
Scale-free adaptive planning for deterministic dynamics & discounted rewards

Peter L. Bartlett¹ Victor Gabillon² Jennifer Healey³ Michal Valko⁴

Abstract

We address the problem of planning in an environment with deterministic dynamics and stochastic discounted rewards under a limited numerical budget where the ranges of both rewards and noise are unknown. We introduce **PlaT γ POOS**, an adaptive, robust, and efficient alternative to the OLOP (open-loop optimistic planning) algorithm. Whereas OLOP requires a priori knowledge of the ranges of both rewards and noise, **PlaT γ POOS** dynamically adapts its behavior to both. This allows **PlaT γ POOS** to be immune to two vulnerabilities of OLOP: failure when given underestimated ranges of noise and rewards and inefficiency when these are overestimated. **PlaT γ POOS** additionally adapts to the global smoothness of the value function. **PlaT γ POOS** acts in a provably more efficient manner vs. OLOP when OLOP is given an overestimated reward and show that in the case of no noise, **PlaT γ POOS** learns exponentially faster.

1. Introduction

We consider the problem of planning in a general *stochastic environment with deterministic dynamics and discounted rewards*. Our goal is to recommend the best first action for an agent to take from a given state. We envision that the discount factor γ is known and that our learner has a limited allocation of n interactions to spend querying a generative model of the environment. The objective is to maximize the sum of discounted rewards of the best sequence of actions following from the recommended first action. This is equivalent to minimizing the *simple regret*. We introduce the algorithm **PlaT γ POOS**, Planning with γ Plus an Online Optimization Strategy, as a robust and efficient scale-free alternative to the OLOP algorithm (open-loop optimistic planning, Bubeck & Munos, 2010; Leurent & Maillard, 2019) for this

setting. Our algorithm implements a scale-free function optimization strategy similar to Sequ00L (Bartlett et al., 2019) rather than an upper-confidence-bound approach which allows our algorithm to efficiently adapt to the problem space *without prior knowledge of the ranges of the noise or the rewards*.

Planning in a stochastic environment is an important setting often modeled by Markov decision processes (MDPs, Puterman, 1994; Bertsekas & Tsitsiklis, 1996). One approach to solving these settings is to find the optimal policy that maximizes the expected sum of rewards and then generate an action recommendation according to that optimal policy. Unfortunately, in most practical settings where we are limited by computational resources, finding this optimal policy is often not possible, especially when the state space becomes large. Therefore, instead of trying to estimate the optimal policy of the MDP, we focus only on finding the *best first action* given our budget. We evaluate the performance of the recommendation in terms of the simple regret, the difference in reward between choosing the optimal first action vs. choosing our recommended first action and then in both cases choosing an optimal sequence of actions following the first action. This metric is often used to evaluate planning strategies that optimize numeric budgets (Bubeck & Munos, 2010; Buşoniu & Munos, 2012; Grill et al., 2016), in contrast to the cumulative regret where we are penalized during the search for querying sub-optimal actions. Once the agent takes the first action and moves to the next state, our evaluation can be repeated with a new budget allocation and the following best first action can be recommended. This allows us to approximately following an optimal policy, action by action, in an online way. Previously, there have been several strategies proposed on how to efficiently allocate a numeric budget to search for an optimal value in a stochastic space. Many of these have been successfully implemented using methods based on upper confidence bounds (UCBs) such as UCT (Upper Confidence Trees, Kocsis & Szepesvári, 2006). This approach has been proven to be very efficient in practice (Coulom, 2007; Gelly et al., 2006; Silver et al., 2016), however, UCT can badly misbehave on some problems (Coquelin & Munos, 2007) and more theoretically sound approaches have been proposed (Hren & Munos, 2008; Bubeck & Munos, 2010; Buşoniu & Munos, 2012; Feldman & Domshlak, 2014; Szörényi et al., 2014;

¹University of California, Berkeley, USA ²Noah's Ark Lab, Huawei Technologies, London, UK ³Adobe Research, San Jose, USA ⁴SequeL team, INRIA Lille - Nord Europe, France. Correspondence to: Victor Gabillon <victor.gabillon@huawei.com>.

Kaufmann & Koolen, 2017; Shah et al., 2019). Some of these methods are connected to the ones from function optimization (Bubeck et al., 2011; Munos, 2011; Valko et al., 2013) as shown by Munos (2014), however, one key difference is that in planning, as opposed to function optimization, the structure of the reward is a discounted reward, specifically a *sum* of rewards and the *discount factor* γ . This reward structure influences the behavior of the optimizers (Bubeck & Munos, 2010), for example, the discount factor brings smoothness of the value function which in turn makes it easier to optimize. **PlaT γ P00S** exploits the effect of the discount factor to efficiently manage an adaptive planning strategy in the face of unknown ranges of noise and rewards.

This adaptive strategy of **PlaT γ P00S** can make it more robust and efficient in practice than other planning strategies. For example, even though they are theoretically sound, the empirical performance UCB-based approaches depend on the careful tuning of the upper confidence bound. If the upper confidence bound is too large then the UCB-based learner plays very conservatively by overestimating sub-optimal options for many rounds. Moreover, these UCBs might depend on instance parameters that are simply not known such as the *range of the rewards* and the *range of the noise*. We build on the function optimization approach of Bartlett et al. (2019) that does not use UCBs and obtains improved results over the state-of-the-art by adapting to the problem difficulty with a scale-free approach. **PlaT γ P00S** adapts this scale-free optimization approach to planning. This *scale-free* property becomes a more desired feature as machine learning gets closer to applications, whether it is online (Ross et al., 2013; Orabona & Pál, 2018) or deep learning (Orabona & Tommasi, 2017), since in the reality, many parameters are never known.

In terms of planning strategy, the **PlaT γ P00S** algorithm is an adaptive, robust, and efficient alternative to OLOP. Whereas OLOP *requires the knowledge of where the ranges of both rewards and noise*, **PlaT γ P00S** dynamically adapts its behavior to the both ranges, as well as some potential additional global smoothness of the value function. Our algorithm’s ability to adapt allows it to avoid failure in cases where the ranges of noise and rewards are underestimated and to act more efficiently in cases where they are overestimated. **PlaT γ P00S** recovers the results of OLOP while allowing improvements in various classes of problems.

Our contributions We show that **PlaT γ P00S**:

- adapts its behavior to an unknown range of rewards,
- requires no apriori assumptions or knowledge on noise,
- empirically learns much faster than UCB approaches,
- gets the fast rate of deterministic planning in low noise for all regime; in particular, it learns exponentially faster than OLOP when there happens to be no noise,
- adapts also the global smoothness ρ and ν beyond the base smoothness provided by γ .

