

CONVERGENCE OF REGRET MATCHING IN POTENTIAL GAMES AND CONSTRAINED OPTIMIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Regret matching (*RM*)—and its modern variants—is a foundational online algorithm that has been at the heart of many AI breakthrough results in solving benchmark zero-sum games, such as poker. Yet, surprisingly little is known so far in theory about its convergence beyond two-player zero-sum games. For example, whether regret matching converges to Nash equilibria in *potential games* has been an open problem for two decades. Even beyond games, one could try to use RM variants for general constrained optimization problems. Recent empirical evidence suggests that they—particularly *regret matching*⁺ (*RM*⁺)—attain strong performance on benchmark constrained optimization problems, outperforming traditional gradient descent-type algorithms.

We show that *RM*⁺ converges to an ϵ -KKT point after $O_\epsilon(1/\epsilon^4)$ iterations, establishing for the first time that it is a sound and fast first-order optimizer. Our argument relates the KKT gap to the accumulated *regret*, two quantities that are entirely disparate in general but interact in an intriguing way in our setting, so much so that when regrets are bounded, our complexity bound improves all the way to $O_\epsilon(1/\epsilon^2)$. From a technical standpoint, while *RM*⁺ does *not* have the usual one-step improvement property in general, we show that it does in a certain region that the algorithm will quickly reach and remain in thereafter. In sharp contrast, our second main result establishes a lower bound: *RM*, with or without alternation, can take an exponential number of iterations to reach a crude approximate solution even in two-player potential games. This represents the first worst-case separation between *RM* and *RM*⁺. Our lower bound shows that convergence to coarse correlated equilibria in potential games is exponentially faster than convergence to Nash equilibria.

1 INTRODUCTION

Regret matching is a foundational online algorithm for minimizing *regret*. It was famously introduced by Hart & Mas-Colell (2000), although its conception can be traced much further back to the seminal *approachability* framework of Blackwell (1956), which lay the groundwork for online learning and regret minimization. As the name suggests, regret matching prescribes playing each action with probability proportional to the (nonnegative) regret accumulated by that action. Its appeal lies in its simplicity and scalability, being both *parameter free* and *scale invariant*.

Regret matching—and modern versions thereof—has been at the forefront of equilibrium computation in massive two-player zero-sum games. A notable variant with strong empirical performance is *regret matching*⁺, introduced by Tammelin (2014); the only difference is that it truncates all negative coordinates of the regret vector to zero in every iteration. Even so, this variant is typically far superior than its predecessor, and was a central component in AI poker breakthroughs (Bowling et al., 2015; Brown & Sandholm, 2017; 2019b; Moravčík et al., 2017) and a more recent superhuman agent for dark chess (Zhang & Sandholm, 2025).

As such, the regret matching family of algorithms has rightfully been the subject of intense study in contemporary research. Much of this focus has been confined to two-player zero-sum games, where minimizing regret translates to convergence—of the *average* strategies—to minimax (equivalently, Nash) equilibria (Freund & Schapire, 1999). More broadly, in general-sum games, no-regret algo-

054 rithms guarantee convergence to the set of *coarse correlated equilibria* (Moulin & Vial, 1978)—a
 055 more permissive concept than Nash equilibria.

056
 057 In this paper, we examine the convergence of regret matching and its variants in the seminal class
 058 of *potential games*, and, more broadly, nonconvex optimization constrained over a product of sim-
 059 plices. Surprisingly little is known about this question even though it was identified early on as an
 060 important open question in this space (Kleinberg et al., 2009; Marden et al., 2007). Recent empir-
 061 ical evidence brings this question to the fore again: Tewolde et al. (2025) showed that the regret
 062 matching family—and especially regret matching⁺—attains strong performance on a benchmark
 063 suite of constrained optimization problems, significantly outperforming gradient descent-type algo-
 064 rithms. Yet, there is no theory to suggest that regret matching will even asymptotically converge to
 065 approximate KKT points in constrained optimization, which are tantamount to Nash equilibria when
 066 dealing specifically with potential games. We fill this gap in this paper.

067 1.1 OUR RESULTS

068
 069 We analyze the convergence of regret matching (RM) and regret matching (RM⁺) in the general class
 070 of (nonconvex) optimization problems constrained over a product of probability simplices. This
 071 encompasses as a special case Nash equilibria in potential games when the objective is multilinear.
 072 More broadly, to have a unifying treatment of both settings, we think of each probability simplex as
 073 being controlled by a single player who is observing the corresponding part of the gradient.

074 We cover both the simultaneous and the alternating version of RM⁺—whereby players update their
 075 strategies one after the other, akin to coordinate descent. Our result for RM⁺ is summarized below.

076 **Theorem 1.1.** *RM⁺ converges to an ϵ -KKT point of any optimization problem over a product of*
 077 *simplices after $O_\epsilon(1/\epsilon^4)$ iterations.*

078
 079 This theorem confirms that RM⁺ is a sound and efficient first-order optimizer, lending further cre-
 080 dence to the empirical results of Tewolde et al. (2025). We hope that **Theorem 1.1** will help cement
 081 RM⁺ in the optimization arsenal going forward.

082 We remark that for potential games and constrained optimization over a single simplex, the $O_\epsilon(1/\epsilon^4)$
 083 bound holds no matter how the regrets in (alternating) RM⁺ are initialized. For constrained opti-
 084 mization over multiple simplices—with or without alternation—we obtain the same rate by suitably
 085 initializing the regret vectors (**Corollaries 3.11** and **C.11**); for the usual parameter-free version of
 086 RM⁺, in which the regret vectors are initialized at zero, we can only guarantee an inferior bound
 087 growing as $O_\epsilon(1/\epsilon^8)$ (**Theorem 3.12**).

