

Model-Free Online Learning in Unknown Sequential Decision Making Problems and Games

Abstract

1 Regret minimization has proved to be a versatile tool for se- 43
2 quential decision making and extensive-form games. In large 44
3 two-player zero-sum imperfect-information games, modern 45
4 extensions of counterfactual regret minimization (CFR) are 46
5 currently the practical state of the art for computing a Nash 47
6 equilibrium. Most regret-minimization algorithms for sequen- 48
7 tial decision making, including CFR, require (i) an exact *model* 49
8 of the player’s decision nodes, observation nodes, and how 50
9 they are linked, and (ii) full knowledge, at all times t , about 51
10 the payoffs—even in parts of the decision space that are not 52
11 encountered at time t . Recently, there has been growing in- 53
12 terest towards relaxing some of those restrictions and making 54
13 regret minimization applicable to settings for which rein- 55
14 forcement learning methods have traditionally been used—for 56
15 example, those in which only black-box access to the envi- 57
16 ronment is available. We give the first, to our knowledge, 58
17 regret-minimization algorithm that guarantees sublinear regret 59
18 with high probability even when requirement (i)—and thus 60
19 also (ii)—is dropped. We formalize an online learning setting 61
20 in which the strategy space is not known to the agent and gets 62
21 revealed incrementally whenever the agent encounters new 63
22 decision points. We give an efficient algorithm that achieves 64
23 $O(T^{3/4})$ regret with high probability for that setting, even 65
24 when the agent faces an adversarial environment. Our experi- 66
25 ments show it significantly outperforms the prior algorithms 67
26 for the problem, which do not have such guarantees. It can 68
27 be used in any application for which regret minimization is 69
28 useful: approximating Nash equilibrium or quantal response 70
29 equilibrium, approximating coarse correlated equilibrium in 71
30 multi-player games, learning a best response, learning safe 72
31 opponent exploitation, and online play against an unknown 73
32 opponent/environment.

1 Introduction

33
34 A *sequential decision making (SDM)* problem formalizes in 74
35 a tree-form structure the interaction of an agent with an un- 75
36 known and potentially adversarial environment. The agent’s 76
37 tree includes decision nodes, observation nodes, and termi- 77
38 nal nodes. SDM captures the problem that a player faces 78
39 in an extensive-form game. SDM also captures MDPs and 79
40 POMDPs where the agent conditions on observed history, 80
41 but SDM problems are more general because the Markovian 81
42 assumption is not necessarily made. 82

Copyright © 2021, Association for the Advancement of Artificial 83
Intelligence (www.aaai.org). All rights reserved. 84

In SDM, the environment may react adversarially to what 43
the agent does. This is important to take into account, for 44
example, in game-theoretic settings and *multiagent reinforcement 45*
learning (MARL) because the other agents’ learning 46
makes the environment nonstationary for the agent (Sand- 47
holm and Crites 1996; Matignon, Laurent, and Le Fort-Piat 48
2012). This is in contrast to the standard assumption in single- 49
agent reinforcement learning where the environment is *oblivious* 50
to the agent instead of adversarial. Hence, learning strong 51
policies (aka. *strategies*) in SDM problems is especially chal- 52
lenging, and the agent must be careful about exploration 53
because exploration actions can change the environment. 54

Online, regret minimization methods have been success- 55
fully used in SDM. In particular, the *counterfactual regret 56*
minimization (CFR) framework decomposes the overall regret 57
of an agent to local regrets at individual decision nodes 58
(aka. information sets in game theory) (Zinkevich et al. 2007). 59
That enables significantly larger SDMs to be tackled. Many 60
enhancements have been developed on top of the basic CFR 61
framework (Lanctot et al. 2009; Tammelin 2014; Brown 62
and Sandholm 2015; Brown, Kroer, and Sandholm 2017; 63
Brown and Sandholm 2017a, 2019a; Farina, Kroer, and Sand- 64
holm 2019), and have led to major milestones in imper- 65
fect-information games such as poker (Bowling et al. 2015; 66
Moravčík et al. 2017; Brown and Sandholm 2017b, 2019b). 67
Many of those methods *guarantee* low regret even against an 68
adversarial environment—which, in turn, enables the compu- 69
tation of game-theoretic solutions such as Nash equilibrium, 70
coarse correlated equilibrium (Moulin and Vial 1978; Celli 71
et al. 2020), best responses, etc. 72

However, those methods usually come with two drawbacks: 73
(i) they require an explicit upfront *model* of the agent’s deci- 74
sion space, and (ii) depending on the online learning model 75
used, they require full feedback, at all times t , about the 76
payoffs assigned by the environment—even in parts of the 77
decision space not encountered at time t . There has been work 78
towards an online learning setting, called the *bandit optimiza- 79*
tion setting, that drops (ii) (Lattimore and Szepesvári 2020). 80
Most MARL algorithms apply for the unknown game setting 81
and drop both (i) and (ii), often at the cost of invalidating any 82
regret guarantees. In this paper, we give, to our knowledge, 83
the first regret-minimization algorithm that guarantees sublin- 84
ear (specifically $O(T^{3/4} \sqrt{\log 1/p})$) with probability $1 - p$) 85
regret even when requirements (i) and (ii) are both dropped. 86

Closely Related Research

In the rest of this section we discuss how our algorithm relates to other attempts to connect online learning guarantees with model-free MARL. A comparison between our algorithm and other settings in online learning is deferred to Section 3, where we formally introduce our interactive bandit model.

The greatest inspiration for our algorithm is the *online variant of Monte Carlo CFR (online MCCFR)* algorithm proposed in passing by (Lanctot et al. 2009). Unlike the traditional “self-play” MCCFR, *online MCCFR* does not assume that the algorithm have control over the environment. The authors note that in theory online MCCFR could be used to play games in a model-free fashion, provided that a lower bound on the reach probability of every terminal state can be enforced. That lower bound, say η , is necessary for them to invoke their main theorem, which guarantees $O(\eta^{-1}T^{1/2}\sqrt{1/p})$ regret with probability $1 - p$. They suggest perhaps using some form of ϵ -greedy exploration at each decision node to enforce the lower bound, but no guarantee is provided and the authors then move away from this side note to focus on the self-play case. We show that their proposed approach encounters significant hurdles. First, using exploration at each decision node results in a lower bound on the reach of every terminal state on the order of $\eta = \epsilon^d$, where d is the depth of the decision process, thus making the regret bound not polynomial in the size of the decision process itself, but rather exponential. Second, the paper did not provide theoretical guarantees for the online case. In particular, the theory does not take into account the degradation effect caused by the exploration itself, which scales roughly as ηT . So, on the one hand, a large η is needed to keep the term η^{-1} in their regret bound under control, but at the same time a large η results in a ηT term being added to the regret. These hurdles show that it is unlikely that their approach can lead to $O(T^{1/2}\sqrt{1/p})$ regret with high probability $1 - p$ as they hypothesized. We address those issues by using a different type of exploration and being careful about bounding the degradation terms in the regret incurred due to the exploration. Because of the latter, our algorithm incurs $O(T^{3/4}\sqrt{\log(1/p)})$ regret with high probability against adversarial opponents. Because the exponent is less than 1, ours is truly a regret minimization algorithm for SDMs with unknown structure, and to our knowledge, the first. At the same time, our exponent is worse than the originally hypothesized exponent $1/2$. It is unknown whether the latter can be achieved. Finally, our dependence on p is better than in their hypothesized regret bound.

A recent paper by Srinivasan et al. (2018) related policy gradient algorithms to CFR (and, to a much lesser degree, MCCFR). Despite their experiments demonstrating empirical convergence rates for sampled versions of their algorithms, formal guarantees are only obtained for tabularized policy iteration in self-play, and use policy representations that require costly ℓ_2 projections back into the policy space. In contrast, our regret guarantees hold (i) in any SDM setting (not just two-player zero-sum extensive-form games), (ii) in high probability, (iii) with sampling, and (iv) when playing against any environment, even an adversarial one, without

Conceptually, our algorithm has elements of both online bandit optimization and MARL. On the one hand, our regret guarantees hold with high probability no matter how the environment picks its actions or assigns utilities to terminal nodes at all iterations t . In fact, our algorithm is a regret minimizer in the proper online-learning sense: it does not need to know—and makes no assumptions about—the underlying policy of the environment or the utility at the terminal states. Furthermore, those quantities are not assumed to be time independent and they can even be selected adversarially by the environment at each iteration. This is contrast with many *self-play* methods, that require control over the opponent in order to retain guarantees. In particular, because we assume no control over the adversarial environment, every interaction with the environment can lead it to react and change behavior in the next iteration. So, it is impossible to “freeze” the policy of the environment to perform off-policy exploration like many self-play methods require. Because of its strong online guarantees, our algorithm can be used for all the applications in which regret minimization provides benefits—for example, to converge to Nash equilibrium in a two-player zero-sum extensive-form game, to learn a best response against static opponents, to converge to coarse correlated equilibrium in a multiagent setting (Moulin and Vial 1978; Celli et al. 2020), to converge to a quantal-response equilibrium (Ling, Fang, and Kolter 2018; Farina, Kroer, and Sandholm 2019), to compute safe exploitative strategies (Farina, Kroer, and Sandholm 2019; Ponsen, De Jong, and Lanctot 2011), or to play online against an unknown opponent/environment.

On the other hand, ours is *not* an algorithm for the online bandit optimization problem. In bandit optimization, the agent does *not* interact with the environment: at all times t , the agent outputs a policy π^t for the whole decision space, and receives a single real number u^t , representing the gain incurred by π^t , as feedback. Instead, our algorithm operates within a slightly different online learning model that we introduce, which we call the *interactive bandit model*. In it, the decision maker gets to observe the signals (actions) selected by the environment on the path from the root of the decision problem to the agent’s terminal state, in additions to u^t . Hence, we operate within a more specific online learning model than bandit optimization (which applies to, for example, picking a point on a sphere also, as opposed to just SDM), and rather one that is natural to SDM. This comes with a significant advantage. While, to our knowledge, all algorithms for the bandit optimization problem require *a priori* knowledge of the full structure of the sequential decision problem,¹ our algorithm for the interactive bandits model is *model free*. Here, the structure of the sequential decision process is at first unknown and has to be discovered by exploring as part of the decision-making process. Decision and observation nodes are revealed only at the time the agent encounters them for the first time.

¹This is needed, for example, to construct a self-concordant barrier function (Abernethy, Hazan, and Rakhlin 2008) or a global regularizer for the strategy space (Abernethy and Rakhlin 2009; Farina, Schmucker, and Sandholm 2020), to compute a barycentric spanner (Awerbuch and Kleinberg 2004; Bartlett et al. 2008), or to set up a kernel function (Bubeck, Lee, and Eldan 2017).

198 requiring complex projections.

199 A very recent paper (Zhou, Li, and Zhu 2020) on two-
200 player zero-sum games proposes combining full-information
201 regret minimization with posterior sampling (Strens 2000)
202 to estimate the utility function of the player and transition
203 model of chance, both of which they assume to be time inde-
204 pendent unlike our setting. They show in-expectation bounds
205 under the assumption that the agent’s strategy space is known
206 *ex ante*. We operate in a significantly more general setting
207 where the observations and utilities are decided by the en-
208 vironment and can change—even adversarially—between
209 iterations. Like theirs, our algorithm converges to Nash equi-
210 librium in two-player zero-sum games when used in self play.
211 However, unlike theirs, our algorithm is a regret minimizer
212 that can be used for other purposes also, such as finding
213 a coarse correlated equilibrium in multiplayer general-sum
214 games, a quantal-response equilibrium, or safe exploitative
215 strategies. In the latter two applications, the payoff function,
216 in effect, changes as the agent changes its strategy. Our regret
217 guarantees hold in high probability and we do not assume *ex*
218 *ante* knowledge of the agent’s strategy space.

219 A different line of work has combined fictitious play
220 (Brown 1951) with deep learning for function approxima-
221 tion (Heinrich, Lanctot, and Silver 2015; Heinrich and Silver
222 2016). Those methods do not give regret guarantees. Finally,
223 other work has studied how to combine the guarantees of on-
224 line learning with MDPs. Kash, Sullins, and Hofmann (2020)
225 combine the idea of breaking up and minimizing regret lo-
226 cally at each decision point, proper of CFR, with Q-learning,
227 obtaining an algorithm with certain in-the-limit guarantees
228 for MDPs. Even-Dar, Kakade, and Mansour (2009) study
229 online optimization (in the full-feedback setting, as opposed
230 to bandit) in general MDPs where the reward function and
231 the structure of the MDP is known. Neu et al. (2010) study
232 online bandit optimization in MDPs in the *oblivious* setting,
233 achieving $O(T^{2/3})$ regret with high probability, again assum-
234 ing that the MDP’s structure is known and certain conditions
235 are met. Zimin and Neu (2013) give bandit guarantees for
236 episodic MDPs with a fixed known transition function.

