding bound, the estimation error is of the order (neglecting logarithms and  $(1-\gamma)$  dependence) of  $1/\sqrt{m}$ . So for a desired error  $\epsilon$  we need to choose  $h$  of order  $\log(1/\epsilon)/\log(1/\gamma)$ , and  $m$  of order  $1/\epsilon^2$  leading to a sample complexity of order  $m(KN)^h = (1/\epsilon)^{2 + \frac{\log(KN)}{\log(1/\gamma)}}$ . (See also [15]) Note that in the worst-case sense there is no uniformly better strategy than a uniform planning, which is achieved by StOP. However, StOP can also do much better in specific settings, as illustrated next.
- **Case with  $K_0 > 1$  actions at the initial state,  $K_1 = 1$  actions for all other states, and arbitrary transition probabilities.** Now each branch corresponds to a single policy. In that case one has  $\kappa = 1$  (even though  $N > 1$ ) and the sample complexity of StOP is of order  $\tilde{O}(\log(1/\delta)/\epsilon^2)$  with high probability<sup>3</sup>. This is the same rate as a Monte-Carlo evaluation strategy would achieve, by sampling  $O(\log(1/\delta)/\epsilon^2)$  random trajectories of length  $\log(1/\epsilon)/\log(1/\gamma)$ . Notice that this result is surprisingly different from OP-MDP which has a complexity of order  $(1/\epsilon)^{\frac{\log N}{\log(1/\gamma)}}$  (in the case when  $\kappa = N$ , i.e., when all transitions are uniform). Indeed, in the case of uniform transition probabilities, OP-MDP would sample the nodes in breadth-first search way, thus achieving this minimax-optimal complexity. This does not contradict the  $\tilde{O}(\log(1/\delta)/\epsilon^2)$  bound for StOP (and Monte-Carlo) since this bound applies to an individual problem and holds in high probability, whereas the bound for OP-MDP is deterministic and holds uniformly over all problems of this type.

<sup>3</sup>We emphasize the dependence on  $\delta$  here since we want to compare this high-probability bound to the deterministic bound of OP-MDP.

Here we see the potential benefit of using StOP instead of OP-MDP, even though StOP only uses a generative model of the MDP whereas OP-MDP requires a full model.

- **Highly structured policies.** This situation holds when there is a substantial gap between near optimal policies and other sub-optimal policies. For example if along an optimal policy, all immediate rewards are 1, whereas as soon as one deviates from it, all rewards are  $< 1$ . Then only a small proportion of the nodes (the ones that contribute to near-optimal policies) will be expanded by the algorithm. In such cases,  $\kappa$  is very close to 1 and in the limit, we recover the previous case when  $K = 1$  and the sample complexity is  $O(1/\epsilon)^2$ .
- **Deterministic MDPs.** Here  $N = 1$  and we have that  $\kappa \in [1, K]$ . When there is structure in the rewards (like in the previous case), then  $\kappa = 1$  and we obtain a rate  $\tilde{O}(1/\epsilon^2)$ . Now when the MDP is almost deterministic, in the sense that  $N > 1$  but from any state-action, there is one next-state probability which is close to 1, then we have almost the same complexity as in the deterministic case (since the nodes that have a small probability to be reached will not contribute to the set of important nodes  $S^{\epsilon,*}$ , which characterizes  $\kappa$ ).
- **Multi-armed bandit** we essentially recover the result of the Action Elimination algorithm [9] for the PAC setting.

Thus we see that in the worst case StOP is minimax-optimal, and in addition, StOP is able to benefit from situations when there is some structure either in the rewards or in the transition probabilities. We stress that StOP achieves the above mentioned results *having no knowledge about  $\kappa$* .

## 1.4 The structure of the paper

Section 2 describes the algorithm, and introduces all the necessary notions. Section 3 presents the consistency and sample complexity results. Section 4 discusses run time efficiency, and in Section 5 we make some concluding remarks. Finally, the supplementary material provides the missing proofs, the analysis of the special cases, and the necessary fixes for the issues with the run-time complexity.

## 2 StOP: Stochastic Optimistic Planning

Recall that  $N \in \mathbb{N}$  denotes the number of possible next states. That is, for each state  $x \in X$  and each action  $u$  available at  $x$ , it holds that  $P(x, u, x') = 0$  for all but at most  $N$  states  $x' \in X$ . Throughout this section, the state of interest is denoted by  $x_0$ , the requested accuracy by  $\epsilon$ , and the confidence parameter by  $\delta_0$ . That is, the problem to be solved is to output an action  $u$  which is, with probability at least  $(1 - \delta_0)$ , at least  $\epsilon$ -optimal in  $x_0$ .

The algorithm and the analysis make use of the notion of an (infinite) planning tree, policies and trajectories. These notions are introduced in the next subsection.

### 2.1 Planning trees and trajectories

The *infinite planning tree*  $\Pi^\infty$  for a given MDP is a rooted and labeled infinite tree. Its root is denoted  $s_0$  and is labeled by the state of interest,  $x_0 \in X$ . Nodes on even levels are called *action nodes* (the root is an action node), and have  $K_d$  children each on the  $d$ -th level of action nodes: each action  $u$  is represented by exactly one child, labeled  $u$ . Nodes on odd levels are called *transition nodes* and have  $N$  children each: if the label of the parent (action) node is  $x$ , and the label of the transition node itself is  $u$ , then for each  $x' \in X$  with  $P(x, u, x') > 0$  there is a corresponding child, labeled  $x'$ . There may be children with probability zero, but no duplicates.

An *infinite policy* is a subtree of  $\Pi^\infty$  with the same root, where each action node has exactly one child and each transition node has  $N$  children. It corresponds to an agent having fixed all its possible future actions. A (*partial*) *policy*  $\Pi$  is a finite subtree of  $\Pi^\infty$ , again with the same root, but where the action nodes have *at most* one child, each transition node has  $N$  children, and all leaves<sup>4</sup> are on the same level. The number of transition nodes on any path from the root to a leaf is denoted  $d(\Pi)$  and is called the *depth* of  $\Pi$ . A partial policy corresponds to the agent having its possible future actions planned for  $d(\Pi)$  steps. There is a natural partial order over these policies: a policy

<sup>4</sup>Note that leaves are, by definition, always action nodes.

$\Pi'$  is called *descendant policy* of a policy  $\Pi$  if  $\Pi$  is a subtree of  $\Pi'$ . If, additionally, it holds that  $d(\Pi') = d(\Pi) + 1$ , then  $\Pi$  is called the *parent policy* of  $\Pi'$ , and  $\Pi'$  the *child policy* of  $\Pi$ .

A (*random*) *trajectory*, or *rollout*, for some policy  $\Pi$  is a realization  $\tau := (x_t, u_t, r_t)_{t=0}^T$  of the stochastic process that belongs to the policy. A random path is generated from the root by always following, from a non-leaf action node with label  $x_t$ , its unique child in  $\Pi$ , then setting  $u_t$  to the label of this node, from where, drawing first a label  $x_{t+1}$  from  $P(x_t, u_t, \cdot)$ , one follows the child with label  $x_{t+1}$ . The reward  $r_t$  is drawn from the distribution determined by  $R(x_t, u_t, \cdot)$ . The *value of the rollout*  $\tau$  (also called return or payoff in the literature) is  $v(\tau) := \sum_{t=0}^T r_t \gamma^t$ , and the *value of the policy*  $\Pi$  is  $v(\Pi) := \mathbb{E}[v(\tau)] = \mathbb{E}[\sum_{t=0}^T r_t \gamma^t]$ . For an action  $u$  available at  $x_0$ , denote by  $v(u)$  the maximum of the values of the policies having  $u$  as the label of the child of root  $s_0$ . Denote by  $v^*$  the maximum of these  $v(u)$  values. Using this notation, the task of the algorithm is to return, with high probability, an action  $u$  with  $v(u) \geq v^* - \epsilon$ .

## 2.2 The algorithm

STOP (Algorithm 1, see Figure 1 in the supplementary material for an illustration) maintains for each action  $u$  available at  $x_0$  a set of *active policies*  $\text{Active}(u)$ . Initially, it holds that  $\text{Active}(u) = \{\Pi_u\}$ , where  $\Pi_u$  is the shallowest partial policy with the child of the root being labeled  $u$ . Also, for each policy  $\Pi$  that becomes a member of an active set, the algorithm maintains high confidence lower and upper bounds for the value  $v(\Pi)$  of the policy, denoted  $\nu(\Pi)$  and  $b(\Pi)$ , respectively.

In each round  $t$ , an *optimistic policy*  $\Pi_{t,u}^\dagger := \arg\max_{\Pi \in \text{Active}(u)} b(\Pi)$  is determined for each action  $u$ . Based on this, the current *optimistic action*  $u_t^\dagger := \arg\max_u b(\Pi_{t,u}^\dagger)$  and *secondary action*  $u_t^{\dagger\dagger} := \arg\max_{u \neq u_t^\dagger} b(\Pi_{t,u}^\dagger)$  are computed. A policy  $\Pi_t$  to explore is then chosen: if the one that belongs to the secondary action is at least as deeply developed as the one that belongs to the optimistic action, the latter is chosen for exploration, and otherwise the former. Note that a smaller depth is equivalent to a larger gap between lower and upper bound, and vice versa<sup>5</sup>. The set  $\text{Active}(u_t)$  is then updated, replacing the policy  $\Pi_t$  by its children policies. Accordingly, the upper and lower bounds for these policies are computed. The algorithm terminates when  $\nu(\Pi_t^\dagger) + \epsilon \geq \max_{u \neq u_t^\dagger} b(\Pi_{t,u}^\dagger)$ —that is, when, with high confidence, no policies starting with an action different from  $u_t^\dagger$  have the potential to have significantly higher value.