We additionally address a realistic constraint where the agent can only *reset* to the original state and not to any state it wishes. Our results hold for MDPs with deterministic dynamics and can equally be applied to open loop planning problems (as discussed by Munos 2014) where we search for the best sequence of actions, ignoring the actual states that are reached after each action (Bubeck & Munos, 2010).

Related algorithms, where the objective is to find the value of the state rather than to identify the best action, include TrailBlazer (Grill et al., 2016) and StOP (Szörényi et al., 2014). A key difference is that these algorithms are *fixed confidence* and output a value using a small number of samples given an accuracy/probability, whereas our algorithm does exploration under a fixed budget of samples and guarantees *how good* the found action is. Even for simple multi-arm bandits, these two problems have different complexity (Carpentier & Locatelli, 2016) and can only be equivalent under unrealistic side knowledge (Gabillon et al., 2012). These related algorithms are also impractical for our setting. TrailBlazer uses confidence bounds that are humongous and StOP takes exponential time. Similar to OLOP, both also need to know noise and reward ranges.

2. Background

We model our problem with an MDP with state space X , action space A and dynamics such that taking the chosen action a_t at time t deterministically transitions the system from $x_t \in X$ to state $x_{t+1} = f(x_t, a_t)$ generating a reward $r_t = r(x_t, a_t) + \varepsilon_t$, with ε_t being the noise. We consider:

deterministic rewards The evaluations are noiseless, that is for all t , $\varepsilon_t \triangleq 0$ and $r_t = r(x_t, a_t)$.

stochastic rewards The evaluations are perturbed by a noise of range $b \in \mathbb{R}_+$: At any round, ε_t is a random variable, independent from noise at previous rounds,

$$\mathbb{E}[r_t|x_t] \triangleq r(x_t, a_t) \text{ and } |r_t - r(x_t, a_t)| \leq b. \quad (1)$$

We assume that all rewards lie in the interval $[0, R_{\max}]$ and while the state space may be large and possibly infinite, that the action space is finite, with K available actions. We treat an infinite time-horizon problem with discounted rewards where the discount factor ($0 \leq \gamma < 1$) is known. For any possible policy $\pi : X \rightarrow A$, we define the value function $V^\pi : X \rightarrow \mathbb{R}$ associated to π as $V^\pi(x) \triangleq \sum \gamma^t r(x_t, \pi(x_t))$, where x_t is the state of the system at time t when starting from x (i.e., $x_0 = x$) and following policy π . In the next definition, we also define the Q -value function $Q^\pi : X \times A \rightarrow \mathbb{R}$ associated to policy π , for each state-action pair (x, a) , as the value of playing action a in state x and the following π thereafter.

Definition 1. The Q -value function Q^π of policy π is

$$Q^\pi(x, a) \triangleq r(x, a) + \gamma V^\pi(f(x, a)).$$

Notice that $V^\pi(x) = Q^\pi(x, \pi(x))$. We define the *optimal* value function, and Q -value function respectively, as $V^*(x) \triangleq \sup_{\pi} V^\pi(x)$ and $Q^*(x, a) \triangleq \sup_{\pi} Q^\pi(x, a)$, which corresponds to playing a first and optimally after. From the dynamic programming, we have the Bellman equations (Bertsekas & Tsitsiklis, 1996; Puterman, 1994),

$$\begin{aligned} V^*(x) &= \max_{a \in A} r(x, a) + \gamma V^*(f(x, a)), \\ Q^*(x, a) &= r(x, a) + \gamma \max_{b \in A} Q^*(f(x, a), b). \end{aligned}$$

Let $[a : c] = \{a, a+1, \dots, c\}$ with $a, c \in \mathbb{N}$, $a \leq c$, and $[a] = [1 : a]$ and let \log_d be the logarithm in base d , $d \in \mathbb{R}$ and \log without a subscript be the natural logarithm.

2.1. Optimistic planning under finite numerical budget

We assume that we have a generative model of f and r that generates simulated transitions and rewards. We want to make the best possible use of this model in order to recommend a best next action $a(n)$ such that the sum of the rewards resulting from playing $a(n)$ and then optimally afterwards is as close as possible to playing optimally from the beginning. For that purpose, we define the performance loss r_n as

$$r_n \triangleq \max_{a \in A} Q^*(x, a) - Q^*(x, a(n)).$$

2.2. The planning tree

For a given initial state x , consider the (infinite) planning tree defined by all possible sequences of actions (thus all possible reachable states starting from x). Let A^∞ be the set of infinite sequences (a_0, a_1, a_2, \dots) where $a_t \in A$. The branching factor of this tree is the number of actions $|A| = K$. Since the dynamics are deterministic, to each finite sequence $a \in A^d$ of length d we assign a state that is reachable starting from x by following this sequence of d actions. Using standard notation for alphabets, we write $A^0 = \{\emptyset\}$ and A^\bullet for the set of finite sequences. For $a \in A^\bullet$ we let $h(a)$ be the length of a , and $aA^h = \{aa', a' \in A^h\}$, where aa' denotes the sequence a followed by a' . We identify the set of finite sequences $a \in A^\bullet$ with the set of nodes of the tree. With $h' \leq h$ and $a \in A^h$ we denote $a_{[h']}$ the sequence of action composed of the h' first actions from a , i.e., $\{a_0, \dots, a_{h'-1}\}$. We fix $a_{[0]} \triangleq \emptyset$. The value $v(a)$ of an infinite sequence $a \in A^\infty$ is the discounted sum of rewards along the trajectory starting from the initial state x and defined by the choice of this sequence of actions,

$$v(a) \triangleq \sum_{t \geq 0} \gamma^t r(x_t, a_t), \text{ where } x_0 = x; x_{t+1} = f(x_t, a_t).$$

Now, for any finite sequence $a \in A^\bullet$, or node, we define the value $v(a) = \sup_{a' \in A^\infty} v(aa')$. We write $v^* = v(\emptyset) = \sup_{a \in A^\infty} v(a)$ for the optimal value at the initial state which

is the root of the tree, $v^* = V^*(x)$. We denote the set of optimal infinite sequence of action as A^* which contains any $a \in A^\infty$ such that $v(a) = v^*$. We note the set of optimal finite sequence of actions of depth h as $A^{*,h}$ which contains any $a \in A^h$ such that $v(a) = v^*$. We also define the u - and b -values for the lower- and upper- bounds on $v(a)$ as

$$u(a) \triangleq \sum_{t=0}^{h(a)-1} \gamma^t r(x_t, a_t), \text{ and } b(a) = u(a) + \frac{\gamma^{h(a)} R_{\max}}{1 - \gamma}.$$