237 2 Our Model for (Unknown) Sequential 238 Decision Making and Games

239 In this section, we introduce the notation for SDM problems
240 that we will be using the rest of the paper.

241 A *sequential decision making (SDM) problem* is structured
242 as a tree made of three types of nodes: (i) *decision nodes*
243 j , in which the agent acts by selecting an action from the
244 finite set A_j (different decision nodes can admit different
245 sets of allowed actions); (ii) *observation points* k , where the
246 agent observes one out of set S_k of finitely many possible
247 signal from the environment (different observation points can
248 admit different sets of possible signals); and (iii) *terminal*
249 *nodes*, corresponding to the end of the decision process. We
250 denote the set of decision nodes in the sequential decision
251 making problem as \mathcal{J} , the set of observation points as \mathcal{K} ,
252 and the set of terminal nodes as Z . Furthermore, we let ρ
253 denote the dynamics of the game: selecting action $a \in A_j$
254 at decision node $j \in \mathcal{J}$ makes the process transition to the

255 next state $\rho(j, a) \in \mathcal{J} \cup \mathcal{K} \cup Z$, while the process transitions
256 to $\rho(k, s) \in \mathcal{J} \cup \mathcal{K} \cup Z$ whenever the agent observes signal
257 $s \in S_k$ at observation point $k \in \mathcal{K}$.

258 Our algorithm operates in the difficult setting where the
259 structure of the SDM problem is at first unknown and can be
260 discovered only through exploration. Decisions and observa-
261 tion nodes are revealed only at the time the agent encounters
262 them for the first time. As soon as a decision node j is re-
263 vealed for the first time, its corresponding set of actions A_j
264 is revealed with it.

265 **Sequences** In line with the game theory literature, we call
266 a *sequence* a decision node-action pair; each sequence (j, a)
267 uniquely identifies a path from the root of the decision pro-
268 cess down to action a at decision node j , included. Formally,
269 we define the set of sequences as $\Sigma := \{(j, a) : j \in \mathcal{J}, a \in$
270 $A_j\} \cup \{\emptyset\}$, where the special element \emptyset is called the *empty*
271 *sequence*. Given a decision node $j \in \mathcal{J}$, its *parent sequence*,
272 denoted p_j , is the last sequence (that is, decision node-action
273 pair) encountered on the path from the root of the decision
274 process down to j . If the agent does not act before j (that is,
275 only observation points are encountered on the path from the
276 root to j), we let $p_j = \emptyset$.

277 Given a terminal node $z \in Z$ and a sequence $(j, a) \in \Sigma$,
278 we write $(j, a) \rightsquigarrow z$ to mean that the path from the root of the
279 decision process to z passes through action a at decision node
280 j . Similarly, given a terminal node $z \in Z$ and an observation
281 node-signal pair (k, s) ($s \in S_k$), we write $(k, s) \rightsquigarrow z$ to
282 indicate that the path from the root of the decision process
283 to z passes through signal s at observation node k . Finally,
284 we let $\sigma(z)$ be the last sequence (decision node-action pair)
285 on the path from the root of the decision process to terminal
286 node $z \in Z$.

287 **Strategies** Conceptually, a strategy for an agent in a se-
288 quential decision process specifies a distribution $x_j \in \Delta^{|A_j|}$
289 over the set of actions A_j at each decision node $j \in \mathcal{J}$. We
290 represent a strategy using the *sequence-form representation*,
291 that is, as a vector $\mathbf{q} \in \mathbb{R}_{\geq 0}^{|\Sigma|}$ whose entries are indexed by
292 Σ . The entry $q[j, a]$ contains the product of the probabili-
293 ties of all actions at all decision nodes on the path from
294 the root of the process down to and including action a at
295 decision node $j \in \mathcal{J}$. A vector $\mathbf{q} \in \mathbb{R}_{\geq 0}^{|\Sigma|}$ is a valid sequence-
296 form strategy if and only if it satisfies the constraints (i)
297 $\sum_{a \in A_j} q[j, a] = q[p_j]$ for all $j \in \mathcal{J}$; and (ii) $x[\emptyset] = 1$ (Ro-
298 manovskii 1962; Koller, Megiddo, and von Stengel 1994;
299 von Stengel 1996). We let Q denote the set of valid sequence-
300 form strategies. Finally, we let $\Pi \subseteq Q$ denote the subset of
301 sequence-form strategies whose entries are only 0 or 1; a
302 strategy $\pi \in \Pi$ is called a *pure* sequence-form strategy, as it
303 assigns probability 1 to exactly one action at each decision
304 node.

305 3 Online Learning and Our Interactive 306 Bandit Model

307 In online learning, an agent interacts with its environment
308 in this order: (i) The environment chooses a (secret) gradi-
309 ent vector ℓ^t of bounded norm; (ii) The agent picks a pure
310 strategy $\pi^t \in \Pi$. The environment evaluates the reward

311 (gain) of the agent as $(\ell^t)^\top \pi^t \in \mathbb{R}$; (iii) The agent observes
312 some feedback about her reward. The feedback is used by
313 the agent to learn to output good strategies over time. The
314 learning is online in the sense that the strategy π^t at time t
315 is output before any feedback for it (or future strategies) is
316 available. A standard quality metric for evaluating an agent
317 in this setting is the *regret* that she accumulates over time:
318 $R^T(\hat{\pi}) := \sum_{t=1}^T (\ell^t)^\top \hat{\pi} - \sum_{t=1}^T (\ell^t)^\top \pi^t$. This measures the
319 difference between the total reward accumulated up to time
320 T , compared to the reward that would have been accumulated
321 had the oracle output the fixed strategy $\hat{\pi} \in \Pi$ at all times. A
322 “good” agent, that is, a *regret minimizer*, is one whose regret
323 grows sublinearly: $R^T(\hat{\pi}) = o(T)$ for all $\hat{\pi} \in \Pi$.

324 Online learning models vary based on the type and extent
325 of feedback that is made available to the agent. We will
326 focus on two existing models—namely, the full-feedback²
327 setting and the bandit linear optimization setting. Then we
328 will introduce a third model that is especially natural for
329 SDM.

330 **Full-Feedback² Setting** Here, the environment always re-
331 veals the full gradient vector ℓ^t to the agent (after the strategy
332 has been output). This is the same setting that was pro-
333 posed in the landmark paper by Zinkevich (2003) and is
334 the most well-studied online optimization setting. Efficient
335 agents that guarantee $O(T^{1/2})$ regret with high probability
336 in the full-feedback setting are known well beyond SDM and
337 extensive-form games (e.g., (Shalev-Shwartz 2012; Hazan
338 2016)). In fact, given any convex and compact set \mathcal{X} , it is
339 possible to construct an agent that outputs decisions $x^t \in \mathcal{X}$
340 that achieves $O(T^{1/2})$ regret *in the worst case*, even when
341 the gradient vectors ℓ^t are chosen adversarially by the en-
342 vironment *after* the decision x^t has been revealed. In the
343 specific context of SDM and extensive-form games, the most
344 widely-used oracle in the full-feedback setting is based on
345 the counterfactual regret (CFR) framework (Zinkevich et al.
346 2007). The idea is to decompose the task of computing a strategy
347 for the whole decision process into smaller subproblems
348 at each decision point. The local strategies are then computed
349 via $|\mathcal{J}|$ independent full-information regret minimizers, one
350 per decision node, that at each t observe a specific feedback
351 that guarantees low global regret across the decision process.
352 CFR guarantees $O(T^{1/2})$ regret with high probability against
353 any strategy $\pi \in \Pi$.

354 **Bandit Linear Optimization** Here, the only feedback
355 that the environment reveals to the agent is the utility
356 $(\ell^t)^\top \pi^t$ at each time t (Kleinberg 2004; Flaxman, Kalai,
357 and McMahan 2005). Despite this extremely limited feed-
358 back, $\tilde{O}(T^{1/2})$ regret can still be guaranteed with high prob-
359 ability in some domains (including simplexes (Auer et al.
360 2002) and spheres (Abernethy and Rakhlin 2009)), although,
361 to our knowledge, a polynomial algorithm that guarantees
362 $\tilde{O}(T^{1/2})$ regret with high probability for any convex and

compact domain has not been discovered yet.³ Guarantee-
363 ing $\tilde{O}(T^{1/2})$ *in expectation* is possible for any domain of
364 decisions (Abernethy, Hazan, and Rakhlin 2008), but unfor-
365 tunately in-expectation low regret is not strong enough a
366 guarantee to enable, for instance, convergence to Nash or
367 correlated equilibrium as described above. In the specific
368 case of sequential decision processes, (Farina, Schmucker,
369 and Sandholm 2020) proposed a bandit regret minimizer
370 that achieves $O(T^{1/2})$ regret in expectation compared to
371 any policy and linear-time iterations. Upgrading to in-high-
372 probability $O(T^{1/2})$ regret guarantees while retaining linear-
373 time iterations remains an open problem. 374

375 **Interactive Bandit** We propose *interactive bandits* as a
376 natural online learning model to capture the essence of se-
377 quential decision processes. Here, an agent interacts with
378 the environment until a terminal state is reached, at which
379 point the payoff (a real number) is revealed to the agent. The
380 agent observes the environment’s action (signal) whenever
381 the interaction moves to an observation point. We formalize
382 this as an online learning model as follows. At all times t ,
383 before the agent acts, the environment privately selects (i) a
384 choice of payoff $u^t : Z \rightarrow \mathbb{R}$ for each terminal state $z \in Z$,
385 and (ii) a secret choice of signals $s_k^t \in S_k$ for all observa-
386 tion points $k \in \mathcal{K}$. These choices are hidden, and only the
387 signals relevant to the observation points reached during the
388 interaction will be revealed. Similarly, only the payoff $u^t(z^t)$
389 relative to the terminal state z^t reached in the interaction will
390 be revealed. In other words, the feedback that is revealed
391 to the agent after the interaction is the terminal state z^t that
392 is reached (which directly captures all signals revealed by
393 the environment, as they are the signals encountered on the
394 path from the root of the decision process to z^t) and its corre-
395 sponding payoff $u^t(z^t)$, which can be equivalently expressed
396 as $u^t(z^t) = (\ell^t)^\top \pi^t$, where the gradient vector ℓ^t is defined
397 as the (unique) vector such that for all strategies $x \in Q$,

$$(\ell^t)^\top x = \sum_{z \in Z} u^t(z) \left(\prod_{(k,s) \rightsquigarrow z} \mathbf{1}[s_k^t = s] \right) x[\sigma(z)]. \quad (1)$$

398 We assume that the environment is adversarial, in the sense
399 that the environment’s choices of payoffs u^t and signals s_k^t
400 at time t can depend on the previous actions of the agent.

401 Our term “interactive bandits” comes from the fact that an
402 algorithm for this setting can be thought of as interacting with
403 the environment until a terminal state of the decision process
404 is reached and a corresponding payoff is revealed. In other
405 words, while for modeling purposes it is convenient to think
406 about online learning algorithms as outputting strategies π
407 for the *whole* strategy space, one can think of an interactive
408 bandits algorithm as one that instead only outputs one action

361 ³Several algorithms are able to guarantee one or two out of the
362 three requirements (i) applicable to any convex domain, (ii) poly-
363 nomial time per iteration, (iii) $\tilde{O}(T^{1/2})$ regret with high probability.
364 For example, Bartlett et al. (2008) achieve (i) and (iii) by extending
365 an earlier paper by Dani, Kakade, and Hayes (2008), and (Györfgy
366 et al. 2007) achieves (ii) for the set of flows with suboptimal regret
367 guarantees.

²In online learning, “full-feedback” is typically called “full-
information”. We use “full-feedback” to avoid confusion with full-
information games, that is, games where the full state is available to
all players at all times.

at a time as the interaction moves throughout the decision process.

Since the interactive bandit model requires that the gradient vector ℓ^t be expressible as in (1), it makes more assumptions on the gradient vector than either the bandit linear optimization model or the full-feedback model. In terms of feedback, it is an intermediate model: it receives a superset of the feedback that the bandit linear optimization framework receives, but significantly less than the full-feedback model. So, theoretically, one could always use an algorithm for the bandit linear optimization model to solve a problem in the interactive bandit model. However, as we show in this paper, one can design a regret-minimization algorithm for the interactive bandit model that achieves sublinear $O(T^{3/4})$ regret with high probability *even in decision processes and extensive-form games whose decision space is at first unknown*. To our knowledge, no algorithm guaranteeing sublinear regret when the decision space is at first unknown has been designed for bandit linear optimization.

4 Algorithm for Unknown Sequential Decision Making Problems

We now describe a regret minimizer for the interactive bandit model. At all time t , it goes through two phases: first, the *rollout* phase, and then the *regret update* phase. During the rollout phase, the algorithm plays until a terminal state z^t and its corresponding payoff u^t are revealed. During the regret update phase, the algorithm rewinds through the decision nodes j encountered during that trajectory, and updates certain parameters at those decision nodes based on the newly-observed payoff u^t .⁴ Like the CFR framework, our algorithm picks actions at each decision node j by means of a *local* full-information regret minimizer \mathcal{R}_j for the action set A_j at that decision node.

Rollout Phase: Playing the Game

We describe two alternative algorithms for the rollout phase (namely the “upfront-flipping” and the “on-path-flipping” rollout variants), which differ in the way the mix exploration and exploitation. Both variants are theoretically sound, and yield to the same sublinear in-high-probability regret bound (Section 5), even when different variants are used at different times t while playing against the adversarial opponent.

We start from the “upfront-flipping” variant, which is arguably the conceptually simpler variant, although we find it to usually perform worse in practice.