### 2.2.1 Number and length of trajectories needed for one partial policy

Fix some integer  $d > 0$  and let  $\Pi$  be a partial policy of depth  $d$ . Let, furthermore,  $\Pi'$  be an infinite policy that is a descendant of  $\Pi$ . Note that

$$0 \leq v(\Pi') - v(\Pi) \leq \frac{\gamma^d}{1-\gamma}. \quad (2)$$

The value of  $\Pi$  is a  $\frac{\gamma^d}{1-\gamma}$ -accurate approximation of the value of  $\Pi'$ . On the other hand, having  $m$  trajectories for  $\Pi$ , their average reward  $\hat{v}(\Pi)$  can be used as an estimate of the value  $v(\Pi)$  of  $\Pi$ . From the Hoeffding bound, this estimate has, with probability at least  $(1 - \delta)$ , accuracy  $\frac{1-\gamma^d}{1-\gamma} \sqrt{\frac{\ln(1/\delta)}{2m}}$ .

With  $m := m(d, \delta) := \lceil \frac{\ln(1/\delta)}{2} (\frac{1-\gamma^d}{\gamma^d})^2 \rceil$  trajectories,  $\frac{\gamma^d}{1-\gamma} \geq \frac{1-\gamma^d}{1-\gamma} \sqrt{\frac{\ln(1/\delta)}{2m}}$  holds, so with probability at least  $(1 - \delta)$ ,  $b(\Pi) := \hat{v}(\Pi) + \frac{\gamma^d}{1-\gamma} + \frac{1-\gamma^d}{1-\gamma} \sqrt{\frac{\ln(1/\delta)}{2m}} \leq \hat{v}(\Pi) + 2\frac{\gamma^d}{1-\gamma}$  and  $\nu(\Pi) := \hat{v}(\Pi) - \frac{1-\gamma^d}{1-\gamma} \sqrt{\frac{\ln(1/\delta)}{2m}} \geq \hat{v}(\Pi) - \frac{\gamma^d}{1-\gamma}$  bound  $v(\Pi')$  from above and below, respectively. This choice balances the inaccuracy of estimating  $v(\Pi')$  based on  $v(\Pi)$  and the inaccuracy of estimating  $v(\Pi)$ .

Let  $d^* := d^*(\epsilon, \gamma) := \lceil (\ln \frac{6}{(1-\gamma)\epsilon}) / \ln(1/\gamma) \rceil$ , the smallest integer satisfying  $3\frac{\gamma^{d^*}}{1-\gamma} \leq \epsilon/2$ . Note that if  $d(\Pi) = d^*$  for any given policy  $\Pi$ , then  $b(\Pi) - \nu(\Pi) \leq \epsilon/2$ . Because of this, it follows (see Lemma 3 in the supplementary material) that  $d^*$  is the maximal length the algorithm ever has to develop a policy.

<sup>5</sup>This approach of using secondary actions is based on the UGapE algorithm [10].

---

**Algorithm 1**  $\text{StOP}(s_0, \delta_0, \epsilon, \gamma)$ 

---

```
1: for all  $u$  available from  $x_0$  do ▷ initialize
2:    $\Pi_u :=$  smallest policy with the child of  $s_0$  labeled  $u$ 
3:    $\delta_1 := (\delta_0/d^*) \cdot (K_0)^{-1}$  ▷  $d(\Pi_u) = 1$ 
4:    $(\nu(\Pi_u), b(\Pi_u)) := \text{BoundValue}(\Pi_u, \delta_1)$ 
5:    $\text{Active}(u) := \{\Pi_u\}$  ▷ the set of active policies that follow  $u$  in  $s_0$ 
6: for round  $t=1, 2, \dots$  do
7:   for all  $u$  available at  $x_0$  do
8:      $\Pi_{t,u}^\dagger := \text{argmax}_{\Pi \in \text{Active}(u)} b(\Pi)$ 
9:      $\Pi_t^\dagger := \Pi_{t,u_t^\dagger}^\dagger$ , where  $u_t^\dagger := \text{argmax}_u b(\Pi_{t,u}^\dagger)$ , ▷ optimistic action and policy
10:     $\Pi_t^{\dagger\dagger} := \Pi_{t,u_t^{\dagger\dagger}}^\dagger$ , where  $u_t^{\dagger\dagger} := \text{argmax}_{u \neq u_t^\dagger} b(\Pi_{t,u}^\dagger)$ , ▷ secondary action and policy
11:    if  $\nu(\Pi_t^\dagger) + \epsilon \geq \max_{u \neq u_t^\dagger} b(\Pi_{t,u}^\dagger)$  then ▷ termination criterion
12:      return  $u_t^\dagger$ 
13:    if  $d(\Pi_t^{\dagger\dagger}) \geq d(\Pi_t^\dagger)$  then ▷ select the policy to evaluate
14:       $u_t := u_t^\dagger$  and  $\Pi_t := \Pi_t^\dagger$ 
15:    else
16:       $u_t := u_t^{\dagger\dagger}$  and  $\Pi_t := \Pi_t^{\dagger\dagger}$  ▷ action and policy to explore
17:     $\text{Active}(u_t) := \text{Active}(u_t) \setminus \{\Pi_t\}$ 
18:     $\delta := (\delta_0/d^*) \cdot \prod_{\ell=0}^{d(\Pi_t)-1} (K_\ell)^{-N^\ell}$  ▷  $\prod_{\ell=0}^{d-1} (K_\ell)^{N^\ell} = \#$  of policies of depth at most  $d$ 
19:    for all child policy  $\Pi'$  of  $\Pi_t$  do
20:       $(\nu(\Pi'), b(\Pi')) := \text{BoundValue}(\Pi', \delta)$ 
21:       $\text{Active}(u_t) := \text{Active}(u_t) \cup \{\Pi'\}$ 
```

---

## 2.2.2 Samples and sample trees

Algorithm  $\text{StOP}$  aims to aggressively reuse every sample for each transition node and every sample for each state-action pair, in order to keep the sample complexity as low as possible. Each time the value of a partial policy is evaluated, all samples that are available for any part of it from previous rounds are reused. That is, if  $m$  trajectories are necessary for assessing the value of some policy  $\Pi$ , and there are  $m'$  complete trajectories available and  $m''$  that end at some inner node of  $\Pi$ , then  $\text{StOP}$  (more precisely, another algorithm,  $\text{Sample}$ , called from  $\text{StOP}$ ) samples rewards (using  $\text{SampleReward}$ ) and transitions ( $\text{SampleTransition}$ ) to generate continuations for the  $m''$  incomplete trajectories and to generate  $(m - m' - m'')$  new trajectories, as described in Section 2.1, where

- $\text{SampleReward}(s)$  for some action node  $s$  samples a reward from the distribution  $R(x, u, \cdot)$ , where  $u$  is the label of the parent of  $s$  and  $x$  is the label of the grandparent of  $s$ , and
- $\text{SampleTransition}(s)$  for some transition node  $s$  samples a next state from the distribution  $P(x, u, \cdot)$ , where  $u$  is the label of  $s$  and  $x$  is the label of the parent of  $s$ .

To compensate for the sharing of the samples, the confidences of the estimates are increased, so that with probability at least  $(1 - \delta_0)$ , all of them are valid<sup>6</sup>. The samples are organized as a collection of sample trees, where a *sample tree*  $\mathcal{T}$  is a (finite) subtree of  $\mathbf{\Pi}^\infty$  with the property that each transition node has exactly one child, and that each action node  $s$  is associated with some reward  $r^\mathcal{T}(s)$ . Note that the intersection of a policy  $\Pi$  and a sample tree  $\mathcal{T}$  is always a path. Denote this path by  $\tau(\mathcal{T}, \Pi)$  and note that it necessarily starts from the root and ends either in a leaf or in an internal node of  $\Pi$ . In the former case, this path can be interpreted as a complete trajectory for  $\Pi$ , and in the latter case, as an initial segment. Accordingly, when the value of a new policy  $\Pi$  needs to be estimated/bounded, it is computed as  $\hat{v}(\Pi) := \frac{1}{m} \sum_{i=1}^m v(\tau(\mathcal{T}_i, \Pi))$  (see Algorithm 2:  $\text{BoundValue}$ ), where  $\mathcal{T}_1, \dots, \mathcal{T}_m$  are sample trees constructed by the algorithm. For terseness, these are considered to be global variables, and are constructed and maintained using algorithm  $\text{Sample}$  (Algorithm 3).

---

<sup>6</sup>In particular, the confidence is set to  $1 - \delta_{d(\Pi)}$  for policy  $\Pi$ , where  $\delta_d = (\delta_0/d^*) \prod_{\ell=0}^{d-1} K_\ell^{-N^\ell}$  is  $\delta_0$  divided by the number of policies of depth at most  $d$ , and by the largest possible depth—see section 2.2.1.

---

**Algorithm 2** BoundValue( $\Pi, \delta$ )

---

**Ensure:** with probability at least  $(1 - \delta)$ , interval  $[\nu(\Pi), b(\Pi)]$  contains  $v(\Pi)$

- 1:  $m := \left\lceil \frac{\ln(1/\delta)}{2} \left( \frac{1-\gamma^{d(\Pi)}}{\gamma^{d(\Pi)}} \right)^2 \right\rceil$
  - 2:  $\text{Sample}(\Pi, s_0, m)$   $\triangleright$  Ensure that at least  $m$  trajectories exist for  $\Pi$
  - 3:  $\hat{v}(\Pi) := \frac{1}{m} \sum_{i=1}^m v(\tau(\mathcal{T}_i, \Pi))$   $\triangleright$  empirical estimate of  $v(\Pi)$
  - 4:  $\nu(\Pi) := \hat{v}(\Pi) - \frac{1-\gamma^{d(\Pi)}}{1-\gamma} \sqrt{\frac{\ln(1/\delta)}{2m}}$   $\triangleright$  Hoeffding bound
  - 5:  $b(\Pi) := \hat{v}(\Pi) + \frac{\gamma^{d(\Pi)}}{1-\gamma} + \frac{1-\gamma^{d(\Pi)}}{1-\gamma} \sqrt{\frac{\ln(1/\delta)}{2m}}$   $\triangleright \dots$  and (2)
  - 6: **return**  $(\nu(\Pi), b(\Pi))$
- 

---

**Algorithm 3** Sample( $\Pi, s, m$ )

---

**Ensure:** there are  $m$  sample trees  $\mathcal{T}_1, \dots, \mathcal{T}_m$  that contain a complete trajectory for  $\Pi$  (i.e.  $\tau(\mathcal{T}_i, \Pi)$  ends in a leaf of  $\Pi$  for  $i = 1, \dots, m$ )

- 1: **for**  $i := 1, \dots, m$  **do**
  - 2:     **if** sample tree  $\mathcal{T}_i$  does not yet exist **then**
  - 3:         let  $\mathcal{T}_i$  be a new sample tree of depth 0
  - 4:     let  $s$  be the last node of  $\tau(\mathcal{T}_i, \Pi)$   $\triangleright$   $s$  is an action node
  - 5:     **while**  $s$  is not a leaf of  $\Pi$  **do**
  - 6:         let  $s'$  be the child of  $s$  in  $\Pi$  and add it to  $\mathcal{T}$  as a new child of  $s$
  - 7:          $s'' := \text{SampleTransition}(s')$ ,  $\triangleright$   $s'$  is a transition node
  - 8:         add  $s''$  to  $\mathcal{T}$  as a new child of  $s'$
  - 9:          $s := s''$
  - 10:         $r^{\mathcal{T}}(s'') := \text{SampleReward}(s'')$
- 

### 3 Analysis

Recall that  $v^*$  denotes the maximal value of any (possibly infinite) policy tree. The following theorem formalizes the consistency result for StOP (see the proof in Section C).