088 Our argument proceeds by parameterizing the rate of convergence of RM⁺ as a function of the
 089 accumulated regret, so much so that if the regret with respect to each individual simplex remains
 090 bounded, the rate is improved all the way to $T^{-1/2}$.

091 **Theorem 1.2.** *Suppose that the regret of RM⁺ on each individual simplex grows as at most T^α for*
 092 *some $\alpha \in [0, 1/2]$. Then RM⁺ converges to an ϵ -KKT point after $O_\epsilon(1/\epsilon^{2/(1-\alpha)})$ iterations.*

093
 094 RM⁺ always guarantees regret growing as \sqrt{T} , so **Theorem 1.1** is implied by **Theorem 1.2**. What
 095 makes the latter theorem surprising is that, in general, regret is a fundamentally disparate property
 096 compared to KKT gap: as we point out in **Proposition 3.2**, a sequence can incur zero regret while
 097 having an $\Omega(1)$ KKT gap in each iteration. Even so, **Theorem 1.2** directly relates the KKT gap in
 098 terms of the regret. In particular, the non-asymptotic rate of **Theorem 1.1** is a consequence of the
 099 fact that RM⁺ has the no-regret property! In the special case of potential games, regret is known
 100 to drive the rate of convergence to *coarse correlated equilibria* (CCE); **Theorem 1.2** shows, for the
 101 first time, that regret can also govern the rate of convergence to Nash equilibria. In particular, if
 102 convergence to CCE happens at a rate of $T^{-(1-\alpha)}$, for some $\alpha \in [0, 1/2]$, the rate of convergence to
 103 Nash equilibria is no slower than $T^{-\frac{1-\alpha}{2}}$.

104 From a technical standpoint, the key challenge is that RM⁺ does *not* have a one-step improvement
 105 property: even if one initializes RM⁺ close to a KKT point, RM⁺ can still grossly overshoot. And, of
 106 course, it is a parameter-free algorithm, so the usual treatment of gradient descent-type algorithms
 107 that relies on appropriately tuning the learning rate falls short. In this context, our starting observa-
 tion is that, at least when the utility function is linear, RM⁺ is bound to improve the utility, although

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

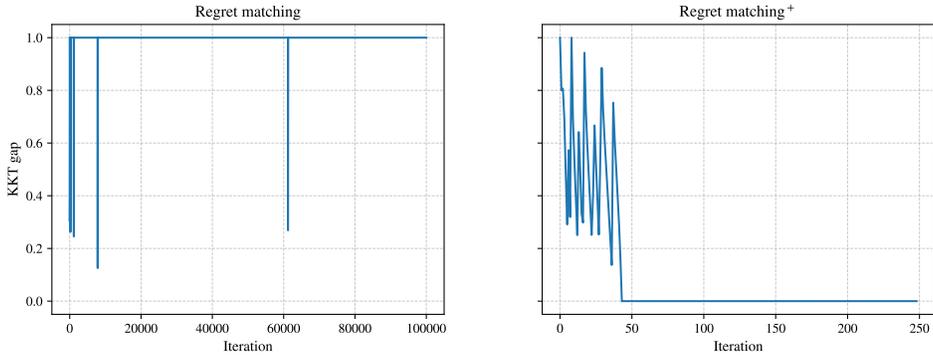


Figure 1: Illustration of our main results: RM^+ always converges fast to a KKT point while RM can take exponential time even in two-player identical-interest games, constructed in Section 4.

the improvement is inversely proportional to the norm of the regret vector (Lemma 3.3). This key property already suffices to show that alternating RM^+ will converge to Nash equilibria in potential games. For the more challenging setting where the updates are simultaneous or the objective is not multilinear, we first show that one-step improvement holds conditional on the norm of the regret vector being sufficiently large (Lemmas 3.7 and C.10). To conclude the argument, we combine this property with the crucial insight that the ℓ_2 norm of the regret vector is monotonically increasing proportionally to the KKT gap (Lemma 3.8). This means that RM^+ will never get stuck in a cycle: the regret vector would quickly grow in norm, at which point the one-step improvement promised by Lemma 3.7 kicks in.

Does RM share the same convergence properties as RM^+ ? As a reminder, the only difference is that RM refrains from truncating negative regrets to zero. Even so, we find that this seemingly innocuous difference gives rise to an exponential gap in the performance of RM vis-à-vis RM^+ , manifested even in two-player identical-interest games—a special case of potential games (Figure 1).

Theorem 1.3. *There is a two-player $m \times m$ identical-interest game where RM , with or without alternation, requires $m^{\Omega(m)}$ iterations to converge to an $m^{-\Theta(1)}$ -approximate Nash equilibrium.*

This lower bound holds not just under an adversarial initialization, but also when players initially mix uniformly at random, which is the most common initialization in practice.

Theorem 1.3 constitutes the first worst-case separation—let alone an exponential one—between RM and RM^+ . Indeed, in zero-sum games, it is known that RM and RM^+ both attain a rate no faster than $T^{-1/2}$ (Farina et al., 2023), even though RM^+ typically performs much better in practice. Theorem 1.3 provides further justification for opting for RM^+ instead of RM , albeit in a fundamentally different setting.