Upfront-Flipping Rollout Variant When the *upfront-flipping rollout variant* is used at time t , at the beginning of the rollout phase and before any action is played, a biased coin is tossed to decide the algorithm to use to play out the interaction:

⁴While our algorithm shares many of the building blocks on Monte Carlo Tree Search (MCTS)—incremental tree-building, back-propagation, rollouts—it is *not* an anytime search algorithm, at least not in the sense of traditional game-tree search like the one employed by the Online Outcome Sampling algorithm (Lisý, Lanctot, and Bowling 2015).

• With probability β^t , the EXPLORE routine is used. It ignores the recommendations of the local regret minimizers at each decision node and instead plays according to the *exploration function* $h^t : \Sigma \rightarrow \mathbb{R}_{>0}$ (more details are below). In particular, at every decision node j encountered during the rollout, the agent picks action $a \in A_j$ at random according to the distribution $h^t(j, a) / (\sum_{a' \in A_j} h^t(j, a'))$.

• With probability $1 - \beta^t$, the EXPLOIT routine is used. With this routine, at every decision node j encountered during the rollout, the agent picks a decision by sampling from the distribution $x_j^t \in \Delta^{|A_j|}$ recommended by the regret minimizer \mathcal{R}_j .

In both cases, a regret minimizer \mathcal{R}_j for decision node j is created when j is first discovered.

Our upfront-flipping rollout strategy differs from the ϵ -greedy strategy in that the coin is tossed for the entire trajectory, not at each decision point.

On-Path-Flipping Rollout Variant In the *on-path-flipping* rollout variant, there is no single coin toss to distinguish between exploration and exploitation, and the two are interleaved throughout the rollout. Before any action is picked, the two *reach* quantities r^t, \hat{r}^t are both set to 1. Then, the rollout phase begins, and eventually the agent will be required to make a decision (pick an action) at some decision point j . Let $x_j^t \in \Delta^{|A_j|}$ be the distribution over actions recommended by the regret minimizer \mathcal{R}_j . In the on-path-flipping rollout variant, the agent picks an action $a \in A_j$ at j with probability proportional to

$$(1 - \beta^t)r^t \cdot x_j^t[a] + \beta^t \hat{r}^t \cdot \frac{h^t(j, a)}{\sum_{a' \in A_j} h^t(j, a')}.$$

Let a^* be the chosen action; r^t and \hat{r}^t are updated according to the formulas $r^t \leftarrow r^t \cdot x_j^t[a^*]$ and $\hat{r}^t \leftarrow \hat{r}^t \cdot \frac{h^t(j, a^*)}{\sum_{a' \in A_j} h^t(j, a')}$.

The agent keeps using this specific way of selecting actions and updating the reach quantities r, \hat{r} for all decision points encountered during the rollout.

The role of h^t In both variants, the role of the exploration function h^t is to guide exploration of different parts of the decision process.⁵ The optimal choice for h^t is to have $h^t(j, a)$ measure the number of terminal states in the subtree rooted at (j, a) . When this information is not available, a heuristic can be used instead. If no sensible heuristic can be devised, the uniform exploration strategy $h^t(j, a) = 1$ for all $(j, a) \in \Sigma$ is always a valid fallback. In Theorem 1 below, we give guarantees about the regret cumulated by our algorithm that apply to any $h^t : \Sigma \rightarrow \mathbb{R}_{>0}$.

Regret Update Phase: Propagating the Payoff up the Tree

In the regret update phase, the revealed feedback (that is, the revealed utility $u^t(z^t)$ and the terminal state z^t that was

⁵Despite the positive exploration term induced by h^t , it is not guaranteed that all decision points will be discovered as $T \rightarrow \infty$ as the adversarial environment might prevent so. Nonetheless, our algorithm guarantees sublinear regret with high probability.

reached in the rollout phase) is used to construct suitable *local* gradient vectors ℓ_j^t for each of the local regret minimizers \mathcal{R}_j on the path from the terminal state z^t up to the root. Let $(j_1, a_1), \dots, (j_m, a_m)$ be the sequence of decision nodes and actions that were played, in order, that ultimately led to terminal state z^t during the repetition of the game at time t . We start to construct local gradient vectors from decision node j_m , where we set the three quantities

$$\gamma^t := (1 - \beta^t) \cdot \prod_{i=1}^m x_{j_i}^t[a_i] + \beta^t \cdot \prod_{i=1}^m \frac{h^t(j_i, a_i)}{\sum_{a' \in A_{j_i}} h^t(j_i, a')},$$

$$\hat{u}_{j_m}^t := \frac{u^t(z^t)}{\gamma^t}, \quad \ell_{j_m}^t := \hat{u}_{j_m}^t e_{a_m},$$

where we used the notation $e_{a_m} \in \Delta^{A_j}$ to denote the a_m -th canonical vector, that is the vector whose components are all 0 except for the a_m -th entry, which is 1. Then, for all $i = 1, \dots, m-1$, we recursively let $\hat{u}_{j_i}^t := x_{j_i}^t[a_i] \cdot \hat{u}_{j_{i+1}}^t$, $\ell_{j_i}^t := \hat{u}_{j_i}^t e_{a_i}$. Finally, for all $i = 1, \dots, m$, the gradient vector $\ell_{j_i}^t$ is revealed as feedback to the local full-feedback regret minimizer \mathcal{R}_{j_i} at decision point j_i .

Average Policy

When regret minimizers are used to solve a convex-concave saddle point problem (such as a Nash equilibrium in a two-player zero-sum game), only the profile of *average* policies produced by the regret minimizers are guaranteed to converge to the saddle point. For this reason, it is crucial to be able to represent the *average* policy of an agent. Since we are assuming that the structure of the decision problem is only partially known, this operation requires more care in our setting. As we now show, it is possible to modify the algorithm so that the average policy can be extracted.

In order to maintain the average policy, we maintain an additional vector \bar{x}_j at each discovered decision node j . Intuitively, these vector will be populated with entries from the cumulative sum of all partial sequence-form strategies recommended so far by the \mathcal{R}_j 's. As soon as j is discovered for the first time (say, time t), we create its \bar{x}_j . If j 's parent sequence is the empty sequence (that is, j is one of the possible first decision nodes in the SDM process, i.e., j is preceded only by observation nodes), we simply set $\bar{x}_j[a]^t := t/|A_j|$ for all $a \in A_j$. Otherwise, let $p_j = (j', a')$ be the parent sequence of j , and we set $\bar{x}_j[a]^t := \bar{x}_{j'}[a']/|A_j|$ for all $a \in A_j$. Then, at all times t , after the feedback has been received but before the regret update phase has started, we introduce a new *average policy update phase*. In it, we iterate through all the decision nodes j_i that have been discovered so far (that is, the union of all decision points discovered up to time t), in the order they have been discovered. For each of them, we update \bar{x}_{j_i} according to the following rule. Let x_{j_i} be the policy strategy returned by the local full-feedback regret minimizer \mathcal{R}_{j_i} . If j_i 's parent sequence is the empty sequence, we simply set $\bar{x}_{j_i}^{t+1} := \bar{x}_{j_i}^t + x_{j_i}^t$ and $r_{j_i}^t := x_{j_i}^t$. Otherwise, let $p_{j_i} = (j', a')$ be the parent sequence of j_i , and we set $\bar{x}_{j_i}^{t+1} := \bar{x}_{j_i}^t + r_{j_i}^t[a'] \cdot x_{j_i}^t$, and $r_{j_i}^t[a] := r_{j_i}^t[a'] \cdot x_{j_i}^t[a]$ for all $a \in A_{j_i}$.

In order to play the average policy, it is enough to play actions proportionally to \bar{x}_j^{t+1} at all discovered decision nodes

j , and actions picked uniformly at random at undiscovered decision nodes.

Observation 1. *In all phases, the agent performs an amount of operations at most linear in the number of actions $|A_j|$ at each decision point j discovered up to time t . So, the average policy update phase requires work at most linear in the size of the underlying SDM.*

5 Guarantees on Regret

We now present regret bounds that hold at all times T for our algorithm from Section 4. The bounds do not depend on which rollout variant is chosen at each time t , and different variants can be chosen freely at different times.

Theorem 1 provides a general bound that holds for any choice of local regret minimizers \mathcal{R}_j ($j \in \mathcal{J}$), (non-increasing) stepsizes β^t , and explorations functions h^t at all times t . It will then be the basis for Theorem 2, which provides a way to set the algorithm parameters to achieve sublinear $O(T^{3/4})$ regret.

Theorem 1. *Let $R_j^T(\hat{\pi}_j)$ denote the regret cumulated by the local full-feedback regret minimizer \mathcal{R}_j compared to a generic strategy $\hat{\pi}_j$, and Δ be the maximum range of payoffs that can be selected by the environment at all times. Then, for all $T \geq 1$, $p \in (0, 1)$, and $\hat{\pi} \in \mathcal{Q}$, with probability at least $1 - p$ the regret cumulated by our algorithm satisfies*

$$R^T(\hat{\pi}) \leq \max_{\mathbf{q} \in \Pi} \left\{ \sum_{j \in \mathcal{J}} q[p_j] \cdot \max_{\hat{\pi}_j \in \Delta^{|A_j|}} R_j^T(\hat{\pi}_j) \right\} + \frac{\Delta}{\beta^T} (1 + \nu) \sqrt{2T \log \frac{2}{p}} + \Delta \sum_{t=1}^T \beta^t,$$

where

$$\nu := \sqrt{\frac{1}{T} \sum_{t=1}^T \max_{z \in \mathcal{Z}} \left\{ \prod_{(j,a) \rightsquigarrow z} \left(\frac{\sum_{a' \in A_j} h^t(j, a')}{h^t(j, a)} \right)^2 \right\}}. \quad (2)$$

(Proofs are in the appendix.) When h^t measures the number of terminal states reachable under any sequence (j, a) , the constant ν satisfies $\nu \leq |\Sigma| - 1$. For the uniform exploration function ($h^t(j, a) = 1$ for all j, a), ν is upper bounded by the product of the number of actions at all decision nodes, a polynomial quantity in ‘‘bushy’’ decision processes. In Theorem 2 we operationalize Theorem 1 by showing sensible choices of stepsizes β^t and local regret minimizers \mathcal{R}_j .

Theorem 2. *Let the local full-feedback regret minimizers \mathcal{R}_j guarantee $O(T^{1/2})$ regret in the worst case⁶, and let $p \in (0, 1)$. Furthermore, let the exploration probabilities be $\beta^t := \min\{1, k \cdot t^{-1/4}\}$ for all t , where $k > 0$ is an arbitrary constant. Then, there exists a (decision-problem-dependent) constant c independent of p and T such that for all $T \geq 1$*

$$\mathbb{P} \left[\max_{\hat{\pi} \in \Pi} R^T(\hat{\pi}) \leq c \cdot T^{3/4} \Delta \sqrt{\log \frac{2}{p}} \right] \geq 1 - p.$$

⁶Valid choices include the following algorithms: regret matching (Hart and Mas-Colell 2000), regret matching⁺ (Tammelin et al. 2015), follow-the-regularized-leader, online mirror descent, exponential weights, hedge, and others.

597 When h^t is an exact measure of the number of terminal states,
 598 c is polynomial in $|\Sigma|$. Otherwise, it is linear in the constant
 599 ν defined in (2), which depends on the specific exploration
 600 functions used.

601 Since Theorem 2 guarantees sublinear regret with high
 602 probability, our algorithm can be used for all purposes de-
 603 scribed in Section 3, including computing an approximate
 604 Nash equilibrium in a two-player zero-sum extensive-form
 605 game whose structure is *a priori* unknown.

606 6 Empirical Evaluation

607 In our experiments, we used our algorithm to compute an
 608 approximate Nash equilibrium. We compared our method to
 609 established model-free algorithms from the multiagent rein-
 610 forcement learning and computational game theory literature
 611 for this setting: *neural fictitious self-play* (NFSP) (Heinrich
 612 and Silver 2016), the *policy gradient* (PG) approach of Sriniv-
 613 asan et al. (2018), and the *online variant of Monte-Carlo*
 614 *CFR* (*online MCCFR*) mentioned in (Lanctot et al. 2009). In
 615 line with prior empirical evaluations of those methods, we
 616 compare the algorithms on two standard benchmark games:
 617 Kuhn poker (Kuhn 1950) and Leduc poker (Southey et al.
 618 2005). The games are reviewed in Appendix C.