**Theorem 1.** *With probability at least  $(1 - \delta_0)$ , StOP returns an action with value at least  $v^* - \epsilon$ .*

Before stating the sample complexity result, some further notation needs to be introduced.

Let  $u^*$  denote an optimal action available at state  $x_0$ . That is,  $v(u^*) = v^*$ . Define for  $u \neq u^*$

$$\mathcal{P}_u^\epsilon := \left\{ \Pi : \Pi \text{ follows } u \text{ from } s_0 \text{ and } v(\Pi) + 3\frac{\gamma^{d(\Pi)}}{1-\gamma} \geq v^* - 3\frac{\gamma^{d(\Pi)}}{1-\gamma} + \epsilon \right\},$$

and also define

$$\mathcal{P}_{u^*}^\epsilon := \left\{ \Pi : \Pi \text{ follows } u^* \text{ from } s_0, v(\Pi) + 3\frac{\gamma^{d(\Pi)}}{1-\gamma} \geq v^* \text{ and } v(\Pi) - 6\frac{\gamma^{d(\Pi)}}{1-\gamma} + \epsilon \leq \max_{u \neq u^*} v(u) \right\}.$$

Then  $\mathcal{P}^\epsilon := \mathcal{P}_{u^*}^\epsilon \cup \bigcup_{u \neq u^*} \mathcal{P}_u^\epsilon$  is the set of “important” policies that potentially need to be evaluated in order to determine an  $\epsilon$ -optimal action. (See also Lemma 8 in the supplementary material.)

Let now  $p(s)$  denote the product of the probabilities of the transitions on the path from  $s_0$  to  $s$ . That is, for any policy tree  $\Pi$  containing  $s$ , a trajectory for  $\Pi$  goes through  $s$  with probability  $p(s)$ . When estimating the value of some policy  $\Pi$  of depth  $d$ , the expected number of trajectories going through some nodes  $s$  of it is  $p(s)m(d, \delta_d)$ . The sample complexity therefore has to take into consideration for each node  $s$  (at least for the ones with “high”  $p(s)$  value) the maximum  $\ell(s) = \max\{d(\Pi) : \Pi \in \mathcal{P}^\epsilon \text{ contains } s\}$  of the depth of the relevant policies it is included in. Therefore, the expected number of trajectories going through  $s$  in a given run of StOP is

$$p(s) \cdot m(\ell(s), \delta_{\ell(s)}) = p(s) \left\lceil \frac{\ln(1/\delta_{\ell(s)})}{2} \left( \frac{1-\gamma^{\ell(s)}}{\gamma^{\ell(s)}} \right)^2 \right\rceil \quad (3)$$

If (3) is “large” for some  $s$ , it can be used to deduce high confidence upper bound on the number of times  $s$  gets sampled. To this end, let  $S^\epsilon$  denote the set of nodes of the trees in  $\mathcal{P}^\epsilon$ , let  $\mathcal{N}^\epsilon$  denote the

smallest positive integer  $\mathcal{N}$  satisfying  $\mathcal{N} \geq |\{s \in \mathcal{S}^\epsilon : p(s) \cdot m(\ell(s), \delta_{\ell(s)}) \geq (8/3) \ln(2\mathcal{N}/\delta_0)\}|$  (obviously  $\mathcal{N}^\epsilon \leq |\mathcal{S}^\epsilon|$ ), and let

$$\mathcal{S}^{\epsilon,*} := \{s \in \mathcal{S}^\epsilon : p(s) \cdot m(\ell(s), \delta_{\ell(s)}) \geq (8/3) \ln(2\mathcal{N}^\epsilon/\delta_0)\}$$

Then  $\mathcal{S}^\epsilon$  is the set of important nodes (since  $\mathcal{P}^\epsilon$  is the set of ‘‘important’’ policies), and  $\mathcal{S}^{\epsilon,*}$  consists of the important nodes which, with high probability, are not sampled more than twice they are expected to be. (This high probability is  $1 - \frac{\delta_0}{2\mathcal{N}^\epsilon}$  according to the Bernstein bound, and so these upper bounds hold jointly with probability at least  $(1 - \frac{\delta_0}{2})$ , as  $\mathcal{N}^\epsilon = |\mathcal{S}^{\epsilon,*}|$ . See also Appendix D.)

The number of times some  $s' \in \mathcal{S}^\epsilon \setminus \mathcal{S}^{\epsilon,*}$  gets sampled has too large variance compared to its expected value (3), so a different approach is needed in order to derive high confidence upper bounds. To this end, for a transition node  $s$ , let  $p^\circ(s) := p^\circ(s, \epsilon) := \sum \{p(s') : s' \text{ is a child of } s \text{ with } p(s') \cdot m(\ell(s'), \delta_{\ell(s')}) < (8/3) \ln(2\mathcal{N}^\epsilon/\delta_0)\}$ , and

$$B(s) := B(s, \epsilon) := \begin{cases} 0, & \text{if } p^\circ(s) \leq \frac{\delta}{2\mathcal{N}^\epsilon m(\ell(s), \delta_{\ell(s)})} \\ \max(6 \ln(\frac{2\mathcal{N}^\epsilon}{\delta_0}), 2p^\circ(s)m(\ell(s), \delta_{\ell(s)})) & \text{otherwise} \end{cases}$$

As it will be shown in the proof of Theorem 2 (in Section D), this is a high confidence upper bound on the number of trajectories that go through some child  $s' \in \mathcal{S}^\epsilon \setminus \mathcal{S}^{\epsilon,*}$  of some  $s \in \mathcal{S}^{\epsilon,*}$ .

**Theorem 2.** *With probability at least  $(1 - 2\delta)$ ,  $\text{StOP}$  outputs a policy of value at least  $(v^* - \epsilon)$  after generating at most  $\sum_{s \in \mathcal{S}^{\epsilon,*}} \left( 2p(s)m(\ell(s), \delta_{\ell(s)}) + B(s) \sum_{d=d(s)+1}^{\ell(s)} \prod_{\ell=d(s)+1}^d K_\ell \right)$  samples, where  $d(s) = \min\{d(\Pi) : s \text{ appears in policy } \Pi\}$  is the depth of node  $s$ .*

Finally, the bound discussed in Section 1 is obtained by setting  $\kappa := \limsup_{\epsilon \rightarrow 0} \max(\kappa_1, \kappa_2)$ , where  $\kappa_1 := \kappa_1(\epsilon, \delta_0, \gamma) := \left( \sum_{s \in \mathcal{S}^{\epsilon,*}} \frac{\epsilon^2(1-\gamma)^2}{\ln(1/\delta_0)} 2p(s)m(\ell(s), \delta_{\ell(s)}) \right)^{1/d^*}$  and  $\kappa_2 := \kappa_2(\epsilon, \delta_0, \gamma) := \left( \frac{\epsilon^2(1-\gamma)^2}{\ln(1/\delta_0)} \sum_{s \in \mathcal{S}^{\epsilon,*}} B(s) \sum_{d=d(s)}^{\ell(s)} \prod_{\ell=d(s)}^d K_\ell \right)^{1/d^*}$ .

## 4 Efficiency

$\text{StOP}$ , as presented in Algorithm 1, is not efficiently executable. First of all, whenever it evaluates an optimistic policy, it enumerates all its children policies, which has typically exponential time complexity. Besides that, the sample trees are also treated in an inefficient way. An efficient version of  $\text{StOP}$  with all these issues fixed is presented in Appendix F of the supplementary material.

## 5 Concluding remarks

In this work, we have presented and analyzed our algorithm,  $\text{StOP}$ . To the best of our knowledge,  $\text{StOP}$  is currently the only algorithm for optimal (i.e. closed loop) online planning with a generative model that provably benefits from local structure both in reward as well as in transition probabilities. It assumes no knowledge about this structure other than access to the generative model, and does not impose any restrictions on the system dynamics.

One should note though that the current version of  $\text{StOP}$  does not support domains with infinite  $N$ . The sparse sampling algorithm in [14] can easily handle such problems (at the cost of a non-polynomial (in  $1/\epsilon$ ) sample complexity), however,  $\text{StOP}$  has much better sample complexity in case of finite  $N$ . An interesting problem for future research is to design adaptive planning algorithms with sample complexity independent of  $N$  ([21] presents such an algorithm, but the complexity bound provided there is the same as the one in [14]).

## Acknowledgments

This work was supported by the French Ministry of Higher Education and Research, and by the European Community’s Seventh Framework Programme (FP7/2007-2013) under grant agreement n° 270327 (project ComplACS). Author two would like to acknowledge the support of the BMBF project ALICE (01IB10003B).