The basic flaw of RM that underpins Theorem 1.3 is that, even with a linear utility, it is not guaranteed to improve the utility even when it has a large best-response gap; specifically, as we show in Lemma 3.6, the improvement is conditional on a good-enough action having nonnegative regret. But herein lies the problem: it could take many iterations before the regret resurfaces to being positive. What happens in the construction behind Theorem 1.3 is that it takes longer and longer—exponentially so—for the regret of the unique good-enough action to be positive; before then, RM is entirely stalled without making any progress. At the same time, RM is guaranteed to converge to the set of coarse correlated equilibria (CCE) at a rate of $T^{-1/2}$, simply because it always has the no-regret property. This leads to the following interesting consequence.

Corollary 1.4. *There is a class of two-player potential games in which RM converges to an ϵ -CCE in $O_\epsilon(1/\epsilon^2)$ rounds but it takes $\exp(\Omega(1/\epsilon))$ rounds to converge to an ϵ -Nash equilibrium.*

To be clear, convergence to a CCE is meant in terms of the average correlated distribution of play, whereas convergence to Nash equilibria is in terms of the individual iterates produced by RM .

We defer further discussion on related work to Section A.

2 PRELIMINARIES

We begin by introducing potential games and the more general problem of constrained optimization over a product of simplices (Section 2.1), and then recall regret matching⁽⁺⁾ in Section 2.2.

2.1 POTENTIAL GAMES AND CONSTRAINED OPTIMIZATION

Normal-form games Our first key focus in this paper is on *potential games*, which we represent in the usual normal form. Here, we have n players, each of whom is to select an action a_i from a finite set \mathcal{A}_i , with $m_i := |\mathcal{A}_i|$ and $m = \max_{1 \leq i \leq n} m_i$. Under a joint action profile $\mathbf{a} = (a_1, \dots, a_n) \in \mathcal{A}_1 \times \dots \times \mathcal{A}_n$, each player $i \in [n]$ receives a payoff given by a *utility function* $u_i : (a_1, \dots, a_n) \mapsto u_i(a_1, \dots, a_n) \in \mathbb{R}$ with range bounded by 1. A player $i \in [n]$ can randomize by specifying a *mixed strategy* $\mathbf{x}_i \in \Delta(\mathcal{A}_i) := \{\mathbf{x}_i \in \mathbb{R}_{\geq 0}^{\mathcal{A}_i} : \sum_{a_i \in \mathcal{A}_i} \mathbf{x}_i[a_i] = 1\}$. Player i strives to maximize its *expected* utility, given by $u_i(\mathbf{x}_1, \dots, \mathbf{x}_n) := \sum_{(a_1, \dots, a_n) \in \mathcal{A}_1 \times \dots \times \mathcal{A}_n} u_i(a_1, \dots, a_n) \prod_{i'=1}^n \mathbf{x}_{i'}[a_{i'}]$. A key fact is that the expected utility is *multi-linear*, in that $u_i(\mathbf{x}_1, \dots, \mathbf{x}_n) = \langle \mathbf{x}_i, \mathbf{u}_i(\mathbf{x}_{-i}) \rangle$ for some utility vector $\mathbf{u}_i(\mathbf{x}_{-i}) \in \mathbb{R}^{\mathcal{A}_i}$ that does not depend on \mathbf{x}_i ; namely, $\mathbf{u}_i(\mathbf{x}_{-i}) = (\sum_{(a_1, \dots, a_{i-1}, a_{i+1}, \dots, a_n)} u_i(a_1, \dots, a_n) \prod_{i' \neq i} \mathbf{x}_{i'}[a_{i'}])_{a_i \in \mathcal{A}_i}$. Here and throughout, we use the shorthand notation $\mathbf{x}_{-i} = (\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n)$, while we recall that $\langle \cdot, \cdot \rangle$ denotes the inner product. Further, we use the shorthand notation $\text{BRGap}_i(\mathbf{x}_i, \mathbf{u}_i) := \max_{\mathbf{x}'_i \in \Delta(\mathcal{A}_i)} \langle \mathbf{x}'_i - \mathbf{x}_i, \mathbf{u}_i \rangle$ for the best-response gap.

The predominant solution concept in game theory is the *Nash equilibrium* (Nash, 1950).

Definition 2.1. A strategy profile $(\mathbf{x}_1, \dots, \mathbf{x}_n) \in \Delta(\mathcal{A}_1) \times \dots \times \Delta(\mathcal{A}_n)$ is an ϵ -Nash equilibrium if for any player $i \in [n]$ and unilateral deviation $\mathbf{x}'_i \in \Delta(\mathcal{A}_i)$,

$$u_i(\mathbf{x}'_i, \mathbf{x}_{-i}) \leq u_i(\mathbf{x}_i, \mathbf{x}_{-i}) + \epsilon.$$

A standard relaxation of the Nash equilibrium is the *coarse correlated equilibrium* (Definition B.1), which can be attained by no-regret algorithms (Proposition B.2). While finding a Nash equilibrium is hard even in two-player general-sum games (Daskalakis et al., 2008; Chen et al., 2009), our focus is on *potential games*—equivalently, *congestion games* (Monderer & Shapley, 1996).

Potential games This is a seminal class that goes back to the work of Rosenthal (1973). The defining property is the admission of a global, player-independent function—the *potential*—whose difference reflects the benefit of any unilateral deviation.