619 We used the implementations of PG and NFSP provided in
 620 *OpenSpiel* (Lanctot et al. 2019). They internally use Tensor-
 621 flow. For PG, we tested the RPG and QPG policy gradi-
 622 ent formulations, but not the RM formulation (it performed worst in
 623 the original paper (Srinivasan et al. 2018)). We implemented
 624 online MCCFR and our algorithm in C++ (online MCCFR
 625 is not implemented in OpenSpiel). We ran every algorithm
 626 with five random seeds. In the figures below, we plot the
 627 average exploitability (a standard measure of closeness to
 628 equilibrium) of the players averaged across the five seeds.
 629 The shaded areas indicate the *maximum* and *minimum* over
 630 the five random seeds. For NFSP we used the hyperparameter
 631 recommended by the OpenSpiel implementation. For Kuhn
 632 poker, we used the settings for PG that were tuned and found
 633 to work the best by Srinivasan et al. (2018)—they are pub-
 634 licly available through OpenSpiel. For PG in Leduc poker,
 635 we performed a hyperparameter search and selected for the
 636 two PG plot (RPG and QPG formulation) the best combina-
 637 tion hyperparameters (full details are in the appendix). For
 638 both online MCCFR and our algorithm, we used RM+ (Tamelin
 639 2014) as the local (full-feedback) regret minimizer. For our
 640 algorithm, we only show performance for the on-
 641 path-flipping variant. The upfront-flipping variant performed
 642 significantly worse and data is available in the appendix. We
 643 tested $k \in \{0.5, 1, 10, 20\}$ and set h^t to either the uniform ex-
 644 ploration function (h^t constant) or the theoretically-optimal
 645 exploration function h^t that measures the number of terminal
 646 nodes as explained in Section 4. The performance of the two
 647 exploration functions was nearly identical, so in Figure 1 we
 648 show our algorithm with the uniform exploration function.
 649 We chose $k = 10$ since that performed well on both games.
 650 The plots for all other hyperparameter combinations for our
 651 algorithm are in the appendix. For online MCCFR, the only
 652 hyperparameter is ϵ , which controls the ϵ -greediness of the
 653 exploration term added before sampling and outputting the

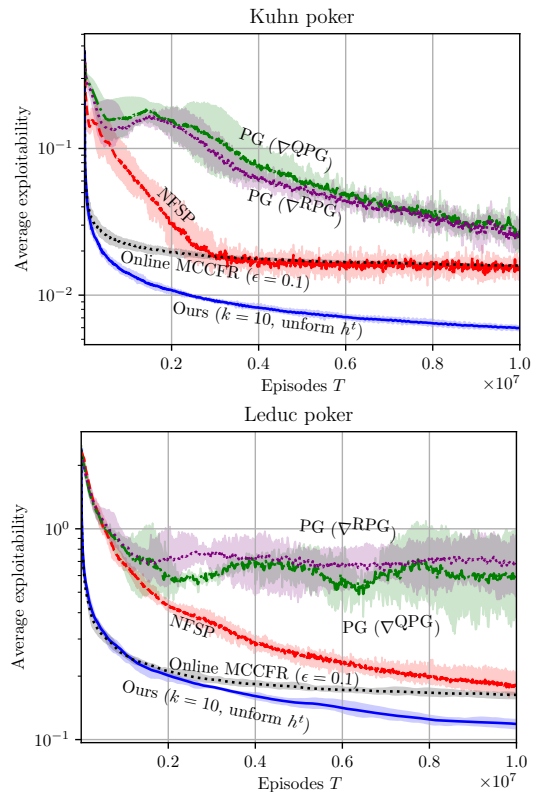


Figure 1: Comparison of the algorithms.

654 strategy at each time t . We tested $\epsilon = 0.6$ (which was found
 655 useful for the different, *self-play* MCCFR algorithm (Lanctot
 656 et al. 2009)), 0.1, and 0.0 (which corresponds to pure
 657 exploitation); Figure 1 shows the setting that performed best.

658 Out of all the algorithms, ours is the only one that guar-
 659 antees sublinear regret with high probability. This superior
 660 guarantee appears to translate into superior practical perfor-
 661 mance as well. In both benchmark games, our algorithm has
 662 lowest exploitability, often by a factor 2x-4x.

663 7 Conclusions and Future Research

664 We introduced a new online learning model, which we coined
 665 the *interactive bandit model*, that captures sequential deci-
 666 sion making. We developed an algorithm that guarantees
 667 sublinear $O(T^{3/4})$ regret with high probability in this model,
 668 even when the structure of the underlying decision problem
 669 is at first unknown to the agent and must be explored as
 670 part of the learning process. This is, to our knowledge, the
 671 first in-high-probability regret minimizer for this setting. It
 672 can be used for multiagent reinforcement learning. Its regret
 673 guarantee enables it to be used in any application for which
 674 regret minimization is useful: approximating Nash equilib-
 675 rium or quantal response equilibrium (Ling, Fang, and Kolter
 676 2018; Farina, Kroer, and Sandholm 2019) in two-player zero-
 677 sum games, approximating coarse correlated equilibrium in
 678 multi-player games (Moulin and Vial 1978; Celli et al. 2020),
 learning a best response, safe opponent exploitation (Farina,
 Kroer, and Sandholm 2019), online play against an unknown
 opponent/environment, etc. It is open whether better than
 $O(T^{3/4})$ regret is achievable in this important setting.

References

- 683
- 684 Abernethy, J.; Hazan, E.; and Rakhlin, A. 2008. Competing in
685 the dark: An efficient algorithm for bandit linear optimization.
686 In *In Proceedings of the 21st Annual Conference on Learning
687 Theory (COLT)*.
- 688 Abernethy, J. D.; and Rakhlin, A. 2009. Beating the adaptive
689 bandit with high probability. *2009 Information Theory and
690 Applications Workshop* .
- 691 Auer, P.; Cesa-Bianchi, N.; Freund, Y.; and Schapire, R. E.
692 2002. The Nonstochastic Multiarmed Bandit Problem. *SIAM
693 Journal of Computing* 32: 48–77.
- 694 Awerbuch, B.; and Kleinberg, R. D. 2004. Adaptive routing
695 with end-to-end feedback: distributed learning and geometric
696 approaches. In *Proceedings of the Annual Symposium on
697 Theory of Computing (STOC)*. ACM.
- 698 Azuma, K. 1967. Weighted sums of certain dependent ran-
699 dom variables. *Tohoku Mathematical Journal* 19(3): 357–
700 367.
- 701 Bartlett, P. L.; Dani, V.; Hayes, T.; Kakade, S.; Rakhlin, A.;
702 and Tewari, A. 2008. High-probability regret bounds for
703 bandit online linear optimization. In *Conference on Learning
704 Theory (COLT)*.
- 705 Bowling, M.; Burch, N.; Johanson, M.; and Tammelin, O.
706 2015. Heads-up Limit Hold'em Poker is Solved. *Science*
707 347(6218).
- 708 Brown, G. W. 1951. Iterative Solutions of Games by Fic-
709 titious Play. In Koopmans, T. C., ed., *Activity Analysis of
710 Production and Allocation*, 374–376. John Wiley & Sons.
- 711 Brown, N.; Kroer, C.; and Sandholm, T. 2017. Dynamic
712 Thresholding and Pruning for Regret Minimization. In *AAAI
713 Conference on Artificial Intelligence (AAAI)*.
- 714 Brown, N.; and Sandholm, T. 2015. Regret-Based Pruning in
715 Extensive-Form Games. In *Proceedings of the Annual Con-
716 ference on Neural Information Processing Systems (NIPS)*.
- 717 Brown, N.; and Sandholm, T. 2017a. Reduced Space and
718 Faster Convergence in Imperfect-Information Games via
719 Pruning. In *International Conference on Machine Learn-
720 ing (ICML)*.
- 721 Brown, N.; and Sandholm, T. 2017b. Superhuman AI for
722 heads-up no-limit poker: Libratus beats top professionals.
723 *Science* eaa01733.
- 724 Brown, N.; and Sandholm, T. 2019a. Solving imperfect-
725 information games via discounted regret minimization. In
726 *AAAI Conference on Artificial Intelligence (AAAI)*.
- 727 Brown, N.; and Sandholm, T. 2019b. Superhuman AI for
728 multiplayer poker. *Science* 365(6456): 885–890.
- 729 Bubeck, S.; Lee, Y. T.; and Eldan, R. 2017. Kernel-based
730 methods for bandit convex optimization. In *Proceedings
731 of the Annual Symposium on Theory of Computing (STOC)*,
732 72–85.
- 733 Celli, A.; Marchesi, A.; Bianchi, T.; and Gatti, N. 2020.
734 Learning to Correlate in Multi-Player General-Sum Sequen-
735 tial Games.
- Dani, V.; Kakade, S. M.; and Hayes, T. P. 2008. The Price of
736 Bandit Information for Online Optimization. In *Proceedings
737 of the Annual Conference on Neural Information Processing
738 Systems (NIPS)*. 739
- Even-Dar, E.; Kakade, S. M.; and Mansour, Y. 2009. On-
740 line Markov decision processes. *Mathematics of Operations
741 Research* 34(3): 726–736. 742
- Farina, G.; Kroer, C.; and Sandholm, T. 2019. Online Con-
743 vex Optimization for Sequential Decision Processes and
744 Extensive-Form Games. In *AAAI Conference on Artificial
745 Intelligence*. 746
- Farina, G.; Kroer, C.; and Sandholm, T. 2020. Stochastic
747 Regret Minimization in Extensive-Form Games. In *ArXiv
748 preprint*. 749
- Farina, G.; Schmucker, R.; and Sandholm, T. 2020. Counterfactual-Free Regret Minimization for Sequential Decision Making and Extensive-Form Games. In *Reinforcement Learning in Games (LRG) Workshop at AAAI 2020*. 750
- Flaxman, A.; Kalai, A.; and McMahan, B. 2005. Online
751 Convex Optimization in the Bandit Setting: Gradient Descent
752 Without a Gradient. In *Annual ACM-SIAM Symposium on
753 Discrete Algorithms (SODA)*. 754
- György, A.; Linder, T.; Lugosi, G.; and Ottucsák, G. 2007.
755 The on-line shortest path problem under partial monitoring.
756 *Journal of Machine Learning Research* 8(Oct): 2369–2403. 757
- Hart, S.; and Mas-Colell, A. 2000. A Simple Adaptive Pro-
758 cedure Leading to Correlated Equilibrium. *Econometrica* 68:
759 1127–1150. 760
- Hazan, E. 2016. Introduction to online convex optimization.
761 *Foundations and Trends in Optimization* 2(3-4): 157–325. 762
- Heinrich, J.; Lanctot, M.; and Silver, D. 2015. Fictitious self-
763 play in extensive-form games. In *International Conference
764 on Machine Learning*, 805–813. 765
- Heinrich, J.; and Silver, D. 2016. Deep reinforcement learn-
766 ing from self-play in imperfect-information games. *arXiv
767 preprint arXiv:1603.01121* . 768
- Hoëffding, W. 1963. Probability Inequalities for Sums of
769 Bounded Random Variables. *Journal of the American Statis-
770 tical Association* 58(301): 13–30. 771
- Kash, I. A.; Sullins, M.; and Hofmann, K. 2020. Combining
772 No-regret and Q-learning. In *Autonomous Agents and Multi-
773 Agent Systems*. 774
- Kleinberg, R. 2004. Nearly tight bounds for the continuum-
775 armed bandit problem. In *Proceedings of the Annual Confer-
776 ence on Neural Information Processing Systems (NIPS)*. 777
- Koller, D.; Megiddo, N.; and von Stengel, B. 1994. Fast
778 algorithms for finding randomized strategies in game trees.
779 In *Proceedings of the 26th ACM Symposium on Theory of
780 Computing (STOC)*. 781
- Kuhn, H. W. 1950. A Simplified Two-Person Poker. In Kuhn,
782 H. W.; and Tucker, A. W., eds., *Contributions to the Theory
783 of Games*, volume 1 of *Annals of Mathematics Studies*, 24,
784 97–103. Princeton, New Jersey: Princeton University Press. 785