## References

- [1] P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning Journal*, 47(2-3):235–256, 2002.
- [2] Dimitri P. Bertsekas. *Dynamic Programming and Optimal Control*. Athena Scientific, 2001.
- [3] S. Bubeck and R. Munos. Open loop optimistic planning. In *Conference on Learning Theory*, 2010.
- [4] Sébastien Bubeck and Nicolò Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 5(1):1–122, 2012.
- [5] Lucian Buşoniu and Rémi Munos. Optimistic planning for markov decision processes. In *Proceedings 15th International Conference on Artificial Intelligence and Statistics (AISTATS-12)*, pages 182–189, 2012.
- [6] E. F. Camacho and C. Bordons. *Model Predictive Control*. Springer-Verlag, 2004.
- [7] Nicolo Cesa-Bianchi and Gabor Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, New York, NY, USA, 2006.
- [8] Rémi Coulom. Efficient selectivity and backup operators in Monte-Carlo tree search. In *Proceedings Computers and Games 2006*. Springer-Verlag, 2006.
- [9] E. Even-Dar, S. Mannor, and Y. Mansour. Action elimination and stopping conditions for reinforcement learning. In T. Fawcett and N. Mishra, editors, *Proceedings of the Twentieth International Conference on Machine Learning (ICML-2003)*, pages 162–169, 2003.
- [10] Victor Gabillon, Mohammad Ghavamzadeh, and Alessandro Lazaric. Best arm identification: A unified approach to fixed budget and fixed confidence. In Peter L. Bartlett, Fernando C. N. Pereira, Christopher J. C. Burges, Lon Bottou, and Kilian Q. Weinberger, editors, *NIPS*, pages 3221–3229, 2012.
- [11] Sylvain Gelly, Yizao Wang, Rémi Munos, and Olivier Teytaud. Modification of UCT with Patterns in Monte-Carlo Go. Rapport de recherche RR-6062, INRIA, 2006.
- [12] Eric A. Hansen and Shlomo Zilberstein. A heuristic search algorithm for Markov decision problems. In *Proceedings Bar-Ilan Symposium on the Foundation of Artificial Intelligence*, Ramat Gan, Israel, 23–25 June 1999.
- [13] J-F. Hren and R. Munos. Optimistic planning of deterministic systems. In *Recent Advances in Reinforcement Learning*, pages 151–164. Springer LNAI 5323, European Workshop on Reinforcement Learning, 2008.
- [14] M. Kearns, Y. Mansour, and A.Y. Ng. A sparse sampling algorithm for near-optimal planning in large Markovian decision processes. In *Machine Learning*, volume 49, pages 193–208, 2002.
- [15] Gunnar Kedenburg, Raphael Fonteneau, and Remi Munos. Aggregating optimistic planning trees for solving markov decision processes. In C.J.C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 2382–2390. Curran Associates, Inc., 2013.
- [16] Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In *In: ECML-06. Number 4212 in LNCS*, pages 282–293. Springer, 2006.
- [17] T. L. Lai and H. Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6:4–22, 1985.
- [18] Rémi Munos. From bandits to Monte-Carlo Tree Search: The optimistic principle applied to optimization and planning. *Foundation and Trends in Machine Learning*, 7(1):1–129, 2014.
- [19] N.J. Nilsson. *Principles of Artificial Intelligence*. Tioga Publishing, 1980.
- [20] M.L. Puterman. *Markov Decision Processes — Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, NY, 1994.
- [21] Thomas J. Walsh, Sergiu Goschin, and Michael L. Littman. Integrating sample-based planning and model-based reinforcement learning. In *Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence*, pages 612–617. AAAI Press, 2010.

## A Illustration of the StOP algorithm

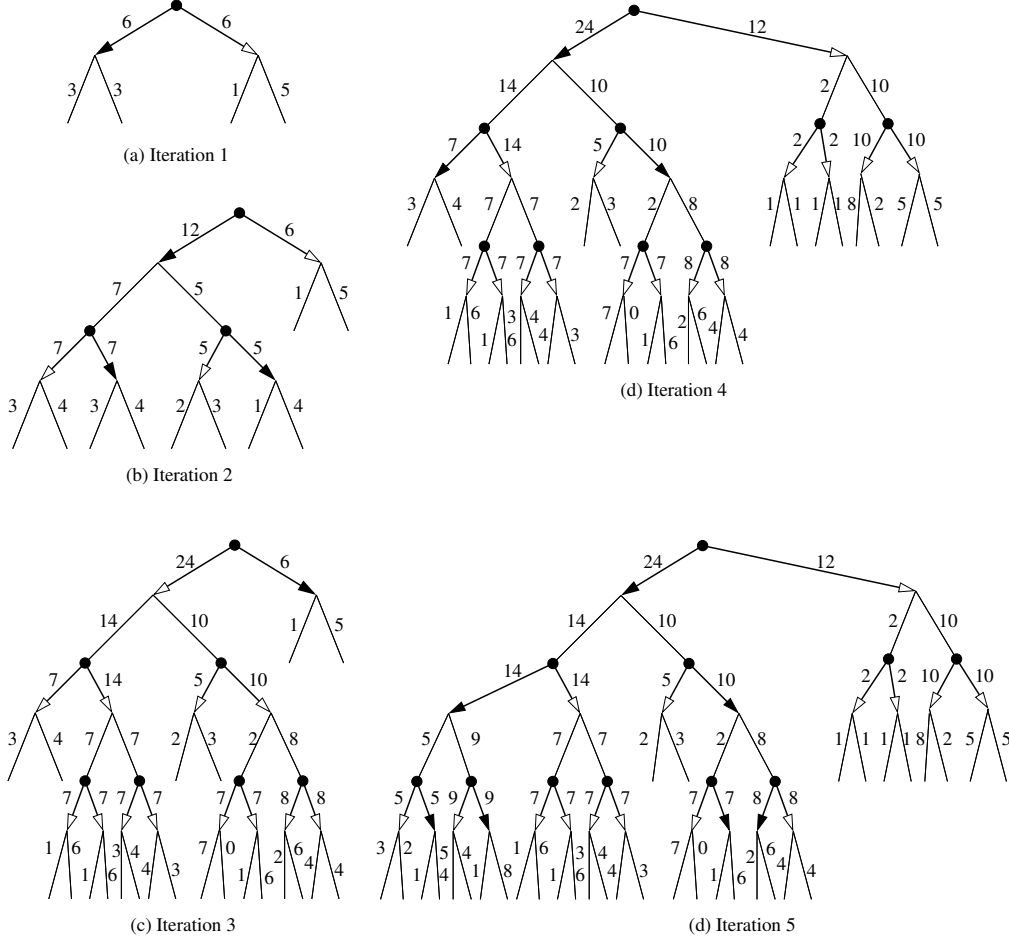


Figure 1: Illustration of the StOP algorithm with  $K = N = 2$ . Black dots represents action-nodes and thick arrows transition-nodes. Thin arrow represents transitions to next action-nodes. The numbers corresponds to the number of samples allocated to each node or transition. For example in Iteration 1, the procedure `Sample` allocated 6 samples to each action. The optimistic policy  $\Pi^\dagger$  is selected (Step 11 of StOP), which is shown by the black arrows. At iteration 2, the leaves of the optimistic policy are expanded and `Sample` generates more samples along the new possible policies. The new optimistic policy is computed. The same process is repeated in later iterations. Notice that the same samples are used to evaluate many policies, and that the leaves of the optimistic policy in Iteration 4 are not all leaves of the whole tree.

## B Chernoff-Hoeffding and Bernstein bounds

This section provides a quick overview of the specific concentration inequalities that are used to obtain high confidence bounds on the values of the policies. The first one is the Hoeffding bound (Corollary A.1 in [7]). It implies that for any given random variable that takes values from the interval  $[0, a]$  and has expected value  $p$ , the average  $p_m$  of  $m$  independent samples satisfy

$$\mathbb{P}\left[\hat{p}_m \leq p - a\sqrt{\frac{\ln(1/\delta)}{2m}}\right] \leq \delta \text{ and } \mathbb{P}\left[\hat{p}_m \geq p + a\sqrt{\frac{\ln(1/\delta)}{2m}}\right] \leq \delta.$$

The second concentration inequality is the Bernstein bound (see e.g. Corollary A.3 in [7]). It implies that for any given  $a > 0$  and for any given Bernoulli variable with parameter  $p$ , the average  $p_m$  of  $m$  independent samples satisfy  $\mathbb{P}[\hat{p}_m > p + a] \leq \exp\left(\frac{-a^2 m}{2p + 2a/3}\right)$  and  $\mathbb{P}[\hat{p}_m < p - a] \leq$

$\exp\left(\frac{-a^2 m}{2p+2a/3}\right)$ . In particular, setting  $a = p$ , one obtains that

$$pm \geq \frac{8}{3} \ln(1/\delta) \Rightarrow \mathbb{P}[\hat{p}_m > 2p] = \mathbb{P}[\hat{p}_m > p + a] \leq \exp\left(\frac{-pm}{8/3}\right) \leq \delta. \quad (4)$$

Similarly, setting  $a = \frac{8 \ln(1/\delta)}{3m}$ , one obtains that

$$pm < \frac{8 \ln(1/\delta)}{3} \Rightarrow \mathbb{P}\left[\hat{p}_m > \frac{16 \ln(1/\delta)}{3m}\right] \leq \mathbb{P}[\hat{p}_m > p + a] \leq \exp\left(\frac{-am}{8/3}\right) = \delta. \quad (5)$$

## C Proof of the consistency result (Theorem 1)

**Lemma 3.** *There can not be an active policy of depth larger than  $d^*$ .*

*Proof.* For a policy with depth larger than  $d^*$  to be in an active policy set, there has to be a round  $t$  with  $d(\Pi_t) = d^*$ . This can only be the case if  $d(\Pi_t^\dagger) = d^*$  or  $d(\Pi_t^{\dagger\dagger}) = d^*$ . However, if  $d(\Pi_t^\dagger) \geq d^*$ , then it holds that  $\nu(\Pi_t^\dagger) + \epsilon/2 \geq b(\Pi_t^\dagger) \geq \max_{u \neq u_t^\dagger} b(\Pi_{t,u}^\dagger)$ , so `StOP` terminates. And since the selection rule for  $u_t$  implies that  $\Pi^{\dagger\dagger}$  is only selected as  $\Pi_t$  if  $d(\Pi_t^\dagger) > d(\Pi_t^{\dagger\dagger})$ , selecting it would mean  $d(\Pi_t^\dagger) > d^*$ , so the algorithm would terminate by the first argument.  $\square$

For convenience, we restate the theorem.

**Theorem 4** (Restatement of the consistency result, Theorem 1). *With probability at least  $(1 - \delta_0)$  `StOP` returns an action with value at least  $v^* - \epsilon$ .*

To prove the consistency of `StOP`, the following guarantee of `BoundValue` is needed.