Definition 2.2 (Potential game). An n -player game is a *potential game* if there exists a function $\Phi : \mathcal{A}_1 \times \dots \times \mathcal{A}_n \rightarrow \mathbb{R}$ such that for any player $i \in [n]$ and actions $a_i, a'_i \in \mathcal{A}_i$, $\mathbf{a}_{-i} \in \times_{i' \neq i} \mathcal{A}_{i'}$,

$$\Phi(a'_i, \mathbf{a}_{-i}) - \Phi(a_i, \mathbf{a}_{-i}) = u_i(a'_i, \mathbf{a}_{-i}) - u_i(a_i, \mathbf{a}_{-i}). \quad (1)$$

A special case of a potential game worth noting is an *identical-interest* game, which means that $u_1(\mathbf{x}_1, \dots, \mathbf{x}_n) = \dots = u_n(\mathbf{x}_1, \dots, \mathbf{x}_n)$ for all $\mathbf{x}_1, \dots, \mathbf{x}_n$. In the presence of only two players, this simplifies to $u_1(\mathbf{x}_1, \mathbf{x}_2) = \langle \mathbf{x}_1, \mathbf{A}\mathbf{x}_2 \rangle = u_2(\mathbf{x}_1, \mathbf{x}_2)$ for a common matrix $\mathbf{A} \in \mathbb{R}^{\mathcal{A}_1 \times \mathcal{A}_2}$.

A (mixed) Nash equilibrium in potential games is amenable to (projected) gradient descent, but is likely hard to compute when the precision $\epsilon > 0$ is exponentially small (Babichenko & Rubinstein, 2021). Our focus will be on algorithms whose complexity is polynomial in $1/\epsilon$.

Constrained optimization More broadly, beyond potential games, we are interested in computing *Karush-Kuhn-Tucker (KKT) points* of a function $u : \mathcal{X} \rightarrow \mathbb{R}$, where $\mathcal{X} := \Delta(\mathcal{A}_1) \times \dots \times \Delta(\mathcal{A}_n)$. We assume that u , which is to be maximized, is differentiable over an open set $\hat{\mathcal{X}} \supset \mathcal{X}$ and L -smooth, meaning that $\|\nabla u(\mathbf{x}) - \nabla u(\mathbf{x}')\|_2 \leq L\|\mathbf{x} - \mathbf{x}'\|_2$ for all $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$; we recall that $\|\mathbf{x}\|_2 := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$ denotes the (Euclidean) ℓ_2 norm. We make the normalization assumption $|\langle \mathbf{x}_i - \mathbf{x}'_i, \nabla_{\mathbf{x}_i} u(\mathbf{x}) \rangle| \leq 1$ for all $i \in [n]$ and $\mathbf{x}_i, \mathbf{x}'_i \in \Delta(\mathcal{A}_i)$. The goal is to minimize the *KKT gap*, which we measure by

$$\text{KKTGap} : \mathcal{X} \ni \mathbf{x} \mapsto \max_{\mathbf{x}' \in \mathcal{X}} \langle \mathbf{x}' - \mathbf{x}, \nabla u(\mathbf{x}) \rangle = \sum_{i=1}^n \text{BRGap}_i(\mathbf{x}_i, \nabla_{\mathbf{x}_i} u(\mathbf{x})). \quad (2)$$

A point with small KKT gap per (2) is also referred to as an approximate *first-order stationary point*, which is an approximate fixed point of the (constrained) gradient descent mapping $\mathbf{x} \mapsto \Pi_{\mathcal{X}}(\mathbf{x} + \eta \nabla u(\mathbf{x}))$, where $\eta \leq 1/L$ and $\Pi_{\mathcal{X}}(\cdot)$ is (Euclidean) projection mapping. A potential game can be seen as the special case in which u is multilinear.

2.2 ONLINE LEARNING AND REGRET MATCHING

Moving on, we now introduce RM and RM^+ within the framework of online learning. Here, a *learner* interacts with an *environment* over a sequence of T rounds. In each round $t \in [T]$, the learner first elects a mixed strategy $\mathbf{x} \in \Delta(\mathcal{A})$. The environment in turn specifies a linear utility function $u^{(t)} : \mathbf{x} \mapsto \langle \mathbf{x}, \mathbf{u}^{(t)} \rangle$ for some utility vector $\mathbf{u}^{(t)} \in \mathbb{R}^{\mathcal{A}}$; it is assumed that $u^{(t)}$ has a range bounded by 1. In the full-feedback setting, $\mathbf{u}^{(t)}$ is revealed to the learner at the end of the t th round. The performance of the learner in this online environment is evaluated through *regret*,

$$\text{Reg}^{(T)} := \max_{\mathbf{x}' \in \Delta(\mathcal{A})} \sum_{t=1}^T \langle \mathbf{x}' - \mathbf{x}^{(t)}, \mathbf{u}^{(t)} \rangle. \quad (3)$$

Two popular algorithms for minimizing regret on the simplex are regret matching (RM) and regret matching⁺ (RM^+), formally defined in Algorithms 1 and 2. They both prescribe playing an action with probability proportional to the nonnegative regret accumulated by that action. Their *only* difference is that RM^+ always truncates the regret to 0 (Algorithm 2); in that line, $\mathbf{1}$ denotes the all-ones vector, whose dimension is omitted as it is clear from the context, and $[\cdot]^+ := \max(\mathbf{0}, \cdot)$ is the nonnegative part.