- 789 Lanctot, M.; Lockhart, E.; Lespiau, J.-B.; Zambaldi, V.;
790 Upadhyay, S.; Pérolat, J.; Srinivasan, S.; Timbers, F.; Tuyls,
791 K.; Omidshafiei, S.; Hennes, D.; Morrill, D.; Muller, P.;
792 Ewalds, T.; Faulkner, R.; Kramár, J.; Vylder, B. D.; Saeta,
793 B.; Bradbury, J.; Ding, D.; Borgeaud, S.; Lai, M.; Schrit-
794 twieser, J.; Anthony, T.; Hughes, E.; Danihelka, I.; and Ryan-
795 Davis, J. 2019. OpenSpiel: A Framework for Reinforce-
796 ment Learning in Games. *CoRR* abs/1908.09453. URL
797 <http://arxiv.org/abs/1908.09453>.
- 798 Lanctot, M.; Waugh, K.; Zinkevich, M.; and Bowling, M.
799 2009. Monte Carlo Sampling for Regret Minimization in
800 Extensive Games. In *Proceedings of the Annual Conference*
801 *on Neural Information Processing Systems (NIPS)*.
- 802 Lattimore, T.; and Szepesvári, C. 2020. *Bandit Algorithms*.
803 Cambridge University Press. ISBN 9781108486828.
- 804 Ling, C. K.; Fang, F.; and Kolter, J. Z. 2018. What game
805 are we playing? End-to-end learning in normal and exten-
806 sive form games. In *Proceedings of the International Joint*
807 *Conference on Artificial Intelligence (IJCAI)*.
- 808 Lisý, V.; Lanctot, M.; and Bowling, M. 2015. Online monte
809 carlo counterfactual regret minimization for search in imper-
810 fect information games. In *Proceedings of the 2015 inter-
811 national conference on autonomous agents and multiagent*
812 *systems*, 27–36.
- 813 Matignon, L.; Laurent, G. J.; and Le Fort-Piat, N. 2012. Inde-
814 pendent reinforcement learners in cooperative markov games:
815 a survey regarding coordination problems. *The Knowledge*
816 *Engineering Review* 27(1): 1–31.
- 817 Moravčík, M.; Schmid, M.; Burch, N.; Lisý, V.; Morrill, D.;
818 Bard, N.; Davis, T.; Waugh, K.; Johanson, M.; and Bowling,
819 M. 2017. DeepStack: Expert-level artificial intelligence in
820 heads-up no-limit poker. *Science* .
- 821 Moulin, H.; and Vial, J.-P. 1978. Strategically zero-sum
822 games: The class of games whose completely mixed equilib-
823 ria cannot be improved upon. *International Journal of Game*
824 *Theory* 7(3-4): 201–221.
- 825 Neu, G.; Antos, A.; György, A.; and Szepesvári, C. 2010.
826 Online Markov decision processes under bandit feedback. In
827 *Advances in Neural Information Processing Systems*, 1804–
828 1812.
- 829 Ponsen, M.; De Jong, S.; and Lanctot, M. 2011. Computing
830 approximate nash equilibria and robust best-responses using
831 sampling. *Journal of Artificial Intelligence Research* 42:
832 575–605.
- 833 Romanovskii, I. 1962. Reduction of a Game with Complete
834 Memory to a Matrix Game. *Soviet Mathematics* 3.
- 835 Sandholm, T.; and Crites, R. 1996. Multiagent Reinforcement
836 Learning in the Iterated Prisoner’s Dilemma. *Biosystems* 37:
837 147–166. Special issue on the Prisoner’s Dilemma. Early
838 version in IJCAI-95 Workshop on Adaptation and Learning
839 in Multiagent Systems.
- 840 Shalev-Shwartz, S. 2012. Online Learning and Online Con-
841 vex Optimization. *Foundations and Trends® in Machine*
842 *Learning* 4(2). ISSN 1935-8237. doi:10.1561/22000000018.
- 843 Southey, F.; Bowling, M.; Larson, B.; Piccione, C.; Burch,
844 N.; Billings, D.; and Rayner, C. 2005. Bayes’ Bluff: Oppo-
845 nent Modelling in Poker. In *Proceedings of the 21st Annual*
846 *Conference on Uncertainty in Artificial Intelligence (UAI)*.
- 847 Srinivasan, S.; Lanctot, M.; Zambaldi, V.; Pérolat, J.; Tuyls,
848 K.; Munos, R.; and Bowling, M. 2018. Actor-Critic Policy
849 Optimization in Partially Observable Multiagent Environ-
850 ments.
- 851 Strens, M. 2000. A Bayesian framework for reinforcement
852 learning. In *International Conference on Machine Learning*
853 *(ICML)*.
- 854 Tammelin, O. 2014. Solving large imperfect information
855 games using CFR+. *arXiv preprint arXiv:1407.5042* .
- 856 Tammelin, O.; Burch, N.; Johanson, M.; and Bowling, M.
857 2015. Solving Heads-up Limit Texas Hold’em. In *Proceed-*
858 *ings of the 24th International Joint Conference on Artificial*
859 *Intelligence (IJCAI)*.
- 860 von Stengel, B. 1996. Efficient Computation of Behavior
861 Strategies. *Games and Economic Behavior* 14(2): 220–246.
- 862 Zhou, Y.; Li, J.; and Zhu, J. 2020. Posterior sampling for
863 multi-agent reinforcement learning: solving extensive games
864 with imperfect information. In *Conference on Learning Rep-*
865 *resentations*.
- 866 Zimin, A.; and Neu, G. 2013. Online learning in episodic
867 Markovian decision processes by relative entropy policy
868 search. In *Advances in neural information processing sys-*
869 *tems*, 1583–1591.
- 870 Zinkevich, M. 2003. Online Convex Programming and Gen-
871 eralized Infinitesimal Gradient Ascent. In *International Con-*
872 *ference on Machine Learning (ICML)*, 928–936. Washington,
873 DC, USA.
- 874 Zinkevich, M.; Bowling, M.; Johanson, M.; and Piccione,
875 C. 2007. Regret Minimization in Games with Incomplete
876 Information. In *Proceedings of the Annual Conference on*
877 *Neural Information Processing Systems (NIPS)*.

A Summary of Online Learning Models

Online learning model	Assumptions on gradient vector ℓ^t (other than bounded norm)	Revealed feedback
Full-feedback ²	unconstrained: $\ell^t \in \mathbb{R}^{ \Sigma }$	ℓ^t
Bandit linear optimization	unconstrained: $\ell^t \in \mathbb{R}^{ \Sigma }$	$(\ell^t)^\top \pi^t \in \mathbb{R}$
Interactive bandit (this paper)	$u^t : Z \rightarrow \mathbb{R}$ choice of payoff at each $z \in Z$ $s_k^t \in S_k$ choice of signal at each $k \in \mathcal{K}$ ℓ^t as in Equation (3)	<ul style="list-style-type: none"> • $z^t \in Z$ terminal state, • $u^t(z^t) = (\ell^t)^\top \pi^t \in \mathbb{R}$ payoff at z^t

Table 1: Comparison between different online learning models and their assumptions on the gradient vector and revealed feedback.

Observation 2. In online learning, the environment picks a secret loss function (gradient vector) ℓ^t at every iteration t . The gradient vector controls the utility of the player, which is defined as the inner product $(\ell^t)^\top \pi^t$, where π^t is the strategy output by the agent at iteration t . In the interactive bandit model, it is perhaps more natural to think of the environment as picking a choice of signal s_k^t at each observation node $k \in \mathcal{K}$, as well as a choice of payoff $u^t(z)$ for each terminal node $z \in Z$. The utility of the agent is then computed as $u^t(z^t)$, where z^t is the (unique) terminal node that is reached when the agent plays according to $\pi^t \in \Pi$, and the environment sends the signals s_k^t .

The two points of view can be reconciled. In particular, at every iteration t it is possible to identify a gradient vector ℓ^t for the interactive bandit setting such that $u^t(z^t) = (\ell^t)^\top \pi^t$, as follows. As mentioned in the body, let $\sigma(z)$ be the last sequence (decision node-action pair) on the path from the root of the decision process to terminal node $z \in Z$, and consider the gradient vector ℓ^t , defined as the (unique) vector such that

$$(\ell^t)^\top \mathbf{x} = \sum_{z \in Z} u^t(z) \left(\prod_{(k,s) \rightsquigarrow z} \mathbb{1}[s_k^t = s] \right) x[\sigma(z)] \quad \forall \mathbf{x} \in Q. \quad (3)$$

For each terminal node $z \in Z$ in the argument of the sum, the product in parenthesis is 1 if and only if signals s_k^t picked by the environment are a superset of those on the path from the root to z . Furthermore, when $\mathbf{x} = \pi^t$, $x[\sigma(z^t)] = 1$ if and only if the decision maker has decided to play all actions on the path from the root to z^t . So, the only terminal node for which the right hand side of Equation (3) is nonzero is $z = z^t$. Consequently, $(\ell^t)^\top \pi^t = u^t(z^t)$.

B Analysis and Proofs

In order to make the proof approachable, at first we will assume that the structure of the decision process is fully known. We will then present and analyze an algorithm that enjoys the properties described in Theorem 1. Finally, we will show that the algorithm can be implemented as described in Section 4, and therefore, that it does not actually need to know the structure of the decision process.

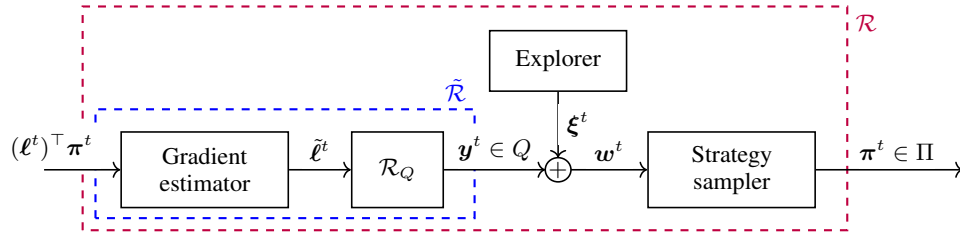


Figure 2: Conceptual construction of our algorithm for the interactive bandit online learning setting.

At a high level, the construction of our interactive bandit regret minimizer works as follows. At each iteration t , we use a full-feedback regret minimizer \mathcal{R}_Q to output a recommendation for the next sequence-form strategy $\mathbf{y}^t \in Q$ to be played. Then, we introduce a bias on \mathbf{y}^t , which can be thought of as an

904 exploration term. The resulting biased strategy is $\mathbf{w}^t \in Q$. Finally, we sample a deterministic policy
 905 $\boldsymbol{\pi}^t \in \Pi$ starting from \mathbf{w}^t . After the playthrough is over and a terminal node $z^t \in Z$ has been reached, we
 906 use the feedback (that is, the terminal node z^t reached in the playthrough, together with its utility $u^t(z^t)$),
 907 to construct an unbiased estimator $\tilde{\boldsymbol{\ell}}^t$ of the underlying gradient vector that was chosen by the environment.
 908 The algorithm that we have just described accumulates regret from the following three sources.

- 909 1. First, it cumulates regret by the full-feedback regret minimizer \mathcal{R}_Q . This scales with the norm of the
 910 gradient estimators $\tilde{\boldsymbol{\ell}}^t$ as $O(\sqrt{T} \cdot \max_{t=1}^T \|\tilde{\boldsymbol{\ell}}^t\|)$.
- 911 2. Second, there is a degradation term due to the fact that we do not exactly follow the recommendations
 912 from \mathcal{R}_Q , and instead bias the recommendations. The biasing is necessary to keep the norm of the
 913 gradient estimators under control.
- 914 3. Third, there is a regret degradation term due to the fact that we sample a pure strategy starting from the
 915 biased recommendation \mathbf{w}^t .

916 Our analysis is split into two main conceptual steps. First, we will quantify the regret degradation that is
 917 incurred in passing from $\tilde{\mathcal{R}}$ to \mathcal{R} (Figure 2). Specifically, we will study how the regret cumulated by $\tilde{\mathcal{R}}$,

$$\tilde{R}^T(\boldsymbol{\pi}) = \sum_{t=1}^T (\boldsymbol{\ell}^t)^\top (\boldsymbol{\pi} - \mathbf{y}^t)$$

918 is related to the regret cumulated by our overall algorithm,

$$R^T(\boldsymbol{\pi}) = \sum_{t=1}^T (\boldsymbol{\ell}^t)^\top (\boldsymbol{\pi} - \boldsymbol{\pi}^t).$$

919 In this case, the degradation term comes from the fact that the strategy output by \mathcal{R} , that is, $\boldsymbol{\pi}^t \in \Pi$, is not
 920 the same as the one, \mathbf{y}^t , that was recommended by $\tilde{\mathcal{R}}$, because an exploration term was added.

921 The second step in the analysis will be the quantification of the regret degradation that is incurred in
 922 passing from \mathcal{R}_Q to $\tilde{\mathcal{R}}$ (Figure 2). Specifically, we will study how the regret cumulated by $\tilde{\mathcal{R}}$ is related to
 923 the regret

$$\mathcal{R}_Q^T(\boldsymbol{\pi}) = \sum_{t=1}^T (\tilde{\boldsymbol{\ell}}^t)^\top (\boldsymbol{\pi} - \mathbf{y}^t),$$

924 which is guaranteed to be $O(T^{1/2})$ with high probability. By summing the degradation bounds, the
 925 asymptotic regret guarantee will change from $T^{1/2}$ to $T^{3/4}$, while still retaining a high-probability bound.

926 Below, we give details about each of the components of our algorithm.

927 **Gradient Estimator** The gradient estimator that we use in our algorithm is given in Lemma 1. It is a
 928 form of importance-sampling estimator that can be constructed starting from the feedback received by
 929 the regret minimizer (that is, the terminal leaf z^t and its corresponding payoff $u^t(z^t)$). It is a particular
 930 instantiation of the *outcome sampling* estimator that appeared in the works by Lanctot et al. (2009). We
 931 follow the formalization of Farina, Kroer, and Sandholm (2020).