Indeed, since all rewards are in $[0, R_{\max}]$, we trivially have that $u(a) \leq v(a) \leq b(a)$. At any finite time t an algorithm has opened a set of nodes, which defines the expanded tree \mathcal{T}_t . We say the learner *opens* (or *expands*) a node a with m evaluations if uses the generative model f and r to generate m transitions and rewards for the K children nodes aA . In the deterministic reward feedback, $m = 1$. The bounds reported in this paper are in terms of the total number of openings n , instead of evaluations. The number of function evaluations is upper bounded by Kn . $T_{x,a}$ denotes the total number of evaluations allocated to action $a \in A$ in state x . We define, especially for the noisy case, the estimated value of the reward $\hat{r}(x, a)$ of action $a \in A$ in state x . Given the $T_{x,a}$ evaluations $r_1, \dots, r_{T_{x,a}}$, $\hat{r}(x, a) \triangleq \frac{1}{T_{x,a}} \sum_{s=1}^{T_{x,a}} r_s$; the empirical average of rewards obtained at when performing action $a \in A$ in state x . To ease notation, for $a \in A^m$ and $h \leq m$, we write $T_a \triangleq \mathbb{E}[T_{x_{h(a)-1}, a_{h(a)-1}} | x_{t+1} \sim P(\cdot | x_t, a_t), x_0 = x]$ for the number of pulls to the last action in a . Similarly, $\hat{r}_h(a) \triangleq \mathbb{E}[\hat{r}(x_h, a_h) | x_{t+1} \sim P(\cdot | x_t, a_t), x_0 = x]$ and $r_h(a) \triangleq \mathbb{E}[r(x_h, a_h) | x_{t+1} \sim P(\cdot | x_t, a_t), x_0 = x]$. In the case of deterministic dynamics, x_h is such that $x_{t+1} \sim P(\cdot | x_t, a_t)$ and $x_0 = x$ is a fixed state from which we can sample from if we have a full access to the generative model. Hence, for a finite sequence $a \in A^\bullet$ or node,

$$\hat{u}(a) \triangleq \sum_{t=0}^{h(a)-1} \gamma^t \hat{r}(x_t, a_t) = \sum_{t=0}^{h(a)-1} \gamma^t \hat{r}_t(a).$$

We the existence of at least one $a^* \in A^\infty$ for which we have $V^*(x) = \sup_{a \in A^\infty} v(a)$ and define a smoothness for v .

Proposition 1. *There exists $\nu \in (0, R_{\max}/(1 - \gamma)]$ and $\rho \in (0, \gamma]$ such that $\forall h \geq 0, \forall a \in A^h, u(a) \geq v(a) - \nu \rho^h$.*

Note that this holds automatically for $\nu = R_{\max}/(1 - \gamma)$ and $\rho = \gamma$. For some problems with an extra regularity this might also hold for some $\nu < R_{\max}/(1 - \gamma)$ and $\rho < \gamma$. Note that our results automatically adapt to ρ without knowing its value. Note that while having a smoothness ρ means having rewards diminishing geometrically with depth with a ratio of ρ , the constant ν is linked to the scale of variation of the V and which can often be realistically smaller than $\nu < R_{\max}/(1 - \gamma)$. We now define a measure of the quantity of near-optimal sequences for the smoothness ν, ρ .

Definition 2. For any $\nu > 0$ and $\rho \in (0, 1)$, the **branching factor** $\kappa^v(\nu, \rho)$ and with associated constant C , is

$$\kappa^v(\nu, \rho) \triangleq \inf \{ \kappa \geq 1 : \exists C > 1, \forall h \geq 0, \mathcal{N}_h^v(3\nu\rho^h) \leq C\kappa^h \},$$

where $\mathcal{N}_h^v(\varepsilon)$ is the number of nodes $a \in A^h$ of depth h such that $v(a) \geq v^* - \varepsilon$.

We also define a related but different quantity than $\kappa^v(\nu, \rho)$. In particular, we define $\kappa^u(\nu, \rho)$ that uses $\mathcal{N}_h^u(\varepsilon)$ instead of $\mathcal{N}_h^v(\varepsilon)$ where $\mathcal{N}_h^u(\varepsilon)$ is the number of nodes $a \in A^h$ of depth h such that $u(a) \geq v^* - \varepsilon$. Our results use the new quantity $\kappa^u(\nu, \rho)$, recover the previous results using $\kappa^v(\nu, \gamma)$ in the method of Bubeck & Munos (2010) and go beyond them. Indeed, theirs are only formulated for $\rho = \gamma$ while we prove the following claim in Appendix A.

Proposition 2. $\kappa^u(\nu/2, \gamma) \leq \kappa^v(\nu, \gamma) \leq \kappa^u(2\nu, \gamma)$.

3. Optimization vs planning

Our **PlaTγPOOS** approach is a very close sibling to the flat optimization algorithm, **Stroqu00L** (Bartlett et al., 2019), however, we explain how the structure of the planning setting and the discount factor γ shaped our approach. The optimal application of an optimization algorithm to the planning setting is often not straightforward as discussed by Bubeck & Munos (2010). In their Section 2.2, Bubeck & Munos (2010) show that optimization can be applied to the planning problem either in a naïve way or a good way. The authors take as an example the uniform planning problem. The naïve and good strategies are evaluated by comparing the uncertainty $|u(a) - \hat{u}(a)|$ of their estimates $\hat{u}(a)$. Both strategies collect rewards identically, evaluating $u(a)$ for all the K^H nodes $a \in A^H$ at a fixed depth $H = h(a)$ by allocating one episode (of length H) for each a and receiving $r_{h,a}$, $1 \leq h \leq H$. In the naïve version, for all sequences a , the estimation of $u(a)$ uses only the H samples collected in the one episode related to a . Here $\hat{u}(a) = \sum_{h=0}^{h(a)-1} \gamma^h r_{h,a}$ and $T_{a[h]} = 1$. In contrast, the good planning strategy *reuses* estimates. For two distinct sequences a and a' , the good strategy re-uses any samples $r_{h,a}$ for the estimation of both $u(a)$ and $u(a')$ if $a[h] = a'[h]$. Here $\hat{u}(a) = \sum_{h=0}^{h(a)-1} \gamma^h \frac{1}{K^{H-h}} \sum_{a': a[h]=a'[h]} r_{h,a'}$ and $T_{a[h]} = K^{H-h}$. This comparison is used to demonstrate the advantage of the use of the cross-sequence information to concentrate the estimate of the mean reward associated with each action more efficiently. It is by using this type of cross-sequence information that OLOP is able to obtain an improved regret over a naïve application of H00 (Bubeck et al., 2011) to the planning problem.

The previous discussion on cross-sequence information is tied to the case of uniform exploration strategies. Good uniform strategies guarantee $T_{a[h]} = K^{H_u-h}$ by exploring

Input: n, A

Initialization: Open $t_{0,1}$. $h_{\max} = \lfloor n/\sqrt{\log(n)} \rfloor$.

For $h = 1$ to h_{\max}

Open $\lfloor h_{\max}/h \rfloor$ nodes $a_{h,i}$ of depth h
with largest values $u(a_{h,i})$.