**Claim 5.** *With probability at least  $(1 - \delta)$ , `BoundValue`( $\Pi, \delta$ ) sets  $\hat{v}(\Pi)$  to some value in the interval  $\left[v(\Pi) - \frac{1-\gamma^{d(\Pi)}}{1-\gamma} \sqrt{\frac{\ln(1/d)}{2m}}, v(\Pi) + \frac{1-\gamma^{d(\Pi)}}{1-\gamma} \sqrt{\frac{\ln(1/d)}{2m}}\right]$ .*

*Proof.* As discussed in Section 2.2.2, each  $\tau(\mathcal{T}_i, \Pi)$  for  $i = 1, \dots, m$  can be interpreted as trajectories for  $\Pi$  that are independent (because the samples are also independent of each other). Therefore, the average of their value (return)  $\hat{v}(\Pi) = (1/m) \sum_{i=1}^m v(\tau(\mathcal{T}_i, \Pi))$  is an unbiased estimate of  $v(\Pi)$ . What is more, according to the Hoeffding bound (recall Section 2.2.1), the accuracy of this estimate is  $\frac{1-\gamma^{d(\Pi)}}{1-\gamma} \sqrt{\frac{\ln(1/d)}{2m}} \leq \frac{\gamma^{d(\Pi)}}{1-\gamma}$ , with probability at least  $1 - \delta$ .  $\square$

Based on this it is now easy to show that the estimates used by the algorithm are all correct with high probability.

**Corollary 6.** *The event that for every round  $t$  throughout the run of the algorithm, for each action  $u$  available at  $x_0$ , for each  $\Pi \in \text{Active}_t(u)$ , and for each descendant  $\Pi'$  of  $\Pi$  (allowing  $\Pi' = \Pi$ ), the value  $v(\Pi')$  of  $\Pi'$  belongs to the interval  $[v(\Pi), b(\Pi)]$  has probability at least  $(1 - \delta_0)$ , and implies  $\nu(\Pi_{t,u}^\dagger) \leq v(u) \leq b(\Pi_{t,u}^\dagger)$ .*

*Proof.* If `BoundValue` is ever called for some policy  $\Pi$ , then it is called with confidence parameter  $\delta$  set to  $\delta_d = (\delta_0/d^*) \prod_{\ell=1}^d K_\ell$ , where  $d = d(\Pi)$  is the depth of  $\Pi$ . Note also that  $\prod_{\ell=0}^{d-1} (K_\ell)^{N^\ell}$  is the number of partial policies of depth  $d$ , and therefore, based on Claim 5 and Lemma 3, with probability at least  $1 - \sum_{d=1}^{d^*} \delta_d \prod_{\ell=0}^{d-1} (K_\ell)^{N^\ell} = 1 - \delta_0$ , for every  $\Pi$  that ever belongs to the set of active policies,  $v(\Pi) \in \left[\hat{v}(\Pi) - \frac{1-\gamma^{d(\Pi)}}{1-\gamma_{a\Pi}} \sqrt{\frac{\ln(1/d)}{2m}}, \hat{v}(\Pi) + \frac{1-\gamma^{d(\Pi)}}{1-\gamma_{d(\Pi)}} \sqrt{\frac{\ln(1/d)}{2m}}\right]$ . The claimed result now follows from (2).  $\square$

The consistency result of Theorem 1 follows immediately from Corollary 6, Lemma 3 and the termination condition of `StOP`.

## D Proof of the sample complexity (Theorem 2)

For convenience, we restate the theorem.

**Theorem 7** (Restatement of the sample complexity bound, Theorem 2). *With probability at least  $(1 - 2\delta)$ ,  $StOP$  outputs a policy of value at least  $(v^* - \epsilon)$  after generating at most*

$$\sum_{s \in \mathcal{S}^{\epsilon, *}} \left( 2p(s)m(\ell(s), \delta_{\ell(s)}) + B(s) \sum_{d=d(s)+1}^{\ell(s)} \prod_{\ell=d(s)+1}^d K_{\ell} \right) \quad (6)$$

*samples, where  $d(s) = \min\{d(\Pi) : s \text{ appears in policy } \Pi\}$  is the depth of node  $s$ .*

For the proof we need that  $\mathcal{P}^{\epsilon}$  does indeed contain, with high probability, all the important policies. The following lemma is essential for this.

**Lemma 8.** *Assume that for each  $t \geq 0$ , for each action available at  $x_0$ , for each policy  $\Pi \in \text{Active}_t(u)$ ,  $\nu(\Pi) \leq v(\Pi) \leq b(\Pi)$ . Then  $\Pi_t \in \mathcal{P}^{\epsilon}$  for every  $t \geq 1$  throughout the whole run of the algorithm, except for maybe the last round.*

*Proof.* Note that, whenever a policy is removed from the set of active policies, it is, actually, replaced by its children policies. So, as  $\Pi_{u^*} \in \text{Active}(u^*)$  initially, in every subsequent step there will be some  $\Pi \in \text{Active}(u^*)$  having a descendant policy of value  $v^*$ . Therefore, by the assumption of the lemma and by Corollary 6,  $b(\Pi_{t, u^*}^{\dagger}) \geq v^*$ , and therefore

$$b(\Pi_t^{\dagger}) \geq b(\Pi_{t, u^*}^{\dagger}) \geq v^* \quad (7)$$

Additionally, the selection rule of  $\Pi_t$  implies

$$d(\Pi_t) \leq \min \left\{ d(\Pi_t^{\dagger}), d(\Pi_t^{\dagger, \dagger}) \right\} \quad (8)$$

For some  $u \neq u^*$  this implies that, whenever  $\Pi_t = \Pi_{t, u}^{\dagger}$  and the termination criterion is not met,

$$\begin{aligned} v(\Pi_t) + 3\frac{\gamma^{d(\Pi_t)}}{1-\gamma} - \epsilon &\geq \nu(\Pi_t) + 3\frac{\gamma^{d(\Pi_t)}}{1-\gamma} - \epsilon && \text{by the assumption} \\ &\geq b(\Pi_t) - \epsilon && \text{by the definition of } b \text{ and } \nu \\ &\geq \max_{u \neq u_t^{\dagger}} b(\Pi_{t, u}^{\dagger}) - \epsilon && \text{by the choice of } \Pi_t \\ &> \nu(\Pi_t^{\dagger}) && \text{termination criterion is not met} \\ &\geq b(\Pi_t^{\dagger}) - 3\frac{\gamma^{d(\Pi_t^{\dagger})}}{1-\gamma} && \text{by the definition of } b \text{ and } \nu \\ &\geq v^* - 3\frac{\gamma^{d(\Pi_t^{\dagger})}}{1-\gamma} && \text{by (7)} \\ &\geq v^* - 3\frac{\gamma^{d(\Pi_t)}}{1-\gamma} && \text{by (8)} \end{aligned}$$

Consequently  $\Pi_t \in \mathcal{P}^{\epsilon}$ .

Similarly, when  $\Pi_t = \Pi_{t, u^*}^{\dagger}$  then  $\{u_t^{\dagger}, u_t^{\dagger, \dagger}\} = \{u^*, u'\}$  for some  $u'$ , and, if the termination criterion is not met, then

$$\begin{aligned} \max_{u \neq u^*} v(u) + 3\frac{\gamma^{d(\Pi_t)}}{1-\gamma} &\geq \max_{u \neq u^*} \nu(\Pi_{t, u}^{\dagger}) + 3\frac{\gamma^{d(\Pi_t)}}{1-\gamma} && \text{by the assumption} \\ &\geq \max_{u \neq u^*} \nu(\Pi_{t, u}^{\dagger}) + 3\frac{\gamma^{d(\Pi_{t, u'}^{\dagger})}}{1-\gamma} && \text{because of (8) and } \{u_t^{\dagger}, u_t^{\dagger, \dagger}\} = \{u^*, u'\} \\ &\geq \nu(\Pi_{t, u'}^{\dagger}) + 3\frac{\gamma^{d(\Pi_{t, u'}^{\dagger})}}{1-\gamma} && \text{because } u' \neq u^* \\ &\geq b(\Pi_{t, u'}^{\dagger}) && \text{by the definition of } b \text{ and } \nu \\ &= \max_{u \neq u^*} b(\Pi_{t, u}^{\dagger}) && \text{because } \{u_t^{\dagger}, u_t^{\dagger, \dagger}\} = \{u^*, u'\} \end{aligned}$$

$$\begin{aligned}
&\geq \max_{u \neq u_t^\dagger} b(\Pi_{t,u}^\dagger) && \text{by the choice of } u_t^\dagger \\
&\geq \nu(\Pi_t^\dagger) + \epsilon && \text{termination criterion is not met} \\
&\geq b(\Pi_t^\dagger) - 3\gamma \frac{d(\Pi_t^\dagger, u)}{1-\gamma} + \epsilon && \text{by the definition of } b \text{ and } \nu \\
&\geq b(\Pi_t) - 3\gamma \frac{d(\Pi_t^\dagger, u)}{1-\gamma} + \epsilon && \text{by the choice of } \Pi_t \\
&\geq b(\Pi_t) - 3\gamma \frac{d(\Pi_t)}{1-\gamma} + \epsilon && \text{by (8)} \\
&\geq v(\Pi_t) - 3\gamma \frac{d(\Pi_t)}{1-\gamma} + \epsilon && \text{by the assumption}
\end{aligned}$$

This, combined with (7), implies that  $\Pi_t \in \mathcal{P}_{u^*}^\epsilon$ .  $\square$

*Proof of Theorem 6.* In the proof it is assumed that  $\Pi_t \in \mathcal{P}^\epsilon$  for every  $t$  throughout the algorithm, except for maybe the last round. According to Lemma 8 and Corollary 6 this holds with probability at least  $(1 - \delta_0)$ .