932 **Lemma 1.** *Let z^t be the terminal state in which the process ends at iteration t , that is the terminal state
 933 reached when the decision maker plays according to the pure strategy $\boldsymbol{\pi}^t$ and the environment plays
 934 according to the (partially hidden) pure strategy \mathbf{y}^t . Also, let $u^t(z^t)$ be the reward of the decision maker
 935 upon termination and let $\sigma(z^t)$ denote the last sequence (decision node-action pair) on the path from the
 936 root of the decision process to z . Then,*

$$\tilde{\boldsymbol{\ell}}^t := \frac{u^t(z^t)}{w^t[\sigma(z^t)]} \mathbf{e}_{z^t}$$

937 *is an unbiased estimator of $\boldsymbol{\ell}^t$, that is, $\mathbb{E}_t[\tilde{\boldsymbol{\ell}}^t] = \boldsymbol{\ell}^t$.*

938 *Proof.* We will prove the lemma by showing that $\mathbb{E}_t[(\tilde{\boldsymbol{\ell}}^t)^\top \mathbf{x}] = (\boldsymbol{\ell}^t)^\top \mathbf{x}$. By definition, z^t is the (unique)
 939 state such that

$$\pi^t[\sigma(z)] = \prod_{(k,s) \rightsquigarrow z} \mathbb{1}[s_k^t = s] = 1.$$

940 For all other terminal states z , $\pi^t[\sigma(z)] \cdot \prod_{(k,s) \rightsquigarrow z} \mathbb{1}[s_k^t = s] = 0$. Hence,

$$\begin{aligned} \mathbb{E}_t \left[(\tilde{\ell}^t)^\top \mathbf{x} \right] &= \mathbb{E}_t \left[\frac{u^t(z^t)}{w^t[\sigma(z^t)]} (\mathbf{e}_z^\top \mathbf{x}) \right] \\ &= \mathbb{E}_t \left[\sum_{z \in \mathcal{Z}} \left(\pi^t[\sigma(z)] \cdot \prod_{(k,s) \rightsquigarrow z} \mathbb{1}[s_k^t = s] \right) \frac{u^t(z)}{w^t[\sigma(z)]} (\mathbf{e}_z^\top \mathbf{x}) \right] \\ &= \sum_{z \in \mathcal{Z}} \frac{\mathbb{E}_t[\pi^t[\sigma(z)]]}{w^t[\sigma(z)]} u^t(z) \left(\prod_{(k,s) \rightsquigarrow z} \mathbb{1}[s_k^t = s] \right) x[\sigma(z)], \end{aligned}$$

941 Using the hypothesis that π^t is a (conditionally) unbiased estimator of w^t we further obtain

$$\mathbb{E}_t[\pi^t[\sigma(z)]] = \mathbb{E}_t[\pi^t][\sigma(z)] = w^t[\sigma(z)]$$

942 and therefore

$$\mathbb{E}_t \left[(\tilde{\ell}^t)^\top \mathbf{x} \right] = \sum_{z \in \mathcal{Z}} u^t(z) \left(\prod_{(k,s) \rightsquigarrow z} \mathbb{1}[s_k^t = s] \right) x[\sigma(z)] = (\ell^t)^\top \mathbf{x},$$

943 where we used the definition of ℓ^t from (1) in the last equality. \square

944 **The Full-Feedback Regret Minimizer** Our full-feedback regret minimizer \mathcal{R}_Q is the counterfactual
 945 regret minimization (CFR) regret minimizer (Zinkevich et al. 2007). CFR decomposes the task of computing
 946 a strategy $\mathbf{y}^t \in Q$ for the whole decision process into smaller subproblems at each decision point. The
 947 local strategies are then computed via $|\mathcal{J}|$ independent full-feedback regret minimizers \mathcal{R}_j ($j \in \mathcal{J}$), one
 948 per decision node, that at each t observe a specific feedback—the *counterfactual gradient*—that guarantees
 949 low global regret across the decision process. CFR guarantees $O(T^{1/2})$ regret in the worst case against any
 950 strategy $\hat{\mathbf{y}}^t \in Q$.

951 We now review how the counterfactual gradients that are given as feedback at all local regret minimizers
 952 \mathcal{R}_j are constructed. Let $\mathbf{x}_j^t, j \in \mathcal{J}$, be the local strategies output by the \mathcal{R}_j 's, and let $\tilde{\ell}^t$ be the gradient
 953 vector observed at time t by \mathcal{R}_Q . We start by constructing the counterfactual values V_j , indexed over $j \in \mathcal{J}$,
 954 which are recursively defined as

$$V_j^t = \sum_{a \in A_j} x_j^t[a] \left(\tilde{\ell}^t[j, a] + \sum_{j' \in \mathcal{J}: p_{j'} = (j, a)} V_{j'}^t \right).$$

955 Then, for all $j \in \mathcal{J}$ we construct the counterfactual gradient ℓ_j^t by setting $\ell_j^t[a] = \tilde{\ell}^t[j, a] +$
 956 $\sum_{j' \in \mathcal{J}: p_{j'} = (j, a)} V_{j'}^t$ for all $a \in A_j$.

957 Let $R_j^T(\hat{\pi}_j)$ denote the regret cumulated up to time T by the local full-feedback regret minimizer \mathcal{R}_j
 958 compared to strategy $\hat{\pi}_j^t \in \Delta^{|A_j|}$. Then, the regret cumulated by the CFR algorithm is known to satisfy the
 959 regret bound in Lemma 2.

960 **Lemma 2** ((Farina, Kroer, and Sandholm 2019)). *At all T , for all $\hat{\pi} \in Q$,*

$$R_Q(\hat{\pi}) \leq \max_{\mathbf{q} \in Q} \left\{ \sum_{j \in \mathcal{J}} q[p_j] \cdot \max_{\hat{\pi}_j \in \Delta^{|A_j|}} R_j^T(\hat{\pi}_j) \right\}.$$

961 **Exploration Term** The role of the exploration term is to reduce the norm of the gradient estimators. We
 962 bias the strategies \mathbf{y}^t output by \mathcal{R}_Q by taking a convex combination with the *exploration strategy* $\xi^t \in Q$
 963 defined as the sequence-form strategy that picks action $a \in A_j$ at decision node $j \in \mathcal{J}$ with probability
 964 proportional to

$$\xi^t[j, a] \propto \frac{h^t(j, a)}{\sum_{a \in A_j} h^t(j, a)},$$

965 where h^t is a generic exploration function $\Sigma \rightarrow \mathbb{R}_{>0}$.

966 Let $\sigma(z)$ denote the last sequence (decision node-action pair) encountered on the path from the root of
 967 the decision process to terminal node $z \in Z$. A key quantity about ξ^t , that will be important to quantify the
 968 effect of the biasing on the norm of the gradient estimator, is the following:

$$\rho^t := \max_{z \in Z} \frac{1}{\xi^t[\sigma(z)]} = \max_{z \in Z} \prod_{(j,a) \rightsquigarrow z} \left(\frac{\sum_{a' \in A_j} h^t(j, a')}{h^t(j, a)} \right).$$

969 At all times t , we construct the biased strategy w^t starting from the unbiased strategy y^t and the
 970 exploration strategy ξ^t as

$$w^t := (1 - \beta^t) \cdot y^t + \beta^t \cdot \xi^t. \quad (4)$$

971 We call the exploration function h^t that assigns the constant value 1 to all sequences the *uniform*
 972 *exploration strategy*. We also study an exploration strategy for which we can prove better worst-case
 973 convergence speed. Specifically, we call the exploration function h^t that assigns, to each sequence $(j, a) \in$
 974 Σ , the number of terminal nodes under (j, a) the *balanced exploration strategy*. For this strategy, the
 975 minimum reach probability across all terminal nodes, which is important for our worst-case convergence
 976 guarantee, scales linearly in the game size.

977 **Lemma 3** ((Farina, Kroer, and Sandholm 2020)). *When h^t is a perfect measure of the number of terminal*
 978 *nodes under each sequence $(j, a) \in \Sigma$, then $\rho^t \leq |\Sigma| - 1$.*

979 **Strategy Sampler** The role of the strategy sampler is to provide an *unbiased* estimator $\pi^t \in \Pi$ of
 980 w^t . In fact, any unbiased estimator can be used in our framework. We start by describing a sampling
 981 scheme for sequence-form strategies that is folklore in the literature. In the natural sampling scheme,
 982 given a sequence-form strategy $q \in Q$, at all decision nodes $j \in \mathcal{J}$ an action is sampled according to the
 983 distribution $q[j, a]/q[p_j] : a \in A_j$.

984 • *Upfront-flipping strategy sampling* Since by construction $w^t = (1 - \beta^t) \cdot y^t + \beta^t \cdot \xi^t$, sampling from
 985 w^t can be done as follows. First, we flip a biased coin, where the probability of heads is $1 - \beta^t$ and
 986 the probability of tails is β^t . If heads comes up, we sample actions according to the natural sampling
 987 scheme applied to the sequence-form strategy y^t output by \mathcal{R}_Q . Otherwise, we sample actions according
 988 to the natural sampling scheme applied to the exploration strategy ξ^t . We call the strategy sampler just
 989 described the *upfront-flipping* sampling scheme.

990 • *On-path-flipping strategy sampling* We now describe a different sampling scheme, which we coin the
 991 *on-path-flipping* sampling scheme. First, we compute the strategy w^t by explicitly taking the convex
 992 combination between y^t and the exploration term ξ^t . Then we apply the natural sampling scheme
 993 for sequence-form strategies to w^t . In the on-path-flipping sampling scheme, the exploration and the
 994 exploitation are interleaved at each decision node in a counterintuitive way.

995 Because both sampling schemes are unbiased, our theory, without changes, applies to both of these sampling
 996 schemes. Later in this appendix, we report experiments on both of these sampling schemes.

997 Relationship Between Exploration and Norm of Counterfactual Gradients

998 The gradient $\tilde{\ell}^t$ that is given to \mathcal{R}_Q as feedback is given in Lemma 1. It has zero entries everywhere,
 999 except for the sequence $\sigma(z^t)$. Consequently, the counterfactual values V_j^t are 0 everywhere, except for
 1000 the decision nodes that were traversed on the path from the root to the terminal node $z^t \in Z$. In turn, this
 1001 means that all counterfactual gradients are 0, except for the sequences that are traversed. For them, the
 1002 construction of the counterfactual gradients reduces to the *regret update phase* described in Section 4.

1003 The norms of the counterfactual gradients are maximum at the leaves. In particular, the following is
 1004 known.

1005 **Lemma 4** ((Lanctot et al. 2009; Farina, Kroer, and Sandholm 2020)). *Let Δ be the maximum range of*
 1006 *payoffs that can be selected by the environment at all times. All counterfactual gradients have norm upper*
 1007 *bounded as*

$$\|\ell_j^t\|_2 \leq \frac{\Delta \rho^t}{\beta^t} \quad \forall j \in \mathcal{J}.$$

1008 Furthermore, the loss estimate $\tilde{\ell}^t$ satisfies

$$(\tilde{\ell}^t)^\top (x - x') \leq \frac{\Delta \rho^t}{\beta^t} \quad \forall x, x' \in Q.$$

1009 **Relationship Between \mathcal{R} and $\tilde{\mathcal{R}}$**

1010 As a first step in our analysis, we establish an important relationship between the regret cumulated by \mathcal{R}
 1011 and $\tilde{\mathcal{R}}$. Fundamentally, it quantifies the degradation in the regret incurred from playing π^t instead of \mathbf{y}^t as
 1012 recommended. The degradation is kept under control by the fact that the expectation of the output π^t is
 1013 close to \mathbf{y}^t when β is small. Before we state the central result of this subsection (Proposition 1), we need
 1014 the following simple bound.

1015 **Lemma 5.** *Let Δ be the payoff range of the interaction, that is the maximum absolute value of any payoff
 1016 that can be selected by the environment at any time. At all T , it holds that*

$$\sum_{t=1}^T (\ell^t)^\top (\mathbf{y}^t - \mathbf{w}^t) \leq \Delta \sum_{t=1}^T \beta^t.$$

1017 *Proof.* Using simple algebraic manipulations,

$$\begin{aligned} \sum_{t=1}^T (\ell^t)^\top (\mathbf{y}^t - \mathbf{w}^t) &= \sum_{t=1}^T (\ell^t)^\top (\mathbf{y}^t - ((1 - \beta^t)\mathbf{y}^t + \beta^t \mathbf{b}^t)) \\ &= \sum_{t=1}^T \beta^t (\ell^t)^\top (\mathbf{y}^t - \mathbf{b}^t) \leq \Delta \sum_{t=1}^T \beta^t, \end{aligned}$$

1018 where the last inequality follows from the definition of Δ . □

1019 With Lemma 5 we are ready to analyze the relationship between \mathcal{R} and $\tilde{\mathcal{R}}$.

1020 **Proposition 1.** *At all T , for all $p \in (0, 1)$ and $\pi \in \Pi$, with probability at least $1 - p$,*

$$R^T(\pi) \leq \tilde{R}^T(\pi) + \Delta \left(\sqrt{2T \log \frac{1}{p}} + \sum_{t=1}^T \beta^t \right).$$

1021 *Proof.* Introduce the discrete-time stochastic process

$$d^t := (\ell^t)^\top (\pi^t - \mathbf{w}^t), \quad t \in \{1, 2, \dots\}.$$

1022 From the online learning model hypotheses, we have that ℓ^t is conditionally independent from π^t and \mathbf{w}^t
 1023 given all past choices of the algorithm up to time $t - 1$, and since $\mathbb{E}_t[\pi^t] = \mathbf{w}^t$ by construction, d^t is a
 1024 martingale difference sequence. Furthermore, some elementary algebra reveals that

$$\begin{aligned} \sum_{t=1}^T d^t &= \sum_{t=1}^T (\ell^t)^\top (\pi^t - \mathbf{w}^t) \\ &= \sum_{t=1}^T (\ell^t)^\top (\pi - \mathbf{y}^t) - \sum_{t=1}^T (\ell^t)^\top (\pi - \pi^t) + \sum_{t=1}^T (\ell^t)^\top (\mathbf{y}^t - \mathbf{w}^t) \\ &= \tilde{R}^T(\pi) - R^T(\pi) + \sum_{t=1}^T (\ell^t)^\top (\mathbf{y}^t - \mathbf{w}^t) \\ &\leq \tilde{R}^T(\pi) - R^T(\pi) + \Delta \sum_{t=1}^T \beta^t, \end{aligned} \tag{5}$$

1025 where the last inequality follows from Lemma 5. Since $|\delta^t| \leq \Delta$ for all t , using the Azuma-Hoeffding
 1026 inequality (Hoeffding 1963; Azuma 1967) we have that

$$\begin{aligned} 1 - p &\leq \mathbb{P} \left[\sum_{t=1}^T d^t \geq -\Delta \sqrt{2T \log \frac{1}{p}} \right] \\ &\leq \mathbb{P} \left[\tilde{R}^T(\pi) - R^T(\pi) + \Delta \sum_{t=1}^T \beta^t \geq -\Delta \sqrt{2T \log \frac{1}{p}} \right] \\ &= \mathbb{P} \left[R^T(\pi) \leq \tilde{R}^T(\pi) + \Delta \left(\sqrt{2T \log \frac{1}{p}} + \sum_{t=1}^T \beta^t \right) \right], \end{aligned}$$

1027 where the second inequality used (5). □

1028 **Relationship Between $\tilde{\mathcal{R}}$ and \mathcal{R}_Q**

1029 The next proposition establishes the important relationship between the regret cumulated by $\tilde{\mathcal{R}}$ and \mathcal{R}_Q .
 1030 Unlike Proposition 1, where a regret degradation was suffered for playing a strategy different than the
 1031 recommended one, in this case the regret degradation comes from the fact that the gradient that is observed
 1032 by \mathcal{R}_Q is different from that observed by $\tilde{\mathcal{R}}$. However, as we will show the degradation is kept under
 1033 control by the fact that $\tilde{\ell}^t$ is an unbiased estimator of ℓ^t by hypothesis.