Output $x(n) = \arg \max_{a_{h,i} \in \mathcal{T}} u(a_{h,i})$.

Figure 1. Algorithm for free planning with no reset condition

until a reasonably shallow depth H_u but sampling all K^{H_u} nodes and sharing cross sequence information. In the case of OLOP, H00 or more particularly **Stroqu00L**, only a subset S of size $|S| \ll K^{H_u}$ of the most promising nodes are explored but at deeper depth $H_s \gg H_u$. Therefore, obtaining for a given sub-sequence of actions a at depth $h < H_s$ a lower bound on the number of sequence of actions at depth H_s that contains a is complex in general. Actually one can design a problem with two actions that would drastically limit the amount of cross-sequence information for the optimal sequence of actions a^* in **Stroqu00L**. Therefore applying **Stroqu00L** with information sharing might not be enough and we chose to ensure algorithmically that a node at depth $h+1$ will be pulled γ^2 times less than a node at depth h . Indeed, using Chernoff-Hoeffding inequalities, we have $|u(a) - \hat{u}(a)| \leq \sum_{h=0}^{h(a)-1} \gamma^h / \sqrt{T_{a[h]}}$. However, **PlaTγPOOS**, through some simple reparametrization, remains very close to **Stroqu00L** and we leave as an open question whether equivalent theoretical guarantees could be proved directly for **Stroqu00L** applied to planning.

4. Deterministic dynamics and rewards

In this section, we consider a simpler case of deterministic rewards in order to introduce our new ideas. The evaluations are noiseless, that is $\forall t, \varepsilon_t = 0$ and $r_t = r(x_t, a_t)$.

In Figure 1, we provide the **Sequ00L** (Bartlett et al., 2019) algorithm applied to planning. In this case it is straightforward to follow the analysis of **Sequ00L** in order to obtain the same rates, up to logarithmic factors, of simple regret as the state of the art algorithm OPD for the doubly deterministic case (Hren & Munos, 2008; Munos, 2014). and get the result¹ of Theorem 1. This direct usage was already discussed by Munos (2014, Section 5.1). Using **Sequ00L** for planning already permits to have an algorithm that does not use the parameter R_{\max} and that can adapt to extra smoothness in the value function ν, ρ as discussed Section 2. Note that to obtain similar adaptations to R_{\max}, ν and ρ , one could have already used S00 (Munos, 2014). However S00 *does not* come with optimal simple regret (Bartlett et al., 2019).

Theorem 1. For any planning problem with associated (ν, ρ) , and branching factor $\kappa = \kappa^u(\nu, \rho)$, we have, after n

¹ $\sqrt{\log(n)}$ is the n -th harmonic number

Input: n, A
Initialization: Open $t_{0,1}$. $h_{\max} = \lfloor \frac{n}{\log n} \rfloor$.
For $h = 1$ to h_{\max}
 Open $\lfloor h_{\max}/h^2 \rfloor$ nodes $a_{h,i}$ of depth h
 with largest values $u(a_{h,i})$.
Output $x(n) = \arg \max_{a_{h,i} \in \mathcal{T}} u(a_{h,i})$.

Figure 2. Algorithm for constraint planning (restart)

rounds, the simple regret of Sequ00L bounded as:

- If $\kappa = 1$, $r_n \leq \nu \rho^{\frac{1}{C}} \lfloor \frac{n}{\log n} \rfloor$.
- If $\kappa > 1$, $r_n = \mathcal{O} \left(\nu \left(\frac{\log \kappa}{C} \lfloor \frac{n}{\log n} \rfloor \right)^{-\frac{\log 1/\rho}{\log \kappa}} \right)$.

On the reset condition In practice, physical constraints can force the exploration to interact with the unknown environment under a reset condition. This means that there exists a starting state x and that a trajectory can only be sampled starting from this state and following a given policy. Therefore, if we wish to collect a sample from an arbitrary state x' after a reset, we must first reach that state from the starting state x .

Sequ00L is a strategy that explores the MDP deeper and deeper from an initial starting state x . The original Sequ00L strategy did not consider any reset condition, so to make a comparable analysis we will say that to ‘open’ a node a you first need to reach a at a budget cost of $h(a)$ which is equal to the depth of a . This additional cost has the consequence that Sequ00L with reset will not be able to explore as deeply as without. A naïve extension is shown in Figure 2. Under a total limited budget of n , the number of nodes now open at depth h is $\mathcal{O}(n/h^2)$ instead of $\mathcal{O}(n/h)$ and the maximal depth is now of order of \sqrt{n} instead of n . However, this does not influence the simple regret when $\kappa > 1$ by more than numerical constants as shown in Theorem 2. When $\kappa > 1$, the exponentially diminishing simple regret remains but is changed from ρ^n to $\rho^{\sqrt{n}}$.

Theorem 2. For a planning problem with associated (ν, ρ) , and branching factor $\kappa = \kappa^u(\nu, \rho)$, after n rounds, the simple regret of Sequ00L with reset condition verifies

- If $\kappa = 1$, $r_n \leq \nu \rho^{\frac{1}{C}} \lfloor \frac{n}{\log n} \rfloor$.
- If $\kappa > 1$, $r_n \leq \mathcal{O} \left(\nu \left(\frac{\log \kappa}{C} \lfloor \frac{n}{\log n} \rfloor \right)^{-\frac{\log 1/\rho}{\log \kappa}} \right)$.

5. Deterministic dynamics, stochastic rewards

We now describe the **PlaTγP00S** algorithm detailed in Figure 3. In the presence of noise, it is natural to evaluate the cells multiple times, not just one time as in the deterministic case. The amount of times a cell should be evaluated to

Input: n, A
Initialization: Open the root node \emptyset , h_{\max} times.
 $h_{\max} = \lfloor \frac{n}{2(\log_2 n + 1)^2} \rfloor$, $p_{\max} = \lfloor \log_2(h_{\max}) \rfloor$.
For $h = 1$ to h_{\max} ◀ Exploration ▶
 For $p = \lfloor \log_2(h_{\max}/\lceil h^2 \gamma^{2h} \rceil) \rfloor$ down to 0
 Open $\lceil h 2^p \gamma^{2h} \rceil$ times the at most $\lfloor \frac{h_{\max}}{h \lceil h 2^p \gamma^{2h} \rceil} \rfloor$
 non-opened nodes $a^{h,i} \in A^h$ with highest values
 $\hat{u}(a^{h,i})$ and given $T_{a^{h,i}} \geq \lceil (h-1) 2^p \gamma^{2(h-1)} \rceil$.
 For $p \in [0 : p_{\max}]$ ◀ Cross-validation ▶
 Evaluate $(t+1)\gamma^{2t} h_{\max}(1-\gamma^2)^2$ times the actions
 at time t , a_t^p , of the candidates:
 $a^p = \arg \max_{a \in A^p : \forall t \in [2:h(a)], T_{a_t^p} \geq \lceil (t-1) 2^p \gamma^{2(t-1)} \rceil} \hat{u}(a)$.
 Output $a^n = \arg \max_{\{a^p, p \in [0:p_{\max}]\}} \hat{u}(a^p)$