Proposition 2.3 (Zinkevich et al., 2007; Farina et al., 2021). *For any sequence of utilities $(\mathbf{u}^{(t)})_{t=1}^T$, both RM and RM^+ guarantee that the ℓ_2 norm of $[\mathbf{r}^{(T)}]^+$ is at most \sqrt{mT} .*

In particular, for both RM and RM^+ , $\text{Reg}^{(T)} \leq \|[\mathbf{r}^{(T)}]^+\|_{\infty} \leq \|[\mathbf{r}^{(T)}]^+\|_2 \leq \sqrt{mT}$.

Algorithm 1: Regret matching (RM)

```

1 Initialize cumulative regrets  $\mathbf{r}^{(0)} \leftarrow \mathbf{0}$ ;
2 Initialize strategy  $\mathbf{x}^{(0)} \in \Delta(\mathcal{A})$ ;
3 for  $t = 1, \dots, T$  do
4   Set  $\boldsymbol{\theta}^{(t)} \leftarrow [\mathbf{r}^{(t-1)}]^+$ ;
5   if  $\boldsymbol{\theta}^{(t)} \neq \mathbf{0}$  then
6     Compute  $\mathbf{x}^{(t)} \leftarrow \boldsymbol{\theta}^{(t)} / \|\boldsymbol{\theta}^{(t)}\|_1$ ;
7   else
8      $\mathbf{x}^{(t)} \leftarrow \mathbf{x}^{(t-1)}$ ;
9   Output strategy  $\mathbf{x}^{(t)} \in \Delta(\mathcal{A})$ ;
10  Observe utility  $\mathbf{u}^{(t)} \in \mathbb{R}^{\mathcal{A}}$ ;
11   $\mathbf{r}^{(t)} \leftarrow \mathbf{r}^{(t-1)} + \mathbf{u}^{(t)} - \langle \mathbf{x}^{(t)}, \mathbf{u}^{(t)} \rangle \mathbf{1}$ ;
```

Algorithm 2: Regret matching⁺ (RM^+)

```

1 Initialize cumulative regrets  $\mathbf{r}^{(0)} := \mathbf{0}$ ;
2 Initialize strategy  $\mathbf{x}^{(1)} \in \Delta(\mathcal{A})$ ;
3 for  $t = 1, \dots, T$  do
4   Set  $\boldsymbol{\theta}^{(t)} \leftarrow \mathbf{r}^{(t-1)}$ ;
5   if  $\boldsymbol{\theta}^{(t)} \neq \mathbf{0}$  then
6     Compute  $\mathbf{x}^{(t)} \leftarrow \boldsymbol{\theta}^{(t)} / \|\boldsymbol{\theta}^{(t)}\|_1$ ;
7   else
8      $\mathbf{x}^{(t)} \leftarrow \mathbf{x}^{(t-1)}$ ;
9   Output strategy  $\mathbf{x}^{(t)} \in \Delta(\mathcal{A})$ ;
10  Observe utility  $\mathbf{u}^{(t)} \in \mathbb{R}^{\mathcal{A}}$ ;
11   $\mathbf{r}^{(t)} \leftarrow [\mathbf{r}^{(t-1)} + \mathbf{u}^{(t)} - \langle \mathbf{x}^{(t)}, \mathbf{u}^{(t)} \rangle \mathbf{1}]^+$ ;
```

Simultaneous and alternating updates We are interested in the convergence of RM and RM^+ when used by all players; in the constrained optimization setting, we think of having one player acting on each simplex, in direct correspondence with potential games. In this setting, the sequence of utilities $(\mathbf{u}_i^{(t)})_{t=1}^T$ given as input to player $i \in [n]$ is determined by the strategies of the other players. If the updates are *simultaneous*, we have $\mathbf{u}_i^{(t)} = \nabla_{\mathbf{x}_i} u(\mathbf{x}^{(t)})$ for each player $i \in [n]$. (In potential games, the potential function Φ plays the role of u , so that $\mathbf{u}_i^{(t)} = \mathbf{u}_i(\mathbf{x}_{-i}^{(t)})$.) In the alternating setting, we go through the players in a round-robin fashion $i = 1, \dots, n$. In the *lazy* version of the update, for a fixed precision $\epsilon > 0$, we first compute $\mathbf{u}_i^{(t)} = \nabla_{\mathbf{x}_i} u(\mathbf{x}_{i' < i}^{(t+1)}, \mathbf{x}_{i' \geq i}^{(t)})$. If the best-response gap of player $i \in [n]$ is already at most ϵ , we refrain from updating that player, so that $\mathbf{x}_i^{(t+1)} := \mathbf{x}_i^{(t)}$. Otherwise, the player updates its strategy to $\mathbf{x}_i^{(t+1)}$ using $\mathbf{u}_i^{(t)}$. We refer to this scheme as ϵ -*lazy alternating* updates (ϵ -lazy simultaneous updates are defined similarly); one limitation of this lazy variant is that it is not an “anytime algorithm,” in that one needs to specify

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

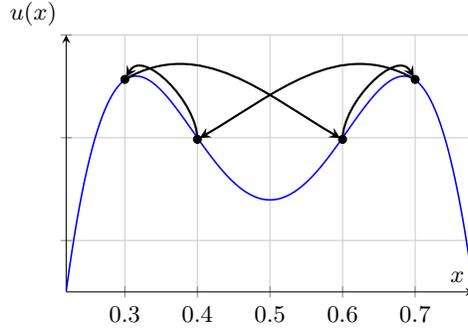


Figure 2: The example corresponding to [Proposition 3.2](#), demonstrating that having zero regret, let alone sublinear, has no implications concerning convergence in terms of KKT gap.

the precision beforehand. In the more common, non-lazy version of alternation, a player is updated regardless of the best-response gap; in the sequel, we obtain convergence results for both variants.