1034 **Proposition 2.** *At all T , for all $p \in (0, 1)$ and $\pi \in \Pi$, with probability at least $1 - p$,*

$$\tilde{R}^T(\pi) \leq R_Q^T(\pi) + \frac{\Delta}{\beta^T} \nu \sqrt{2T \log \frac{1}{p}},$$

1035 *where*

$$\nu := \sqrt{\frac{1}{T} \sum_{t=1}^T (\rho^t)^2} = \sqrt{\frac{1}{T} \sum_{t=1}^T \max_{z \in Z} \left\{ \prod_{(j,a) \rightsquigarrow z} \left(\frac{\sum_{a' \in A_j} h^t(j, a')}{h^t(j, a)} \right)^2 \right\}}.$$

1036 *Proof.* Introduce the discrete-time stochastic process

$$d^t := (\ell^t - \tilde{\ell}^t)^\top (\pi - \mathbf{y}^t), \quad t \in \{1, 2, \dots\}.$$

1037 From the online learning model hypotheses, we have that ℓ^t and $\tilde{\ell}^t$ are conditionally independent from \mathbf{y}^t ,
 1038 and since $\mathbb{E}_t[\tilde{\ell}^t] = \ell^t$ by construction, d^t is a martingale difference sequence. At each t , the conditional
 1039 range of d^t is upper bounded by

$$|d^t| \leq |(\ell^t)^\top (\mathbf{y}^t - \pi)| + |(\tilde{\ell}^t)^\top (\mathbf{y}^t - \pi)| \leq \Delta + \frac{\Delta}{\beta^t} \rho^t \leq \frac{2\Delta}{\beta^t} \rho^t,$$

1040 where the last inequality follows since $\beta^t \leq 1$ and $\rho^t \geq 1$. Hence, the sum of quadratic ranges for the
 1041 martingale difference sequence is

$$r := \sqrt{\frac{1}{T} \sum_{t=1}^T \left(\frac{2\Delta}{\beta^t} \rho^t \right)^2} \leq \sqrt{\frac{1}{T} \sum_{t=1}^T \left(\frac{2\Delta}{\beta^T} \rho^t \right)^2} = \frac{2\Delta}{\beta^T} \nu,$$

1042 where we used the fact that the β^t are (weakly) decreasing. Using the Azuma-Hoeffding inequality, we
 1043 obtain

$$\begin{aligned} 1 - p &\leq \mathbb{P} \left[\sum_{t=1}^T d^t \leq r \sqrt{\frac{T}{2} \log \frac{1}{p}} \right] \\ &= \mathbb{P} \left[\sum_{t=1}^T d^t \leq \frac{\Delta}{\beta^T} \nu \sqrt{2T \log \frac{1}{p}} \right]. \end{aligned} \quad (6)$$

1044 Finally,

$$\begin{aligned} \sum_{t=1}^T d^t &= \sum_{t=1}^T (\ell^t - \tilde{\ell}^t)^\top (\mathbf{y}^t - \pi) \\ &= \sum_{t=1}^T (\ell^t)^\top (\mathbf{y}^t - \pi) - \sum_{t=1}^T (\tilde{\ell}^t)^\top (\mathbf{y}^t - \pi) \\ &= \tilde{R}^T(\pi) - R_Q^T(\pi), \end{aligned} \quad (7)$$

1045 and substituting (7) into (6) we have

$$\mathbb{P} \left[\tilde{R}^T(\pi) - R_Q^T(\pi) \leq \frac{\Delta}{\beta^T} \nu \sqrt{2T \log \frac{1}{p}} \right] \geq 1 - p.$$

1046 Rearranging the terms inside of the square brackets yields the statement. \square

1047 **Regret Analysis for the Overall Algorithm**

1048 Combining Proposition 2 and Proposition 1 using the union bound lemma, we obtain the following (note
1049 that the fractions $1/p$ inside of the logarithms have become $2/p$ as a consequence of the union bound
1050 lemma).

1051 **Corollary 1.** *At all T , for all $p \in (0, 1)$ and $\pi \in \Pi$, with probability at least $1 - p$,*

$$R^T(\pi) \leq R_Q^T(\pi) + \frac{\Delta}{\beta^T}(1 + \nu)\sqrt{2T \log \frac{2}{p}} + \Delta \sum_{t=1}^T \beta^t.$$

1052 where ν is as in Proposition 2.

1053 *Proof.* From Proposition 2 and Proposition 1, respectively, we have that

$$\begin{aligned} \frac{p}{2} &\geq \mathbb{P} \left[R^T(\pi) \geq \tilde{R}^T(\pi) + \Delta \left(\sqrt{2T \log \frac{2}{p}} + \sum_{t=1}^T \beta^t \right) \right] \\ \frac{p}{2} &\geq \mathbb{P} \left[\tilde{R}^T(\pi) \geq R_Q^T(\pi) + \frac{\Delta}{\beta^T} \nu \sqrt{2T \log \frac{2}{p}} \right]. \end{aligned}$$

1054 Summing the two inequalities and using the union bound, we obtain

$$\begin{aligned} p &\geq \mathbb{P} \left[R^T(\pi) \geq \tilde{R}^T(\pi) + \Delta \left(\sqrt{2T \log \frac{2}{p}} + \sum_{t=1}^T \beta^t \right) \right] + \mathbb{P} \left[\tilde{R}^T(\pi) \leq R_Q^T(\pi) + \frac{\Delta}{\beta^T} \nu \sqrt{2T \log \frac{2}{p}} \right] \\ &\geq \mathbb{P} \left[\left(R^T(\pi) \geq \tilde{R}^T(\pi) + \Delta \left(\sqrt{2T \log \frac{2}{p}} + \sum_{t=1}^T \beta^t \right) \right) \vee \left(\tilde{R}^T(\pi) \leq R_Q^T(\pi) + \frac{\Delta}{\beta^T} \nu \sqrt{2T \log \frac{2}{p}} \right) \right] \\ &\geq \mathbb{P} \left[R^T(\pi) + \tilde{R}^T(\pi) \geq \left(\tilde{R}^T(\pi) + \Delta \left(\sqrt{2T \log \frac{2}{p}} + \sum_{t=1}^T \beta^t \right) \right) + \left(R_Q^T(\pi) + \frac{\Delta}{\beta^T} \nu \sqrt{2T \log \frac{2}{p}} \right) \right] \\ &= \mathbb{P} \left[R^T(\pi) \geq R_Q^T(\pi) + \frac{\Delta}{\beta^T} (\beta^T + \nu) \sqrt{2T \log \frac{2}{p}} + \Delta \sum_{t=1}^T \beta^t \right] \\ &\geq \mathbb{P} \left[R^T(\pi) \geq R_Q^T(\pi) + \frac{\Delta}{\beta^T} (1 + \nu) \sqrt{2T \log \frac{2}{p}} + \Delta \sum_{t=1}^T \beta^t \right]. \end{aligned}$$

1055 Taking complements yields the statement. \square

1056 **Theorem 2.** *Let the local full-feedback regret minimizers \mathcal{R}_j guarantee $O(T^{1/2})$ regret in the worst case⁷,
1057 and let $p \in (0, 1)$. Furthermore, let the exploration probabilities be $\beta^t := \min\{1, k \cdot t^{-1/4}\}$ for all t , where
1058 $k > 0$ is an arbitrary constant. Then, there exists a (decision-problem-dependent) constant c independent
1059 of p and T such that for all $T \geq 1$*

$$\mathbb{P} \left[\max_{\hat{\pi} \in \Pi} R^T(\hat{\pi}) \leq c \cdot T^{3/4} \Delta \sqrt{\log \frac{2}{p}} \right] \geq 1 - p.$$

1060 When h^t is an exact measure of the number of terminal states, c is polynomial in $|\Sigma|$. Otherwise, it is linear
1061 in the constant ν defined in (2), which depends on the specific exploration functions used.

1062 *Proof.* Fix $T \in \{1, 2, \dots\}$. At all times $t \in \{1, \dots, T\}$, the norm that enter regret minimizer \mathcal{R}_j has norm
1063 upper bounded by

$$\|\ell_j^t\| \leq \frac{\Delta}{\beta^t} \rho^t \leq \Delta \max_{t=1}^T \{\rho^t\} \max \left\{ \frac{1}{k}, 1 \right\} T^{1/4},$$

1064 where we used the fact that $\beta^T \geq \min\{k, 1\} T^{-1/4}$ for all T . Letting $\tilde{M} := \max_{t=1}^T \rho^t$ and $\tilde{k} :=$
1065 $\max\{\frac{1}{k}, 1\}$, we obtain $\|\ell_j^t\| \leq \Delta \tilde{M} \tilde{k} T^{1/4}$.

⁷Valid choices include the following algorithms: regret matching (Hart and Mas-Colell 2000), regret matching⁺ (Tammelin et al. 2015), follow-the-regularized-leader, online mirror descent, exponential weights, hedge, and others.

1066 Let c_j be the constant in the regret guarantee for the $O(T^{1/2})$ regret of \mathcal{R}_j . Since regret guarantees are
 1067 always with respect to some ball of gradients of bounded norm (say, norm 1) the regret cumulated by \mathcal{R}_j is
 1068 bounded as $\max_{\hat{\pi}_j} R_j^T(\hat{\pi}_j) \leq c_j \|\ell_j^t\| T^{1/2}$. In other words, we need to keep into account the degradation
 1069 factor due to the fact that the norms of the gradients constructed through Lemma 1 might exceed the bound
 1070 for which the regret guarantee for \mathcal{R}_j was given. So, in particular,

$$\max_{\hat{\pi}_j} R_j^T(\hat{\pi}_j) \leq c_j \left(\Delta \tilde{M} \tilde{k} T^{1/4} \right) T^{1/2} = \Delta \tilde{M} \tilde{k} c_j T^{3/4},$$

1071 and so

$$\max_{q \in \mathcal{Q}} \left\{ \sum_{j \in \mathcal{J}} q[p_j] \cdot \max_{\hat{\pi}_j} R_j^T(\hat{\pi}_j) \right\} \leq \sum_{j \in \mathcal{J}} \max \left\{ 0, \max_{\hat{\pi}_j} R_j^T(\hat{\pi}_j) \right\} \leq \Delta \tilde{M} \tilde{k} \left(\sum_{j \in \mathcal{J}} c_j \right) T^{3/4}. \quad (8)$$

1072 Furthermore,

$$\frac{\Delta}{\beta^T} (1 + \nu) \sqrt{2T \log \frac{2}{p}} \leq \Delta \tilde{k} T^{1/4} (1 + \nu) \sqrt{2T \log \frac{2}{p}} = \Delta (1 + \nu) \tilde{k} T^{3/4} \sqrt{\log \frac{2}{p}}. \quad (9)$$

1073 Finally,

$$\Delta \sum_{t=1}^T \beta^t = \Delta \sum_{t=1}^T \min\{1, k \cdot t^{-1/4}\} \leq \Delta k \sum_{t=1}^T t^{-1/4} \leq \Delta k \int_0^T t^{-1/4} dt = \frac{4\Delta k}{3} T^{3/4}. \quad (10)$$

1074 Substituting the bounds (8), (9), and (10) into the general result of Theorem 1 (specifically, (??)), we obtain

$$\begin{aligned} R^T(\pi) &\leq \Delta \tilde{M} \tilde{k} \left(\sum_{j \in \mathcal{J}} c_j \right) T^{3/4} + \Delta (1 + \nu) \tilde{k} T^{3/4} \sqrt{\log \frac{2}{p}} + \frac{3\Delta k}{4} T^{3/4} \\ &= \Delta \left(\frac{4k}{3} + (1 + \nu) \tilde{k} \sqrt{\log \frac{2}{p}} + \tilde{k} \tilde{M} \sum_{j \in \mathcal{J}} c_j \right) T^{3/4} \\ &\leq \Delta \left(\frac{4k}{3 \log 2} + (1 + \nu) \tilde{k} + \frac{\tilde{k}}{\log 2} \tilde{M} \sum_{j \in \mathcal{J}} c_j \right) T^{3/4} \sqrt{\log \frac{2}{p}}, \end{aligned}$$

1075 where we used the fact that $\log(2/p) \geq \log 2$ for all $p \in (0, 1)$. Setting c to be the quantity in parentheses
 1076 and noting that when h^t measures the number of terminal states in each subtree $\nu \leq |\Sigma| - 1$ (Lemma 3),
 1077 we obtain the statement. \square

1078 **Remark 1.** *Theorem 2 shows that our algorithm achieves $O(T^{3/4} \sqrt{\log(1/p)})$ regret with high proba-*
 1079 *bility (specifically, probability at least $1 - p$). Indeed, using the properties of logarithms, $\sqrt{\log(2/p)} =$*
 1080 *$\sqrt{\log(1/p) + \log 2} \leq 2\sqrt{\log(1/p)} = O(\sqrt{\log(1/p)})$ for all $p \leq 1/2$.*

1081 Handling the Unknown Structure

1082 As already noted in Section B, only the local full-feedback regret minimizers \mathcal{R}_j on the path from the
 1083 root to z^t observe a nonzero counterfactual gradient, which is computed as in the *regret update phase*
 1084 described in Section 4. All other local regret minimizers see a zero gradient, and therefore their regret does
 1085 not increase. So, we can safely avoid updating those regret minimizers that are not on the path from the
 1086 root to the most recent terminal leaf.