 Figure 3. The **PlaTγP00S** algorithm

differentiate its value from the optimal value of the function depends on the gap between these two values as well as the range of noise. As we do not want to make *any* assumptions on knowing these quantities, our algorithm tries to be robust to any potential values by not making a fixed choice on the number of evaluations. Intuitively we do this following a path similar to Stroqu00L (Bartlett et al., 2019) by using a modified version of Sequ00L, denoted Sequ00L(m), that allows us to evaluate cells m times, whereas for Sequ00L, $m = 1$. Evaluating cells more times (m large) leads to a better quality of the mean estimates in each cell, however, as a trade-off, it uses more evaluations per depth. This would normally limit us from exploring deep depths of the partition, however, **PlaTγP00S** takes advantage of the knowledge of γ which gives less weight to reward collected deeper in the tree. In order to obtain the concentration results for a node a on $\hat{u}(a) - u(a)$ in Lemma 2, **PlaTγP00S** uses a Chernoff-Hoeffding result that gives with high probability, $\hat{u}(a) - u(a) \leq \sqrt{\sum_{h=0}^{h(a)-1} \gamma^{2h} / T_{a_{[h]}}}$ and balances the range of confidence intervals at different depths. Therefore, **PlaTγP00S** tends to pull less with deeper depth as the number of pull for a fixed m is $\lceil h m \gamma^{2h} \rceil$ where the additional h term ensures that the sum of confidence interval until depth h is bounded for any h . **PlaTγP00S** then *implicitly* performs $\log n$ instances of Sequ00L(m) each with a number of evaluations of $m = 2^p$, where we have $p \in [0 : \log n]$.

The ‘Planning wiTh γ Plus an Online Optimization Strategy’ algorithm **PlaTγP00S** is in Figure 3. Remember that ‘opening’ a node means ‘evaluating’ its children actions. The algorithm opens nodes by sequentially diving them deeper and deeper from the root node $h = 0$ to a maximal depth of h_{\max} . At depth h , we allocate, in an almost decreasing fashion, different number of evaluations

$\lceil h2^p\gamma^{2h} \rceil$ to the nodes with highest value of that depth, with p starting at $\lfloor \log_2(h_{\max}/h) \rfloor$ down to 0. The best node that has been evaluated at least $\mathcal{O}(h_{\max}/h)$ times is opened with $\mathcal{O}(h_{\max}/h)$ evaluations, the two next best cells that have been evaluated at least $\mathcal{O}(h_{\max}/(2h))$ times are opened with $\mathcal{O}(h_{\max}/(2h))$ evaluations, the four next best cells that have been evaluated at least $\mathcal{O}(h_{\max}/(4h))$ times are opened with $\mathcal{O}(h_{\max}/(4h))$ evaluations and so on, until some $\mathcal{O}(h_{\max}/h)$ next best cells that have been evaluated at least once are opened with one evaluation. More precisely, given, p and h , we open, with $\lceil h2^p\gamma^{2h} \rceil$ evaluations, the $\lfloor h_{\max}/(h \lceil h2^p\gamma^{2h} \rceil) \rfloor$ non-previously-opened nodes $a^{h,i} \in A^h$ with highest values $\hat{u}(a^{h,i})$ and given that $T_{a^{h,i}} \geq \lceil (h-1)2^p\gamma^{2(h-1)} \rceil$. The maximum number of evaluations of any node is $2^{p_{\max}}$, with $2^{p_{\max}} = \mathcal{O}(h_{\max})$ as $p_{\max} \triangleq \lfloor \log_2(h_{\max}) \rfloor$. For each $p \in [0 : p_{\max}]$, the candidate output a^p is the node a with highest estimated value such that all actions leading to that node have been evaluated in the following way $\forall t \in [2 : h(a)]$, $T_{a[t]} \geq \lceil (t-1)2^p\gamma^{2(t-1)} \rceil$. We set $h_{\max} \triangleq \lfloor n/(2(\log_2 n + 1)^2) \rfloor$.

5.1. Analysis of P1aTγPOOS

$\perp_{h,p}$ is the depth of the deepest opened node, a with at least $\lceil h2^p\gamma^h \rceil$ evaluations such that there is a $a^* \in A^*$ with $a^* = ab$, with $b \in A^\infty$, at the end of the opening of depth h .

Lemma 1. *For any planning problem with associated (ν, ρ) (see Property 1), on event ξ defined in Appendix C, for any depths $h \in [h_{\max}]$, for any $p \in [0 : \lfloor \log_2(h_{\max}/(h^2\gamma^{2h})) \rfloor]$, we have $\perp_{h,p} = h$ if conditions (1) and (2) simultaneously hold true.*

(1) $b\sqrt{\log(4n/\delta)}/2^{p+1} \leq \nu\rho^h$

(2) We distinguish cases and express the condition in each:

Case 1) $h2^p\gamma^{2h} \leq 1$:

$$\frac{h_{\max}}{h} = \frac{h_{\max}}{h \lceil h2^p\gamma^{2h} \rceil} \geq C\kappa(\nu, \rho)^h$$

And for all $h' \in [h]$, $\frac{h_{\max}}{h'^2 2^{p+1} \gamma^{2h'}} \geq C\kappa(\nu, \rho)^{h'}$.

Case 2) $h2^p\gamma^{2h} \geq 1$:

Case 2.1) $\gamma^2\kappa^u \geq 1$: $\frac{h_{\max}}{h^2 2^{p+1} \gamma^{2h}} \geq C\kappa(\nu, \rho)^h$

Case 2.2) $\gamma^2\kappa^u \leq 1$: $\frac{h_{\max}}{h^2 2^{p+1}} \geq C$.