The assumption implies that all rollouts generated by  $\text{StOP}$  consist of nodes that belong to  $\mathcal{S}^\epsilon$ . It also implies that for any node  $s$  of  $\Pi^\infty$ , the depth of any policy  $\Pi$  that includes  $s$  and is evaluated by  $\text{StOP}$  is bounded by  $\ell(s)$ . The largest amount of samples required by such a policy is thus  $m(\ell(s), \delta_{\ell(s)})$ . Therefore, according to the Bernstein bound (4), for any  $s \in \mathcal{S}^{\epsilon,*}$ , the number of sample trees containing  $s$  is upper bounded by  $2p(s)m(\ell(s), \delta_{\ell(s)})$  with probability at least  $(1 - \delta_0/(2\mathcal{N}^\epsilon))$ , and so this also upper bounds the number of samples that are generated for  $s$ .

It is now only left to upper bound the number of samples that are generated for nodes in  $(\mathcal{S}^\epsilon \setminus \mathcal{S}^{\epsilon,*})$ . For this, first partition these nodes by forming, for each  $s \in \mathcal{S}^{\epsilon,*}$ , a group consisting of all the nodes having  $s$  as their lowest ancestor in  $\mathcal{S}^{\epsilon,*}$ . Note that the probability that a trajectory traverses through this group is  $p^\circ(s)$ , and therefore, according to the Bernstein bound, the number of trajectories that traverses this group is upper bounded by  $B(s)$  with probability at least  $(1 - \delta/(2\mathcal{N}^\epsilon))$ . Indeed, in case  $p^\circ(s)m(\ell(s), \delta_{\ell(s)}) \geq (8/3) \ln(2\mathcal{N}^\epsilon/\delta)$ , the Bernstein bound (4) guarantees the bound  $2p^\circ(s)m(\ell(s), \delta_{\ell(s)})$  with confidence at least  $(1 - \delta/(2\mathcal{N}^\epsilon))$ , otherwise (5) provides the bound  $p^\circ(s)m(\ell(s), \delta_{\ell(s)}) + 3 \ln(2\mathcal{N}^\epsilon/\delta) \leq 6 \ln(2\mathcal{N}^\epsilon/\delta)$ . In fact, when  $p^\circ(s) \leq \delta/(2\mathcal{N}^\epsilon m(\ell(s), \delta_{\ell(s)}))$  then, according to the Bernoulli inequality, with probability at least  $(1 - \delta_0/(2\mathcal{N}^\epsilon))$ , no trajectory traverses through the group. Finally note that a sample tree contains at most  $\sum_{d=d(s)+1}^{\ell(s)} \prod_{\ell=d(s)+1}^d K_\ell$  samples below node  $s$ .  $\square$

## E Worst case bound and special cases

Before we turn to the analysis of the special cases, we discuss shortly the second term in the sample complexity bound (6).

**Claim 9.**  $\sum_{s \in \mathcal{S}^{\epsilon,*}} B(s) \sum_{d=d(s)+1}^{\ell(s)} \prod_{\ell=d(s)+1}^d K_\ell \leq |\mathcal{S}^\epsilon \setminus \mathcal{S}^{\epsilon,*}| \cdot 6 \cdot \ln(\frac{2\mathcal{N}^\epsilon}{\delta_0})$ .

*Proof.* First of all, each  $s \in \mathcal{S}^{\epsilon,*}$  has at least  $p^\circ(s) \cdot (3/8) \cdot m(d, \delta_{\ell(s)})/\ln(2\mathcal{N}^\epsilon/\delta_0)$  children  $s'$  with  $p(s') \cdot m(d, \delta_{\ell(s')}) < (8/3) \ln(2\mathcal{N}^\epsilon/\delta_0)$  (note that  $\ell(s) = \ell(s')$ ), therefore  $\max\left(6 \ln(\frac{2\mathcal{N}^\epsilon}{\delta_0}), 2p^\circ(s)m(\ell(s), \delta_{\ell(s)})\right)$  is upper bounded by the number of these children multiplied by  $6 \ln(2\mathcal{N}^\epsilon/\delta_0)$ . Note also that number of nodes in  $\mathcal{S}^\epsilon$  below  $s'$  is at least  $\sum_{d=d(s)+1}^{\ell(s)} \prod_{\ell=d(s)+1}^d K_\ell$ .

To sum up,  $B(s)$  accounts at most  $6 \ln \frac{2\mathcal{N}^\epsilon}{\delta_0}$  for every  $s' \in \mathcal{S}^\epsilon \setminus \mathcal{S}^{\epsilon,*}$  having  $s$  as its lowest ancestor in  $\mathcal{S}^{\epsilon,*}$ .  $\square$

Now recall that  $d^* = d^*(\epsilon, \gamma) = \left\lceil \frac{\ln((1-\gamma)\epsilon/6)}{\ln \gamma} \right\rceil$ , and also that this implies

$$\epsilon(1-\gamma) \leq 6\gamma^{d^*-1} \quad (9)$$

Defining

$$\begin{aligned}
\kappa_1 &:= \kappa_1(\epsilon, \delta_0, \gamma) := \left( \sum_{s \in \mathcal{S}^{\epsilon, *}} \frac{\epsilon^2(1-\gamma)^2}{\ln(1/\delta_0)} 2p(s)m(\ell(s), \delta_{\ell(s)}) \right)^{1/d^*} \\
&\leq \left( \frac{\epsilon^2(1-\gamma)^2}{\ln(1/\delta_0)} \sum_{s \in \mathcal{S}^{\epsilon, *}} p(s) \cdot \frac{1}{\gamma^{2\ell(s)}} \ln \frac{d^* \prod_{\ell=1}^{\ell(s)} (K_\ell)^{N^\ell}}{\delta_0} \right)^{1/d^*} \\
&\leq \left( \frac{\epsilon^2(1-\gamma)^2}{\gamma^{2d^*}} \sum_{s \in \mathcal{S}^{\epsilon, *}} p(s) \left( \ln d^* + \sum_{\ell=1}^{\ell(s)} N^\ell \ln K_\ell \right) \right)^{1/d^*} \\
&\leq \left( \frac{6}{\gamma^2} \sum_{s \in \mathcal{S}^{\epsilon, *}} p(s) \left( \ln d^* + \sum_{\ell=1}^{\ell(s)} N^\ell \ln K_\ell \right) \right)^{1/d^*} \quad (\text{by 9))}
\end{aligned}$$

one obtains the bound

$$\begin{aligned}
\sum_{s \in \mathcal{S}^{\epsilon, *}} 2p(s)m(\ell(s), \delta_{\ell(s)}) &= \frac{\ln(1/\delta_0)}{(1-\gamma)^2 \epsilon^2} \sum_{s \in \mathcal{S}^{\epsilon, *}} \frac{\epsilon^2(1-\gamma)^2}{\ln(1/\delta_0)} 2p(s)m(\ell(s), \delta_{\ell(s)}) \\
&= \frac{\ln(1/\delta_0)}{\epsilon^2(1-\gamma)^2} \cdot \kappa_1^{d^*} \\
&= \frac{\ln(1/\delta_0)}{\epsilon^2(1-\gamma)^2} \cdot \kappa_1 \frac{\ln((1-\gamma)\epsilon) - \ln 6}{\ln \gamma} \\
&= \left( \ln \frac{1}{\delta_0} \right) \cdot \kappa_1 \frac{\ln 6}{\ln(1/\gamma)} \cdot \left( \frac{1}{(1-\gamma)\epsilon} \right)^{2 + \frac{\ln \kappa_1}{\ln(1/\gamma)}}
\end{aligned}$$

Similarly, defining

$$\begin{aligned}
\kappa_2 &:= \kappa_2(\epsilon, \delta_0, \gamma) := \left( \frac{\epsilon^2(1-\gamma)^2}{\ln(1/\delta_0)} \sum_{s \in \mathcal{S}^{\epsilon, *}} B(s) \sum_{d=d(s)} \prod_{\ell=d(s)}^{\ell(s)} K_\ell \right)^{1/d^*} \\
&= \left( \frac{\epsilon^2(1-\gamma)^2}{\ln(1/\delta_0)} \cdot |\mathcal{S}^\epsilon \setminus \mathcal{S}^{\epsilon, *}| \cdot 6 \cdot \ln \left( \frac{2|\mathcal{S}^{\epsilon, *}|}{\delta_0} \right) \right)^{1/d^*} \quad (\text{by Claim 9))} \\
&\leq \left( \epsilon^2(1-\gamma)^2 \cdot |\mathcal{S}^\epsilon \setminus \mathcal{S}^{\epsilon, *}| \cdot 6 \cdot \ln(2|\mathcal{S}^{\epsilon, *}|) \right)^{1/d^*} \\
&\leq \left( 6\gamma^{2d^*-2} \cdot |\mathcal{S}^\epsilon \setminus \mathcal{S}^{\epsilon, *}| \cdot 6 \cdot \ln(2|\mathcal{S}^{\epsilon, *}|) \right)^{1/d^*} \quad (\text{by 9))}
\end{aligned}$$

one obtains the bound

$$\sum_{s \in \mathcal{S}^{\epsilon, *}} B(s) \sum_{d=d(s)} \prod_{\ell=d(s)}^{\ell(s)} K_\ell = \frac{\ln(1/\delta_0)}{\epsilon^2(1-\gamma)^2} \cdot \kappa_2 \frac{\ln((1-\gamma)\epsilon) - \ln 6}{\ln \gamma} = \left( \ln \frac{1}{\delta_0} \right) \cdot \kappa_2 \frac{\ln 6}{\ln(1/\gamma)} \cdot \left( \frac{1}{(1-\gamma)\epsilon} \right)^{2 + \frac{\ln \kappa_2}{\ln(1/\gamma)}}$$

Finally, defining  $\kappa := \limsup_{\epsilon \rightarrow 0} \max(\kappa_1, \kappa_2)$ , one obtains the following sample complexity bound.