3 CONVERGENCE OF REGRET MATCHING⁺

In this section, we analyze the convergence of RM^+ in potential games ([Section 3.1](#)), and more broadly, constrained optimization ([Section 3.2](#)). A central theme in our analysis of RM and RM^+ is a recurring connection between regret and convergence to KKT points.

Before we proceed, it is worth highlighting that, in general, the no-regret property is fundamentally different from convergence to KKT points in *nonconvex problems*. To begin with, we point out that when the underlying function to be maximized, u , is concave, then the no-regret property does imply convergence to a global optimum, from Jensen’s inequality.

Proposition 3.1 (Under concavity, no-regret implies convergence). *Let u be a smooth concave function. If an online algorithm observes the sequence of utilities $(\nabla u(\mathbf{x}^{(t)}))_{t=1}^T$, then $\frac{1}{T} \sum_{t=1}^T u(\mathbf{x}^{(t)}) \geq \max_{\mathbf{x} \in \mathcal{X}} u(\mathbf{x}) - \frac{1}{T} \text{Reg}^{(T)}$, where $\text{Reg}^{(T)}$ is the regret of the algorithm per (3).*

Thus, if the algorithm has vanishing average regret, $u(\mathbf{x}^{(t)})$ converges to $\max_{\mathbf{x} \in \mathcal{X}} u(\mathbf{x})$ (in density). But beyond concave problems, no-regret algorithms do not necessarily guarantee convergence even to a KKT point, as we point out below.

Proposition 3.2. *For any $T \in \mathbb{N}$ with $T \equiv 0 \pmod{4}$, there exists a polynomial function u in $[0, 1]$ and a sequence of points $(x^{(t)})_{t=1}^T$ such that*

- the regret of the sequence with respect to $(\nabla u(x^{(t)}))_{t=1}^T$ is zero, while
- every point in the sequence has an $\Omega(1)$ KKT gap with respect to the function u .

This is based on the 4-cycle $0.6 \rightarrow 0.7 \rightarrow 0.4 \rightarrow 0.3 \rightarrow 0.6$. If the gradients observed at those points are $0.6 \mapsto 2, 0.7 \mapsto -1, 0.4 \mapsto -2, 0.3 \mapsto 1$, it follows that i) $\sum_{t=1}^T \nabla u(x^{(t)}) = 0$ and ii) $\sum_{t=1}^T x^{(t)} \nabla u(x^{(t)}) = 0$, which in turn implies that this sequence incurs zero regret. But, by construction, the gradients at those interior points have a large magnitude, which in turn implies that the KKT gap is large. (That the average is a local minimum is coincidental.) A polynomial consistent with the above gradients is $90x - 298.3x^2 + 416.6x^3 - 208.3x^4$, leading to [Proposition 3.2](#); we note that the above sequence of iterates is not realizable through an algorithm such as gradient descent.

3.1 POTENTIAL GAMES

We first analyze convergence in potential games. A key property, which paves the way for [Theorem 3.4](#), is that, for a fixed utility vector, RM^+ has a one-step improvement property; the lemma below takes the perspective of a single, arbitrary player in the game.

Lemma 3.3 (One-step improvement for RM^+). *For any $\mathbf{r} \in \mathbb{R}_{\geq 0}^A$ and $\mathbf{u} \in \mathbb{R}^A$, we define $\mathbf{x} := \mathbf{r} / \|\mathbf{r}\|_1$; if $\mathbf{r} = \mathbf{0}$, $\mathbf{x} \in \Delta(\mathcal{A})$ can be arbitrary. If $\mathbf{r}' := [\mathbf{r} + \mathbf{u} - \langle \mathbf{x}, \mathbf{u} \rangle \mathbf{1}]^+ \neq \mathbf{0}$ and $\mathbf{x}' := \mathbf{r}' / \|\mathbf{r}'\|_1$,*

324 then

$$325 \quad \langle \mathbf{x}' - \mathbf{x}, \mathbf{u} \rangle \geq \frac{1}{\|\mathbf{r}'\|_1} \left(\max_{a \in \mathcal{A}} \mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle \right)^2 = \frac{1}{\|\mathbf{r}'\|_1} \text{BRGap}(\mathbf{x}, \mathbf{u})^2. \quad (4)$$

327 If $\mathbf{r}' = \mathbf{0}$, then $\langle \mathbf{x}, \mathbf{u} \rangle = \langle \mathbf{x}', \mathbf{u} \rangle \geq \max_{a \in \mathcal{A}} \mathbf{u}[a]$.

329 The left-hand side of (4) reflects the improvement in utility obtained by updating \mathbf{x} to \mathbf{x}' ; in particular, $\max_{a \in \mathcal{A}} \mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle$ is the best-response gap of \mathbf{x} with respect to \mathbf{u} . Lemma 3.3 implies that the utility is monotonically increasing—unless the current strategy is already a best response to \mathbf{u} . Furthermore, so long as the regret vector is *small enough*, the improvement is bound to be substantial, being proportional to the squared best-response gap. It is worth noting that Lemma 3.3 holds no matter the initial regret vector \mathbf{r} , subject to $\mathbf{r} \in \mathbb{R}_{\geq 0}$; this invariance always holds for RM^+ (by definition of the algorithm in Algorithm 2), but that is not so for RM (cf. Lemma 3.6).