Lemma 1 gives two conditions so that the cell containing a $a^* \in A^*$ is opened at depth h . This holds if (1) P1aTγPOOS opens, with $\lceil h2^p\gamma^{2h} \rceil$ evaluations, more cells at depth h than the number of near-optimal cells at depth h ($h_{\max}/(h^2 2^{p+1} \gamma^{2h}) \geq C\kappa(\nu, \rho)^h$ if $\gamma^2\kappa^u \geq 1$ and $h_{\max}/h2^p \geq C$ if $\gamma^2\kappa^u \leq 1$) and (2) the $\lceil h2^p\gamma^{2h} \rceil$ evaluations are sufficient to discriminate the empirical average of near-optimal cells from the empirical average of sub-

optimal cells ($b\sqrt{\log(4n/\delta)}/2^{p+1} \leq \nu\rho^h$). To state the next theorems, we introduce \tilde{h} , \tilde{h}_1 and \tilde{h}_2 three positive real numbers satisfying respectively the equations:

$$h_{\max}\nu^2\rho^{2\tilde{h}_1} / \left(\tilde{h}_1 b^2 g_{n,b}^{\delta, R_{\max}} \gamma^{2\tilde{h}_1} \right) = C\kappa^{\tilde{h}_1} \text{ and}$$

$$h_{\max}\nu^2\rho^{2\tilde{h}_2} / \left(\tilde{h}_2 b^2 g_{n,b}^{\delta, R_{\max}} \right) = C, \text{ where}$$

$$\tilde{h}_1 = \frac{1}{\log(\gamma^2\kappa/\rho^2)} \log\left(\frac{\bar{n}_1}{\log \bar{n}_1}\right) + o(1) \text{ and}$$

$$\tilde{h}_2 = \frac{1}{\log(1/\rho^2)} \log\left(\frac{\bar{n}_2}{\log \bar{n}_2}\right) + o(1) \text{ with}$$

$$\bar{n}_1 \triangleq \frac{\nu^2 h_{\max} \log(\gamma^2\kappa/\rho^2)}{C b^2 g_{n,b}^{\delta, R_{\max}}}, \bar{n}_2 \triangleq \frac{\nu^2 h_{\max} \log(1/\rho^2)}{C b^2 g_{n,b}^{\delta, R_{\max}}},$$

where $g_{n,b}^{\delta, R_{\max}} \triangleq p_{\max} \log(R_{\max} n^{3/2}/b(1-\delta))$. \tilde{h} is defined similarly in Equation 5 in Appendix D. The quantities \tilde{h}_1 and \tilde{h}_2 give the respective depths of deepest cell opened by P1aTγPOOS that contains a a^* with high probability in the cases $\gamma^2\kappa \geq 1$ and $\gamma^2\kappa \leq 1$. Additionally, \tilde{h}_1 and \tilde{h}_2 also lets us characterize for which regime of the noise range b we recover results similar to the loss of the deterministic case. Discriminating on the noise regimes, we now state two of our results, Theorem 3 for a high noise and Theorem 4 for a low one. A more exhaustive list of results is in the Appendix D or in the Table 1.

Theorem 3. High-noise regime After n rounds, for any problem with associated (ν, ρ) , and branching factor denoted $\kappa = \kappa^u(\nu, \rho)$, if the noise b is high enough to verify both high-noise conditions as defined in the caption of Table 1, the simple regret of P1aTγPOOS obeys

$$\mathbb{E}r_n = \begin{cases} \tilde{\mathcal{O}}\left(\left(\frac{n}{b^2}\right)^{-\frac{1}{2}}\right) & \text{if } \gamma^2\kappa \leq 1, \\ \tilde{\mathcal{O}}\left(\left(\frac{n}{b^2}\right)^{-\frac{\log(1/\rho)}{\log(\gamma^2\kappa/\rho^2)}}\right) & \text{if } \gamma^2\kappa > 1. \end{cases}$$

The proofs are in appendix D. They are quite technical but they are simply based on checking the conditions of Lemma 1 under different b, ρ, γ, κ regimes.

Theorem 4. Low-noise regime After n rounds, for any problem with associated (ν, ρ) , and branching factor denoted $\kappa = \kappa^u(\nu, \rho)$, if the noise b is low enough to verify both high-noise conditions as defined in the caption of Table 1, the simple regret of P1aTγPOOS obeys

$$\mathbb{E}r_n = \begin{cases} \tilde{\mathcal{O}}(\nu\rho^n) & \text{if } \kappa = 1, \\ \tilde{\mathcal{O}}\left(\nu\left(\frac{n}{b^2}\right)^{-\frac{\log(1/\rho)}{\log(\kappa)}}\right) & \text{if } \kappa > 1. \end{cases}$$

Worst-case comparison with OLOP When b is large and known: In Table 1, we give our results in various classes

	$\gamma^2 \kappa \leq 1$		$\gamma^2 \kappa \geq 1$	
	High noise (ii)	Low noise (ii)	High noise (iii)	Low noise (iii)
High noise (i)	$m_{\nu,\gamma} \left(\frac{n}{b^2}\right)^{-\frac{1}{2}}$	$\nu \rho \sqrt{n}$	$m_{\nu,\gamma} \left(\frac{n}{b^2}\right)^{-\frac{\log(1/\rho)}{\log(\gamma^2 \kappa / \rho^2)}}$	$\nu \left(\frac{n}{b^2}\right)^{-\frac{\log(1/\rho)}{\log(\kappa)}}$
Low noise (i)	$m_{\nu,\gamma} \left(\frac{n}{b^2}\right)^{-\frac{\log(1/\rho)}{\log(\kappa)}}$	$\kappa = 1 : \nu \rho^n$	$m_{\nu,\gamma} \left(\frac{n}{b^2}\right)^{-\frac{\log(1/\rho)}{\log(\kappa)}}$	$\nu \left(\frac{n}{b^2}\right)^{-\frac{\log(1/\rho)}{\log(\kappa)}}$
		$\kappa > 1 : \nu \left(\frac{n}{b^2}\right)^{-\frac{\log(1/\rho)}{\log(\kappa)}}$		

Table 1. Rates of the our upper bounds on the simple regret of **PlaTγPOOS** in various classes of problems. The condition on noise (i) is whether the noise b verifies $\nu^2 \rho^{2\tilde{h}} / \left(\gamma^{2\tilde{h}} \tilde{h} b^2 g_{n,b}^{\delta, R_{\max}}\right) \leq 1$. The condition on noise (ii) is whether $\nu^2 \rho^{2\tilde{h}} / \left(b^2 g_{n,b}^{\delta, R_{\max}}\right) \leq 1$. The condition on noise (iii) is whether $\nu^2 \rho^{2\tilde{h}} / \left(b^2 g_{n,b}^{\delta, R_{\max}}\right) \leq 1$. $m_{\nu,\gamma} = \max(1/1 - \gamma^2, \nu)$.

of problems depending on the whether $\gamma^2 \kappa \geq 1$, whether $\kappa = 1$ or $\kappa > 1$, and several conditions on the range of the noise b . The results for OLOP were distinguishing the results based on $\gamma^2 \kappa$ being greater or smaller than 1. For this two cases we recover, as displayed in Table 1 with a grey background color, the same rate in term of n (for instance taking $b = 1$ like in OLOP and having $\rho = \gamma$) for the simple regret. However, we provide more specific treatment for sub-cases with associated improvements that we list and detail next.