**Theorem 10.** *Sample complexity (6) is upper bounded by  $\left( \ln \frac{1}{\delta_0} \right) \cdot C(\kappa, \gamma) \cdot \left( \frac{1}{(1-\gamma)\epsilon} \right)^{2 + \frac{\ln \kappa}{\ln(1/\gamma)}}$ , where  $C(\kappa, \gamma) := 2\kappa \frac{\ln 6}{\ln(1/\gamma)}$ .*

## E.1 Worst case

When  $K_\ell = K > 1$  for each  $\ell > 0$  then  $\sum_{s \in \mathcal{S}^{\epsilon, *}} p(s) = \sum_{s \in \mathcal{S}^\epsilon} p(s) \leq K^{d^*}$ , and so

$$\kappa_1 \leq \left( \frac{6(\ln d^* + N^{d^*} d^* \ln K)}{\gamma^2} \sum_{s \in \mathcal{S}^{\epsilon, *}} p(s) \right)^{1/d^*} \leq \left( \frac{6(\ln d^* + N^{d^*} d^* \ln K) K^{d^*}}{\gamma^2} \right)^{1/d^*}.$$

Therefore,  $\limsup_{\epsilon \rightarrow 0} \kappa_1 \leq KN$ . Similarly, noting that  $|\mathcal{S}^\epsilon| \leq (NK)^{d^*}$ ,

$$\kappa_2 \leq \left( \gamma^{2d^*-2} \cdot (NK)^{d^*} \cdot 6 \cdot d^* \ln(NK) \right)^{1/d^*}$$

implying  $\limsup_{\epsilon \rightarrow 0} \kappa_2 \leq \gamma^2 KN$ .

## E.2 Case $K_0 > 1, K_\ell = 1$ for all $\ell \geq 1$

In this case

$$\sum_{s \in \mathcal{S}^{\epsilon, *}} p(s) \leq d^* K, \quad (10)$$

and so

$$\kappa_1 \leq \left( \frac{6}{\gamma^2} \sum_{s \in \mathcal{S}^{\epsilon, *}} p(s) (\ln d^* + N \ln K) \right)^{1/d^*} = \left( \frac{6}{\gamma^2} (\ln d^* + N \ln K) d^* K \right)^{1/d^*}$$

implying  $\limsup_{\epsilon \rightarrow 0} \kappa_1 \leq 1$ .

To bound  $\kappa_2$  note that  $p^\circ(s) \leq p(s)$  for all  $s$  and that  $\sum_{d=1}^{d^*} \prod_{\ell=d}^{d^*} K_\ell = 1$ , which imply

$$\begin{aligned} \kappa_2 &\leq \left( \frac{\epsilon^2 (1-\gamma)^2}{\ln(1/\delta_0)} \sum_{s \in \mathcal{S}^\epsilon} \left( 2p(s) m(\ell(s), \delta_{\ell(s)}) + 6 \ln\left(\frac{2N^\epsilon}{\delta_0}\right) \right) \right)^{1/d^*} \\ &\leq \left( \kappa_1^{d^*} + \frac{\epsilon^2 (1-\gamma)^2}{\ln(1/\delta_0)} \cdot |\mathcal{S}^{\epsilon, *}| \cdot 6 \ln\left(\frac{2N^\epsilon}{\delta_0}\right) \right)^{1/d^*} \end{aligned}$$

By (10) and the definition of  $\mathcal{S}^{\epsilon, *}$ , the restriction that  $K_\ell = 1$  for all  $\ell > 1$  imply

$$|\mathcal{S}^{\epsilon, *}| \leq K \cdot d^* \frac{3m(d^*, \delta_{d^*})}{8 \ln(2N^\epsilon/\delta_0)} \leq K \cdot d^* \frac{3N \ln(d^* K/\delta_0)}{16\gamma^{2d^*} \ln(1/\delta_0)}$$

Therefore, recalling also (9),

$$\begin{aligned} \kappa_2 &\leq \left( \kappa_1 + \frac{\gamma^{2d^*-2}}{\ln(1/\delta_0)} K \cdot d^* \frac{3N \ln(d^* K/\delta_0)}{16\gamma^{2d^*} \ln(1/\delta_0)} 6d^* \ln\left(\frac{KN}{\delta_0}\right) \right)^{1/d^*} \\ &= \left( \kappa_1 + \frac{(d^*)^2 2NK \ln(d^* K/\delta_0) \ln(KN/\delta_0)}{\gamma^2 \ln^2(1/\delta_0)} \right)^{1/d^*} \end{aligned}$$

Consequently,  $\limsup_{\epsilon \rightarrow 1} \kappa_2 \leq 1$  as well.

## E.3 Bandit case

Again  $K_0 > 1, K_\ell = 1$  for all  $\ell \geq 1$ , but it is also assumed that  $N = 1$  and all the rewards in one branch are the same (they can be different though in different branches). Then, directly from (6), one easily deduces the bound  $O\left(\left(\ln \frac{d^*}{\delta_0}\right) \sum_{u \neq u^*} \left(\frac{1}{(1-\gamma)(v^* - v(u) + \epsilon)}\right)^{-2}\right)$ .

## E.4 Deterministic MDPs

In case  $N = 1$  and  $K_\ell = K > 1$  for  $\ell \geq 0$ , then  $\kappa_1 \leq \left(\frac{6}{\gamma^2} \cdot K^{d^*} \cdot (\ln d^* + d^* \cdot \ln K)\right)^{1/d^*}$ , so  $\limsup_{\epsilon \rightarrow 0} \kappa_1 \leq K$ . Additionally,  $\kappa_2 = 0$ , since in this case  $p(s) = 1$  for each node  $s$ .

Assume now some structure in the rewards: for every action  $u$  on exactly one path in  $\mathbf{\Pi}^\infty$  the rewards are 1; everywhere else they are 0. Then nodes with depth at least  $\log(5)/\log(1/\gamma)$  bigger than their lowest ancestor having nonzero reward do not appear in  $\mathcal{S}^\epsilon$ . Therefore,

$$\begin{aligned} \kappa_1 &\leq \left( \frac{\epsilon^2 (1-\gamma)^2}{\ln(1/\delta_0)} \cdot K \cdot \sum_{d=1}^{d^*} K^{1+\log(5)/\log(1/\gamma)} m(d, \delta_d) \right)^{1/d^*} \\ &\leq \left( \frac{3}{\gamma^2 \ln(1/\delta_0)} \cdot d^* K^{1+\log(5)/\log(1/\gamma)} \ln \frac{d^* K^{d^*}}{\delta_0} \right)^{1/d^*}, \end{aligned}$$

and so  $\limsup_{\epsilon \rightarrow 0} \kappa_1 = 1$ .

## F Efficient version of StOP

This section is devoted to fix all the time-efficiency issues in the previous version of the algorithm. The primary task here is to find a way to solve both the policy evaluation and the construction of the optimistic policies efficiently.

With some abuse of notation let  $\text{Active}_t$  denote the set in round  $t$  consisting of policies  $\Pi$  for which rollout  $\tau(\Pi, \mathcal{T}_i)$  has length  $d(\Pi)$  for  $1 \leq i \leq m(d(\Pi), \delta_{d(\Pi)})$ , and, at the same time, for some child policy  $\Pi'$  of  $\Pi$  some rollout  $\tau(\Pi', \mathcal{T}_i)$  for  $1 \leq i \leq m(d(\Pi), \delta_{d(\Pi)})$  has length less than  $d(\Pi')$ .

### F.1 Evaluating the children of $\Pi_t$

The first problem to solve is to maintain the sample trees without actually going through all the children policies of  $\Pi_t$ .

To this end, define first  $m_d(s)$  as the number of times  $s$  appears in sample trees  $\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_{m(d, \delta_d)}$  at the current round. Similarly, let  $\hat{r}_d(s)$  denote the average of the rewards for  $s$  in  $\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_{m(d, \delta_d)}$  at the current round. These values are easily updated using a simple recursion rule applied in algorithm `Sample-eff`.

**Claim 11.** *Executing `Sample-eff`( $\Pi, s, m$ ) ensures that  $\tau(\mathcal{T}_i, \Pi)$  has length  $d(\Pi)$  (i.e., has full length) for  $i = 1, 2, \dots, m$ , and runs in time  $O(m \cdot d(\Pi))$ .*

As the next step, note that, if the first  $K_d$  children policies of  $\Pi_t$  that `StOP` picks to evaluate in round  $t$  (where  $d = d(\Pi_t)$ ) do not share any leaves, then `BoundValue` will not call `SampleTransition` or `SampleReward` for any other children of  $\Pi_t$ . The reason for this is that the first  $K_d$  trees include all the nodes that appear in any child policy of  $\Pi_t$ .

The above argument shows that the evaluation of a policy  $\Pi$  in `StOP-eff` and in `StOP` are essentially equivalent.

### F.2 Constructing the optimistic policies

Note that, in round  $t$ , for any  $\Pi \in \cup_{t' \leq t} \text{Active}$  it holds that

$$\hat{v}(\Pi) = \sum_{s \in \Pi} \gamma^{d(s)} \cdot m_d(s) \cdot \hat{r}_d(s) .$$

Additionally, as  $b(\Pi) = \hat{v}(\Pi) + 2\frac{\gamma^{d(\Pi)}}{1-\gamma}$ , it holds for any two policies  $\Pi$  and  $\Pi'$  of the same depth that

$$b(\Pi) > b(\Pi') \Leftrightarrow \hat{v}(\Pi) > \hat{v}(\Pi') .$$

It is thus easy to compute the value of any active policy, and also to decide between two policies which one is better. However, it is less obvious how to construct the optimistic policies efficiently.

**Theorem 12.** *For any action  $u$  accessible from  $x_0$ , and any round  $t$ , `ValueTr`( $s_u$ ) returns  $\Pi_{t,u}^\dagger$ , where  $s_u$  is the child of the root labeled  $u$ .*

*Proof.* Let  $a_d(s) = a_{d,t}(s)$  be the indicator that, for some  $1 \leq i \leq m(d, \delta_d)$ , sample tree  $\mathcal{T}_i$  has a leaf below  $s$  with  $d(s) = d$  at iteration  $t$ . Note that for action node  $s$ ,  $a_d(s)$  must be set to 1 if  $m_{d(s)}(s) > 0$  and  $m_{d(s)+1}(s') = 0$  for some child  $s'$  of  $s$ , otherwise it must be set to 0. For node  $s$  of depth  $d(s) < d$ ,  $a_{d,t}(s)$  can be computed based on the simple recursion rule  $a_d(s) := \max_{s' \text{ child of } s} a_d(s')$ .