336 **Convergence in potential games** We now employ Lemma 3.3 to show that alternating RM^+ quickly converges to approximate Nash equilibria in potential games. Using the fact that the game admits a potential (per Definition 2.2), we have that for any round $t \in [T]$, $\Phi(\mathbf{x}_1^{(t+1)}, \dots, \mathbf{x}_n^{(t+1)}) - \Phi(\mathbf{x}_1^{(t)}, \dots, \mathbf{x}_n^{(t)}) \geq \sum_{i=1}^n \frac{1}{\|\mathbf{r}_i^{(t)}\|_1} \text{BRGap}_i(\mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)})^2 \mathbb{1}\{\text{BRGap}_i(\mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)}) > \epsilon\}$, where we used Lemma 3.3 together with the assumption that only players with more than ϵ best-response gap update their strategies. The telescopic summation over $t = 1, \dots, T$ yields

$$344 \quad \Phi_{\text{range}} \geq \sum_{t=1}^T \sum_{i=1}^n \frac{1}{\|\mathbf{r}_i^{(t)}\|_1} \text{BRGap}_i(\mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)})^2 \mathbb{1}\{\text{BRGap}_i(\mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)}) > \epsilon\}, \quad (5)$$

346 where Φ_{range} denotes the range of the potential function. If in every round $t \in [T]$ there is a player $i \in [n]$ such that $\text{BRGap}_i(\mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)}) > \epsilon$, we have $\Phi_{\text{range}} \geq \sum_{t=1}^T \frac{1}{m\sqrt{t}} \epsilon^2 \geq \frac{1}{m} \epsilon^2 \sqrt{T}$, where we used that $\|\mathbf{r}_i^{(t)}\|_1 \leq \sqrt{m} \|\mathbf{r}_i^{(t)}\|_2 \leq m\sqrt{T}$ (Proposition 2.3). We thus arrive at the following result.

350 **Theorem 3.4.** *In any potential game, ϵ -lazy alternating RM^+ requires at most $1 + \frac{(m\Phi_{\text{range}})^2}{\epsilon^4}$ rounds to converge to an ϵ -Nash equilibrium. More broadly, if $\|\mathbf{r}_i^{(t)}\|_1 \leq C(n, m)t^\alpha$ for all $i \in [n]$ and some $\alpha \in [0, 1/2]$, it requires $1 + \frac{(C(n, m)\Phi_{\text{range}})^\beta}{\epsilon^{2\beta}}$ rounds, where $\beta := 1/1-\alpha$.*

354 This provides a convergence rate of $T^{-1/4}$ to Nash equilibria. Notwithstanding Proposition 3.2, an intriguing aspect of Theorem 3.4 is that it connects convergence to Nash equilibria to the regret of RM^+ . In particular, if RM^+ did not have the no-regret property, meaning that $\|\mathbf{r}_i^{(t)}\|_1 = \Omega(t)$, one could only prove an exponential bound since $\sum_{t=1}^T 1/t = \Theta(\log T)$. At the other end, when each player accumulates constant regret, Theorem 3.4 implies an improved convergence rate of $T^{-1/2}$.

360 One caveat of lazy alternating RM^+ prescribed by Theorem 3.4 is that the desired precision should be known in advance in order to execute the algorithm. We address this limitation in Theorem C.7, where we show that the usual, non-lazy version also converges after $O_\epsilon(1/\epsilon^4)$ rounds, albeit at the cost of introducing an additional dependence in the total number of iterations.

365 **Faster rates using discounting** Next, we refine Theorem 3.4 through the use of *discounted* RM^+ , which means that the regret vector is multiplied by a discount factor $\alpha^{(t)} \in (0, 1]$ in each round; we spell out DRM^+ in Algorithm 3. This class of algorithms was introduced by Brown & Sandholm (2019a), who showed that discounting drastically improves empirical performance in zero-sum games. Our next result shows that DRM^+ with geometric discounting, that is, $\alpha^{(t)} = 1 - \gamma$ for some time-invariant $\gamma \in (0, 1)$, attains a rate of $T^{-1/2}$ to Nash equilibria in potential games; this is considerably faster than the $T^{-1/4}$ rate for alternating RM^+ guaranteed by Theorem 3.4. The basic reason is that DRM^+ —with geometric discounting—maintains the norm of the regret vector bounded by $\sqrt{m/\gamma}$ (Lemma C.2 and Corollary C.3), while still enjoying the one-step improvement property.

374 **Corollary 3.5.** *In any potential game, ϵ -lazy alternating DRM^+ with discount factor $1 - \gamma \in (0, 1)$ requires at most $1 + \frac{m\Phi_{\text{range}}}{\epsilon^2\sqrt{\gamma}}$ rounds to converge to an ϵ -Nash equilibrium.*

377 To our knowledge, this is the first time that discounting yields a provable, worst-case improvement over the bound obtained for non-discounted RM^+ .

Regret matching Before we switch gears to the more general constrained optimization setting, it is instructive to examine the behavior of RM. It turns out that one can adjust [Lemma 3.3](#), but with a crucial caveat: the one-step improvement property is now only *conditional*, as specified below.