Adaptation to the range of the noise b without a prior knowledge Our analysis shows that **PlaTγPOOS** adapts favorably to the unknown range of noise. Already, in the standard cases discussed above, where the noise is large, our bound already adapts to the amount of noise as it scales with n/b^2 . OLOP requires an estimate \tilde{b} of b and has a regret scaling with n/\tilde{b}^2 which is problematic in case of a wrong estimate $\tilde{b} \gg b$. Moreover, we give technical conditions on the range of noise that shows when **PlaTγPOOS** gets improved rates. When $\gamma^2 \kappa \geq 1$, OLOP was already obtaining rates that were the same as the rates of deterministic reward case. Therefore, beyond the n/b^2 improvement and the adaptation to extra smoothness that will be discussed latter, no more rate improvement should be expected. In the case $\gamma^2 \kappa \leq 1$, the improvement are even more striking. When the noise is very low instead of the OLOP rate of $\mathcal{O}(1/\sqrt{n})$, we obtain the deterministic rate of OPD (Hren & Munos, 2008; Munos, 2014) which is either $n^{-\frac{\log(1/\rho)}{\log \kappa}}$ or ρ^n . These improved rates *could not* be obtained by OLOP. Indeed, OLOP relies on upper confidence bound (UCB) that uses a range of the noise \tilde{b} as input,

$$\bar{u}(a) = \mathcal{O} \left(\hat{u}(a) + \sum_{h=0}^{h(a)-1} \gamma^h \tilde{b} \sqrt{\frac{1}{T_{a[a]}}} + \frac{R_{\max} \gamma^{h(a)}}{1 - \gamma} \right).$$

This works if $\tilde{b} = b$. However the true b is unknown. If $b = 0$, using any $\tilde{b} > 0$ will not result in an improved rate.

Adaptation to additional smoothness ν and ρ beyond γ As defined in Section 2, we aim to adapt to the true smoothness ν, ρ of V which can go beyond γ . Our bounds show that **PlaTγPOOS** is able to take advantage of ν, ρ in a large portion of cases. In most cases in the rate the γ in OLOP is replaced by ρ in **PlaTγPOOS**. In the case where $\gamma^2 \kappa \geq 1$ we have $r^{\text{OLOP}} = \mathcal{O}(n^{-\frac{\log(1/\gamma)}{\log(\kappa)}}) \leq \mathcal{O}(n^{-\frac{\log(1/\rho)}{\log(\gamma^2 \kappa / \rho^2)}}) = r^{\text{PlaTγPOOS}}$.

Adaptation to the deterministic case and $\kappa = 1$ **PlaTγPOOS** adapts to the branching factor of the problem κ in a way that, under low noise conditions, leads to an exponentially decreasing simple regret $r_n^{\text{PlaTγPOOS}} = \mathcal{O}(\nu \rho \sqrt{n})$. This is a light-years improvement over OLOP in these conditions as OLOP's regret is at best $r_n^{\text{OLOP}} = \mathcal{O}(n^{-1/2})$. This result is possible because **PlaTγPOOS** can explore much deeper than OLOP, as its maximal depth is of order n . Actually, in most scenarios, the actual larger depth explored will be of order \sqrt{n} due to sample limitations. On the other hand, OLOP can only go $\log n$ deep.

Moreover, $\kappa = 1$ is a common case in planning. Indeed, as discussed by Bubeck & Munos (2010), $\kappa = 1$ is equivalent to having near-optimal dimension $d = 0$ in an optimization task (Munos, 2014) which is a common value as shown by Valko et al. (2013). Therefore, we expect the case when $\gamma \kappa^2 \leq 1$, that is, where we get the most significant improvement over OLOP, to be the most common in practice.

The reset condition The effect of the reset condition, as discussed in Section 4, affects, in the stochastic reward case, **PlaTγPOOS** as follows. First, all polynomial rates stay the same. Then, only the exponential rates change. A ρ^n rate without the condition becomes a $\rho \sqrt{n}$ with it. Next, a $\rho \sqrt{n}$

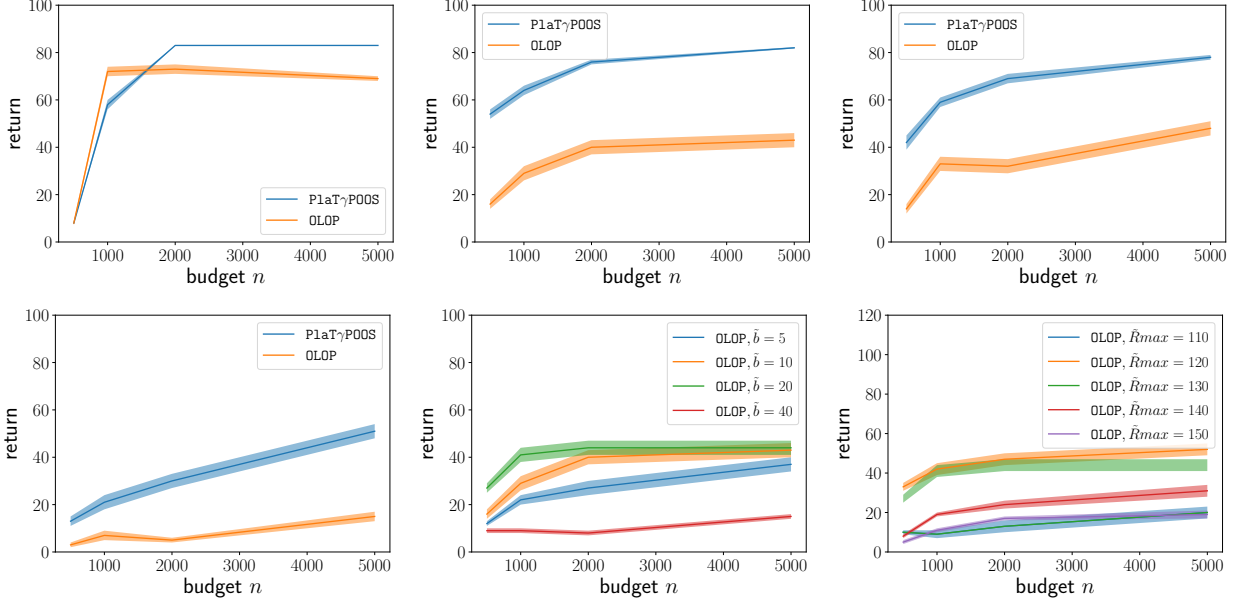


Figure 4. *Top and bottom left*: Expected cumulative discounted return collected by OLOP and P1aTγPOOS with different range of noise, $b = 1$ (top left), $b = 10$ (top center) $b = 20$, (top right) and $b = 50$ (bottom left). *Bottom center*: the sensitivity of OLOP to different \tilde{R}_{\max} parameters. *Bottom right*: the sensitivity of OLOP to different range of the input noise \tilde{b} as parameters while the true b is set to 10.

rate becomes a $\rho^{n^{1/3}}$ one.