Equivalently,  $a_{d,t}(s)$  indicates that, for some policy  $\Pi$  of depth  $d$  containing  $s$ , rollout  $\tau(\mathcal{T}_i, \Pi)$  has length  $d$  (i.e., full length) for  $i = 1, \dots, m(d, \delta_d)$ , but for some child policy  $\Pi'$  of  $\Pi$  and for some  $1 \leq i \leq m(d, \delta_d)$  rollout  $\tau(\Pi', \mathcal{T}_i)$  goes through  $s$  and has length at most  $d$  (instead of  $d+1$ , which would be the maximal possible). On one hand, the extra requirement about the rollout going through  $s$  makes a distinction between  $a_{d,t}(s)$  and the indicator that  $s$  belongs to some policy in  $\text{Active}_t$ , but, at the same time, this is the distinction that makes it easy to compute it efficiently with the recursive rule described above. This is the key insight that is used in constructing the optimistic policies efficiently too.



Now, consider, for each node  $s$  the policies in  $\cup_{t' \leq t} \text{Active}_{t'}$  with  $d(\Pi) = d$ , and denote by  $\Pi_{t,d}^{\text{comp}}(s)$  the one that has the largest cumulated reward below  $s$  in the first  $m(d, \delta_d)$  sample trees. Denote this cumulated reward by  $\hat{v}_d^{\text{comp}}(s)$ , and note that it can be computed recursively by

- setting it to  $\hat{r}_d(s)$  for each action node  $s$  with  $d(s) = d$ ,
- setting it to  $\max_{s' \text{ children of } s} \hat{v}_d^{\text{comp}}(s')$  for all action nodes with  $d(s) < d$ , and
- setting it to  $\hat{v}_d^{\text{compl}}(s) := \gamma \sum_{s': \text{child of } s} (m_d(s') \cdot \hat{v}_d^{\text{compl}}(s'))$  for a transition node  $s$  with  $d(s) \leq d$ .

Finally, consider, for a node  $s$ , those policies in  $\cup_{t' \leq t} \text{Active}_{t'}$  which satisfy that

- $d(\Pi) = d$
- rollout  $\tau(\mathcal{T}_i, \Pi)$  has length  $d$  (i.e., full length) for  $i = 1, \dots, m(d, \delta_d)$ ,
- for some child policy  $\Pi'$  of  $\Pi$  and for some  $1 \leq i \leq m(d, \delta_d)$  rollout  $\tau(\Pi', \mathcal{T}_i)$  goes through  $s$  and has length  $d$  too (instead of  $d + 1$ ).

Denote by  $\Pi_{t,d}^{\text{inc}}(s)$  the one that has the largest cumulated reward below  $s$  in the first  $m(d, \delta_d)$  sample trees, and by  $\hat{v}_d^{\text{inc}}$  this cumulated reward. This value can also be computed efficiently using recursion:

- $\hat{v}_d^{\text{inc}}(s) := \hat{r}_d(s)$  for a transition node  $s$  with  $d(s) = d$
  - $\hat{v}_d^{\text{inc}}(s) := \max_{s' \text{ children of } s \text{ with } a_d(s)=1} \hat{v}_d^{\text{inc}}(s')$  for a transition node  $s$  with  $d(s) < d$ , and
  -
- $$\hat{v}_d^{\text{inc}}(s) := \gamma \max_{s': \text{child of } s \text{ with } a_d(s')=1} \left( m_d(s') \cdot \hat{v}_d^{\text{inc}}(s') + \sum_{s'' \neq s' \text{ child of } s} (m_d(s'') \cdot \hat{v}_d^{\text{compl}}(s'')) \right)$$
- for an action node  $s$  with  $d(s) \leq d$

The claim of the theorem follows by noting that, for any child node  $s$  of the root,  $\Pi_{t,d}^{\text{inc}}(s) = \Pi_{t,u}^{\dagger}$ , where  $u$  is the label of  $s$ .  $\square$

In order to simplify the pseudocode, the construction of the optimistic policies are not implemented. Nevertheless, they can be easily obtained along the same line values  $\hat{v}_d^{\text{comp}}(s)$  and  $\hat{v}_d^{\text{inc}}(s)$  are computed.

Finally note that in step  $t$  only the values belonging to the nodes of policy  $\Pi_t$  require update. Making use of this, an even more significant speed-up is possible.

---

**Algorithm 4** StOP-eff( $s_0, \delta_0, \epsilon, \gamma$ )

---

```
1: for all  $u$  available from  $x_0$  do ▷ Initialize
2:    $\Pi :=$  smallest policy with the child  $s_u$  of  $s_0$  labeled  $u$ 
3:    $\delta_1 := (\delta_0/d^*) \cdot (K_0)^{-1}$  ▷  $d(\Pi) = 1$ 
4:   Sample( $\Pi, s_u, m(1, \delta_1)$ )
5:  $t := 1$ 
6: for round  $t = 1, 2, \dots$  do
7:   for all  $u$  available at  $x_0$  do
8:     ValueTr( $s_u$ )
9:      $\Pi_{t,u}^\dagger := \operatorname{argmax}_{\Pi \in \text{Active}(u)} b(\Pi)$ 
10:     $\Pi_t^\dagger := \Pi_{t,u_t^\dagger}^\dagger$ , where  $u_t^\dagger := \operatorname{argmax}_u b(\Pi_{t,u}^\dagger)$  ▷ optimistic policy and action
11:     $\Pi_t^{\dagger\dagger} := \Pi_{t,u_t^{\dagger\dagger}}^\dagger$ , where  $u_t^{\dagger\dagger} := \operatorname{argmax}_{u \neq u_t^\dagger} b(\Pi_{t,u}^\dagger)$  ▷ secondary policy and action
12:    if  $\nu(\Pi_t^\dagger) + \epsilon \geq \max_{u \neq u_t^\dagger} b(\Pi_{t,u}^\dagger)$  then ▷ termination criterion
13:      return  $u_t^\dagger$ 
14:    if  $d(\Pi_t^{\dagger\dagger}) \geq d(\Pi_t^\dagger)$  then ▷ choose action and policy to explore
15:       $u_t := u_t^\dagger$  and  $\Pi_t := \Pi_t^\dagger$ 
16:    else
17:       $u_t := u_t^{\dagger\dagger}$  and  $\Pi_t := \Pi_t^{\dagger\dagger}$ 
18:    set  $d_t := d(\Pi_t)$ 
19:     $\delta := (\delta_0/d^*) \cdot \prod_{\ell=0}^{d_t-1} (K_\ell)^{-N^\ell}$  ▷ the # of policies of depth at most  $d$  is  $\prod_{\ell=0}^{d-1} (K_\ell)^{N^\ell}$ 
20:    for each of the  $K_{d_t}$  action  $u$  do
21:      let  $\Pi_{t,u}$  be the policy children of  $\Pi$  that follows action  $u$  from each leaf of  $\Pi$ 
22:      set  $a_{d_t}(s) := 1$  for each node  $s$  of  $\Pi_{t,u}$  that are not in  $\Pi_t$ 
23:      Sample( $\Pi_{t,i}, s_{u_t}, m(d_t + 1, \delta_{d_t+1})$ )
24:     $t := t + 1$ 
```

---

---

**Algorithm 5** Sample-eff( $\Pi, s, m$ )

---

```
1: if  $s$  is a leaf of  $\Pi$  then return
2: let  $s'$  be the child node of  $s$  in  $\Pi$ 
3: while  $m_{d(\Pi)}(s') < m$  ▷ make sure that  $s$  has at least  $m$  samples do
4:    $m_{d(\Pi)}(s') := m_{d(\Pi)}(s') + 1$ 
5:    $s'' := \text{SampleTransition}(s')$ 
6:    $\hat{r}_{d(\Pi)}(s'') := \frac{\hat{r}_{d(\Pi)}(s'') \cdot m_{d(\Pi)}(s'') + \text{SampleReward}(s'')}{1 + m_{d(\Pi)}(s'')}$ 
7:    $m_{d(\Pi)}(s'') := m_{d(\Pi)}(s'') + 1$ 
8: for all grandchildren  $s''$  of  $s$  do ▷ ensure that all rollouts going through  $s$  have full length in  $\Pi$ 
9:   Sample-eff( $\Pi, s'', m_{d(\Pi)}(s'')$ )
```

---

---

**Algorithm 6** ValueTr( $s$ )

---

```
1:  $a_d(s) = 0$ 
2: for all children  $s'$  of  $s$  with  $\max_{d=d(s'), \dots, d^*} m_d(s') > 0$  do
3:   ValueAc( $s'$ )
4: for all  $d := d(s) + 1, \dots, d^*$  with  $m_d(s) > 0$  do
5:    $\hat{v}_d^{\text{compl}}(s) := \gamma \sum_{s': \text{child of } s} (m_d(s') \cdot \hat{v}_d^{\text{compl}}(s'))$ 
6:    $a_d(s) := \max_{s' \text{ child of } s} a_d(s')$ 
7:    $\hat{v}_d^{\text{inc}}(s) := \gamma \max_{s': \text{child of } s \text{ with } a_d(s')=1} (m_d(s') \cdot \hat{v}_d^{\text{inc}}(s'))$ 
8:    $+ \sum_{s'' \neq s' \text{ child of } s} (m_d(s'') \cdot \hat{v}_d^{\text{compl}}(s''))$ 
```

---

---

**Algorithm 7** ValueAc( $s$ )

---

- 1: **for all** children  $s'$  of  $s$  **do**
- 2:   ValueTr( $s'$ )
- 3:  $\hat{v}_{d(s)}^{\text{comp}}(s) := \hat{r}_{d(s)}(s)$
- 4: **if**  $m_{d(s)}(s) > 0$  **but**  $m_{d(s)+1}(s') = 0$  for some child  $s'$  of  $s$  **then**
- 5:    $a_{d(s)}(s) := 1$
- 6:    $\hat{v}_{d(s)}^{\text{inc}}(s) := \hat{r}_{d(s)}(s)$
- 7: **for**  $d := d(s) + 1, \dots, d^*$  **do**
- 8:    $\hat{v}_d^{\text{comp}}(s) := \max_{s' \text{ children of } s} \hat{v}_d^{\text{comp}}(s')$
- 9:    $a_d(s) := \max_{s' \text{ children of } s} a_d(s')$
- 10:  $\hat{v}_d^{\text{inc}}(s) := \max_{s' \text{ children of } s \text{ with } a_d(s)=1} \hat{v}_d^{\text{inc}}(s')$

---