1087 The upfront-flipping sampling scheme and the on-path-flipping sampling scheme can both be imple-
 1088 mented so that the actions are sampled incrementally as the interaction with the environment progresses. In
 1089 other words, there is no need to sample actions at all decision nodes upfront in order to interact with the
 1090 environment. In the body of this paper, we only described how to do so for the upfront-flipping sampling
 1091 scheme, which is conceptually the most important, as it showcases well the difference with, for example,
 1092 epsilon-greedy exploration used by MCCFR.

1093 C Additional Details about Experiments

1094 Description of Games

1095 Here, we review the two standard benchmark games that we use in the experiments.

1096 In Kuhn poker (Kuhn 1950), two players put an ante worth 1 into the pot at the beginning of the game.
 1097 Then, each player is privately dealt one card from a deck that contains only three cards—specifically, jack,

1098 queen, and king. Then Player 1 decides to either check or bet 1. If Player 1 checks, Player 2 can decide to
 1099 either check or raise 1. If Player 2 checks, a showdown occurs. If Player 2 raises, Player 1 can fold—and
 1100 the game ends with Player 2 taking the pot—or call, at which point a showdown occurs. Otherwise, if the
 1101 first action of Player 1 is to raise, then Player 2 may fold (the game ends with Player 1 taking the pot) or
 1102 call, at which point a showdown occurs. In a showdown, the player with the higher card wins the pot and
 1103 the game ends.

1104 *Leduc poker* (Southey et al. 2005) is played with a deck of 3 unique ranks, each appearing twice in the
 1105 deck. There are two rounds in the game. In the first round, all players put an ante of 1 in the pot and are
 1106 privately dealt a single card. A round of betting then starts. Player 1 acts first, and at most two bets are
 1107 allowed per player. Then, a card is publicly revealed, and another round of betting takes place, with the
 1108 same dynamics described above. After the two betting round, if one of the players has a pair with the public
 1109 card, that player wins the pot. Otherwise, the player with the higher card wins the pot. All bets in the first
 1110 round are worth 2, while all bets in the second round are 4.

1111 **Hyperparameter Tuning for the Policy Gradient Approach by Srinivasan et al. (2018)**

1112 We acknowledge the help of some of the original authors (Srinivasan et al. 2018) in tuning hyperparameters
 1113 for their algorithm on the Leduc poker benchmark.

1114 For both the RPG and the QPG formulation, the following combinations of hyperparameters was tested:

- 1115 • Critic learning rate: $\{0.001, 0.05, 0.01\}$;
- 1116 • Pi learning rate: $\{0.001, 0.01, 0.05\}$;
- 1117 • Entropy cost (used to multiply the entropy loss): $\{0.01, 0.1\}$;
- 1118 • Batch size (for Q and Pi learning): $\{16, 64, 128\}$;
- 1119 • Num critics (before each Pi update): $\{16, 64, 128\}$.

1120 In Tables 2 and 3 we report the average exploitability (computed across 5 independent runs for each
 1121 choice of hyperparameters) for the top 5 choices of hyperparameters for the RPG and QPG formulation,
 1122 respectively. The experimental results in the body show data for the best combination of eitherparameters
 1123 (first row of each table).

Critic learning rate	Pi learning rate	Entropy cost	Batch size	Num Critics	Avg. exploitability
0.01	0.05	0.01	128	64	0.448680
0.01	0.05	0.01	64	64	0.451409
0.001	0.01	0.01	16	128	0.472729
0.01	0.01	0.01	128	16	0.493694
0.05	0.01	0.01	128	64	0.494840

Table 2: Performance of PG with RPG gradient formulation using the top 5 combinations of hyperparameters.

Critic learning rate	Pi learning rate	Entropy cost	Batch size	Num Critics	Avg. exploitability
0.01	0.05	0.01	64	128	0.421745
0.01	0.01	0.01	64	16	0.451535
0.001	0.01	0.01	16	128	0.457430
0.01	0.01	0.01	64	128	0.471100
0.05	0.01	0.01	128	64	0.524005

Table 3: Performance of PG with QPG gradient formulation using the top 5 combinations of hyperparameters.

1124 **Additional Experiments**

1125 In Figures 3 to 6 we study how the uniform and balanced exploration functions compare in Kuhn and
 1126 Leduc poker, for both the on-path-flipping and the upfront-flipping sampling scheme. We test different
 1127 values of the exploration multiplier k (Theorem 2), specifically $k \in \{0.5, 1, 10, 20\}$.

1128 The experiments show that the on-path-flipping sampling schemes leads to significantly less variance
 1129 than the upfront-flipping sampling scheme. Sometimes the latter leads to better average convergence very
 1130 early on in the learning, but overall has worse convergence in practice.

1131 The experiments show no meaningful different between the uniform and balanced exploration strategies.

Figure 7 shows how MCCFR performed with different exploration parameters.

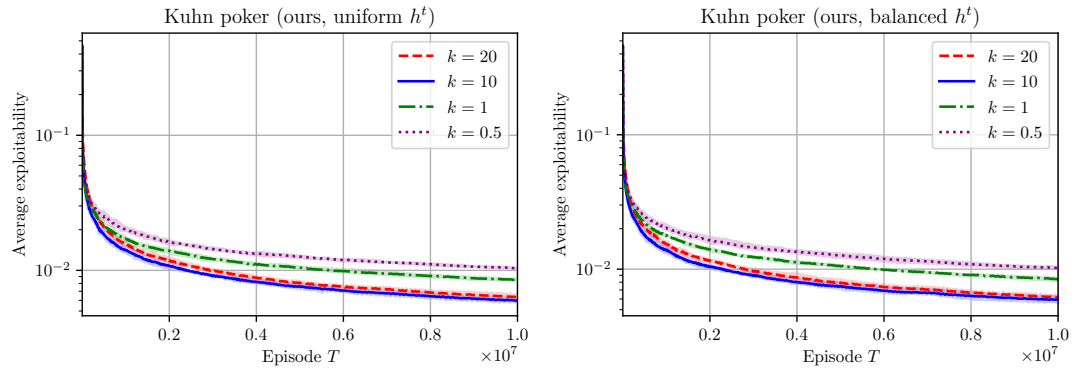


Figure 3: Comparison between uniform and balanced exploration strategies in Kuhn poker, for different values of the exploration multiplier k , when using the on-path-flipping sampling scheme.

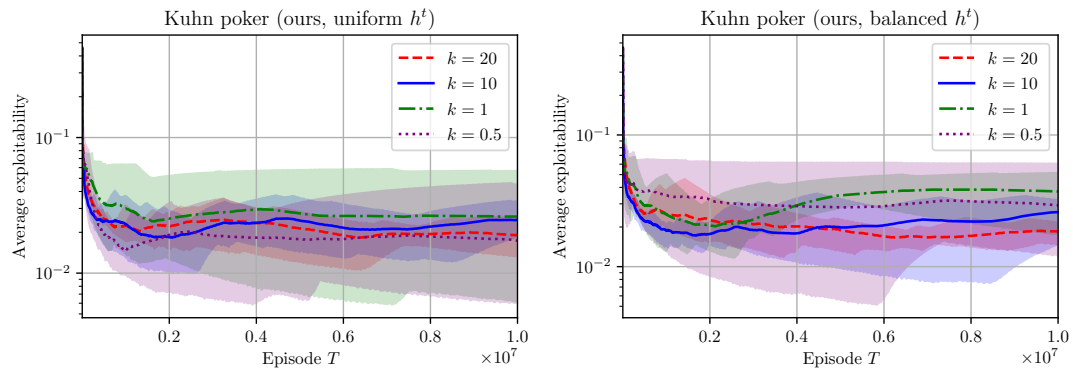


Figure 4: Comparison between uniform and balanced exploration strategies in Kuhn poker, for different values of the exploration multiplier k , when using the upfront-flipping sampling scheme.

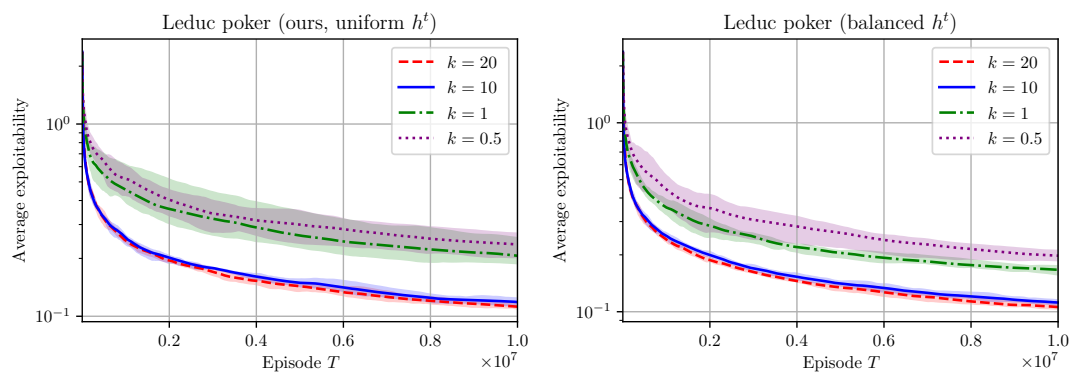


Figure 5: Comparison between uniform and balanced exploration strategies in Leduc poker, for different values of the exploration multiplier k , when using the on-path-flipping sampling scheme.

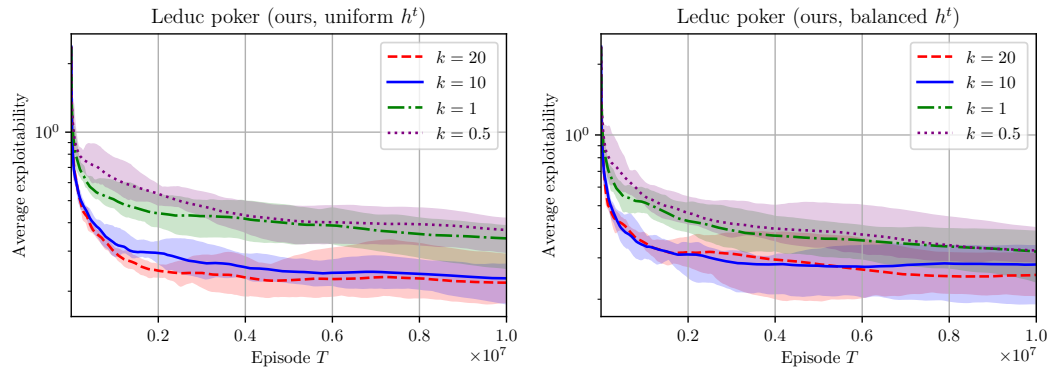


Figure 6: Comparison between uniform and balanced exploration strategies in Leduc poker, for different values of the exploration multiplier k , when using the upfront-flipping sampling scheme.

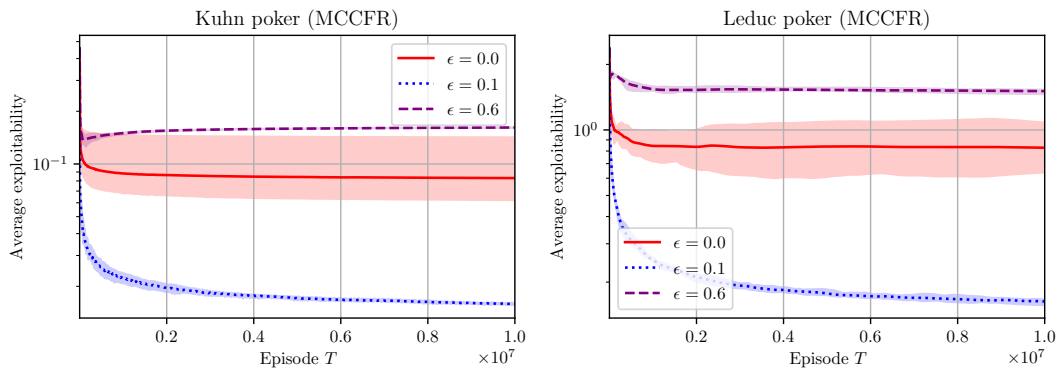


Figure 7: Performance of MCCFR for different amounts of ϵ -greedy exploration.