6. Numerical experiments

In this section, we empirically illustrate the benefits of P1aTγPOOS. We chose a simple MDP, shown in Figure 5. In this MDP, a state $x = (\text{bin}, d)$ is a pair of a binary variable bin and a non-negative integer d . The MDP has two actions that are also a binary variable a . If $\text{bin} \neq a$, the base reward is 2, in which case, the next state is $(a, 0)$. Otherwise, if $\text{bin} = a$, then $r = d$ and the next state is $(a, d + 1)$. The reward is then shifted by adding 100 to it so that the noises with different ranges can be added on top without making the reward negative.

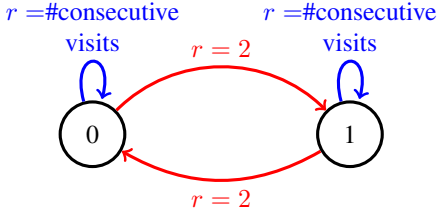


Figure 5. MDP used for experiments

The initial state is $(0, 0)$. Therefore, the agent has a choice. It can, for instance, remain in the same binary state bin , starting with a null reward but see its instant reward growing with time if it keeps taking the same action in the future. Alternatively, it could greedily switch to the other binary

state bin and obtain a reward of 2 but delaying the hope of obtaining growing reward as in the first scenario. We set $\gamma = 0.95$. Therefore $R_{\max} \approx 130$.

Figure 4 reports the results. All the figures show the cumulative discounted return collected by OLOP and P1aTγPOOS after having interacted for 20 steps with the MDP, having chosen each time an action following their planning strategy and then being transferred to the state resulting of applying the recommended action in the current state and therefore also collecting a reward that is composing the final return.

Note that the return reported are shifted in order to not take into account the 100 fixed part of each the reward. The figures in the top row, as well as the figure at the bottom left, reports the comparison between the two returns of OLOP and P1aTγPOOS for different ranges of noise b . P1aTγPOOS is systematically outperforming OLOP while in this case OLOP is given as input parameters the correct \tilde{R}_{\max} , with $\tilde{R}_{\max} = R_{\max}$ and the correct range of the noise \tilde{b} as $\tilde{b} = b$.

In Figure 4, bottom center and right, we illustrate the sensitivity of OLOP to misleading input parameters. Notice that the performance of OLOP is very vulnerable to these mis-specifications while P1aTγPOOS is not using such inputs.

Acknowledgments The research presented was supported by European CHIST-ERA project DELTA, French Ministry of Higher Education and Research, Nord-Pas-de-Calais Regional Council, Inria and Otto-von-Guericke-Universität Magdeburg associated-team north-European project Allo-

cate, and French National Research Agency project BoB (grant n.ANR-16-CE23-0003), FMJH Program PGMO with the support to this program from Criteo.

References

- Bartlett, P. L., Gabillon, V., and Valko, M. A simple parameter-free and adaptive approach to optimization under a minimal local smoothness assumption. In *Algorithmic Learning Theory*, 2019.
- Bertsekas, D. and Tsitsiklis, J. *Neuro-Dynamic Programming*. Athena Scientific, Belmont, MA, 1996.
- Bubeck, S. and Munos, R. Open-loop optimistic planning. In *Conference on Learning Theory*, 2010.
- Bubeck, S., Munos, R., Stoltz, G., and Szepesvári, C. X-armed bandits. *Journal of Machine Learning Research*, 12:1587–1627, 2011.
- Buşoniu, L. and Munos, R. Optimistic planning for Markov decision processes. In *International Conference on Artificial Intelligence and Statistics*, 2012.
- Carpentier, A. and Locatelli, A. **Tight (lower) bounds for the fixed budget best-arm identification bandit problem**. In *Conference on Learning Theory*, 2016.
- Coquelin, P.-A. and Munos, R. Bandit algorithms for tree search. In *Uncertainty in Artificial Intelligence*, 2007.
- Coulom, R. Efficient selectivity and backup operators in Monte-Carlo tree search. *Computers and games*, 4630: 72–83, 2007.
- Feldman, Z. and Domshlak, C. Simple regret optimization in online planning for Markov decision processes. *Journal of Artificial Intelligence Research*, 2014.
- Gabillon, V., Ghavamzadeh, M., and Lazaric, A. Best arm identification: A unified approach to fixed budget and fixed confidence. In *Neural Information Processing Systems*, 2012.
- Gelly, S., Yizao, W., Munos, R., and Teytaud, O. Modification of UCT with patterns in Monte-Carlo Go. Technical report, Inria, 2006.
- Grill, J.-B., Valko, M., and Munos, R. Blazing the trails before beating the path: Sample-efficient monte-carlo planning. In *Neural Information Processing Systems*, 2016.
- Hoorfar, A. and Hassani, M. Inequalities on the lambert w function and hyperpower function. *Journal of Inequalities in Pure and Applied Mathematics*, 9(2):5–9, 2008.
- Hren, J.-F. and Munos, R. Optimistic Planning of Deterministic Systems. In *European Workshop on Reinforcement Learning*, 2008.
- Kaufmann, E. and Koolen, W. M. Monte-carlo tree search by best arm identification. In *Advances in Neural Information Processing Systems*, pp. 4897–4906, 2017.
- Kocsis, L. and Szepesvári, C. Bandit-based Monte-Carlo planning. In *European Conference on Machine Learning*, 2006.
- Leurent, E. and Maillard, O.-A. Practical open-loop optimistic planning. *arXiv preprint arXiv:1904.04700*, 2019.
- Munos, R. Optimistic optimization of deterministic functions without the knowledge of its smoothness. In *Neural Information Processing Systems*, 2011.
- Munos, R. From bandits to Monte-Carlo tree search: The optimistic principle applied to optimization and planning. *Foundations and Trends in Machine Learning*, 7(1):1–130, 2014.
- Orabona, F. and Pál, D. Scale-free online learning. *Theoretical Computer Science*, 2018.
- Orabona, F. and Tommasi, T. Backprop without Learning Rates Through Coin Betting. In *Neural Information Processing Systems*, 2017.
- Puterman, M. L. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, New York, NY, 1994.
- Ross, S., Mineiro, P., and Langford, J. Normalized online learning. In *Uncertainty in Artificial Intelligence*, may 2013.
- Shah, D., Xie, Q., and Xu, Z. On reinforcement learning using monte carlo tree search with supervised learning: Non-asymptotic analysis. *arXiv preprint arXiv:1902.05213*, 2019.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., van den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., Dieleman, S., Grewe, D., Nham, J., Kalchbrenner, N., Sutskever, I., Lillicrap, T., Leach, M., Kavukcuoglu, K., Graepel, T., and Hassabis, D. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- Szörényi, B., Kedenburg, G., and Munos, R. Optimistic planning in Markov decision processes using a generative model. In *Neural Information Processing Systems*, 2014.
- Valko, M., Carpentier, A., and Munos, R. Stochastic simultaneous optimistic optimization. In *International Conference on Machine Learning*, 2013.