Lemma 3.6. *For any $\mathbf{r} \in \mathbb{R}_{\geq 0}^A$ and $\mathbf{u} \in \mathbb{R}^A$, we define $\mathbf{x} := \boldsymbol{\theta} / \|\boldsymbol{\theta}\|_1$, where $\boldsymbol{\theta} := \max(\mathbf{r}, \mathbf{0})$; if $\boldsymbol{\theta} = \mathbf{0}$, $\mathbf{x} \in \Delta(\mathcal{A})$ can be arbitrary. If $\mathbf{r}' := \mathbf{r} + \mathbf{u} - \langle \mathbf{x}, \mathbf{u} \rangle \mathbf{1}$ and $\mathbf{x}' := \boldsymbol{\theta}' / \|\boldsymbol{\theta}'\|_1$, where $\boldsymbol{\theta}' = \max(\mathbf{r}', \mathbf{0}) \neq \mathbf{0}$, then $\langle \mathbf{x}' - \mathbf{x}, \mathbf{u} \rangle \geq \frac{1}{\|\boldsymbol{\theta}'\|_1} \|\boldsymbol{\theta}' - \boldsymbol{\theta}\|_2^2 \geq \frac{1}{\|\boldsymbol{\theta}'\|_1} (\max_{a \in \mathcal{A}} \mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle)^2 \mathbb{1}\{\mathbf{r}[a] \geq 0\}$, where $a \in \arg \max_{a' \in \mathcal{A}} \mathbf{u}[a']$. If $\boldsymbol{\theta}' = \mathbf{0}$, then $\langle \mathbf{x}, \mathbf{u} \rangle = \langle \mathbf{x}', \mathbf{u} \rangle \geq \mathbf{u}[a]$.*

We see that RM’s one-step improvement is conditional on the regret accumulated thus far by a best-response action to be nonnegative. This is not an artifact of our analysis; it alludes to a fundamental discrepancy between RM and RM^+ that will be formally established later on ([Theorem 4.4](#)). The main issue with RM can be seen as follows. If we consider a utility vector $\mathbf{u} = (1, 0)$ and the initial regret vector is, say, $(-R, R)$, it will take RM many iterations—proportionally to the magnitude of $R > 0$ —to finally change strategies, although this will eventually happen with a stationary utility.

3.2 CONSTRAINED OPTIMIZATION AND SIMULTANEOUS UPDATES

We now treat the more general setting where we are maximizing an L -smooth function u .

Single simplex We begin with the special case of a single probability simplex, $\mathcal{X} = \Delta(\mathcal{A})$. Our first goal is to adapt [Lemma 3.3](#). The key challenge is that RM^+ does *not* have a one-step improvement, unlike algorithms such as gradient descent (for a small enough learning rate), even if one initializes RM^+ close to a KKT point. But we observe that if the norm of the regret vector is *large enough*—having small regrets is an obstacle here, in contrast to [Section 3.1](#)—we are guaranteed a one-step improvement in terms of the value of the function ([Lemma 3.7](#)).

To do so, we will use the basic quadratic bound, which yields $u(\mathbf{x}') \geq u(\mathbf{x}) + \langle \nabla u(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle - \frac{L}{2} \|\mathbf{x} - \mathbf{x}'\|_2^2$; we think of \mathbf{x}' as the updated strategy starting from \mathbf{x} . Using a slight refinement of [Lemma 3.3](#), we first have the lower bound $\langle \mathbf{x}' - \mathbf{x}, \nabla u(\mathbf{x}) \rangle \geq \frac{1}{\|\mathbf{r}'\|_1} \|\mathbf{r} - \mathbf{r}'\|_2^2$ ([Lemma C.5](#)).

Also, we observe that $\|\mathbf{x} - \mathbf{x}'\|_1 \leq \|\mathbf{r} - \mathbf{r}'\|_1 \left(\frac{1}{\|\mathbf{r}\|_1} + \frac{1}{\|\mathbf{r}'\|_1} \right)$ ([Lemma C.6](#)). We are now ready to establish a *conditional* one-step improvement when the regret vector has a *sufficiently large* norm.

Lemma 3.7. *Let u be an L -smooth function over $\Delta(\mathcal{A})$. For any $\mathbf{r} \in \mathbb{R}_{\geq 0}^A$ with $\mathbf{r} \neq \mathbf{0}$, we define $\mathbf{x} := \mathbf{r} / \|\mathbf{r}\|_1$. Further, let $\mathbf{r}' := [\mathbf{r} + \nabla u(\mathbf{x}) - \langle \mathbf{x}, \nabla u(\mathbf{x}) \rangle \mathbf{1}]^+ \neq \mathbf{0}$ and $\mathbf{x}' := \mathbf{r}' / \|\mathbf{r}'\|_1$. If $\|\mathbf{r}'\|_2 \geq \max\{2m, 9mL\}$, then*

$$u(\mathbf{x}') - u(\mathbf{x}) \geq \frac{1}{2\|\mathbf{r}'\|_1} \left(\max_{\mathbf{x}^* \in \Delta(\mathcal{A})} \langle \mathbf{x}^* - \mathbf{x}, \nabla u(\mathbf{x}) \rangle \right)^2.$$

[Lemma 3.7](#) only shows a one-step improvement so long as the norm of the regret vector is large enough. But how can we guarantee that? It would seem possible that RM^+ ends up cycling in perpetuity under a regret vector with small norm. The following lemma shows that cannot happen; it turns out that maintaining this monotonicity property for the norm of the regret vector is crucial for designing an optimal variant of RM^+ in zero-sum games ([Zhang et al., 2025](#)).

Lemma 3.8. *For any t , RM^+ guarantees $\|\mathbf{r}^{(t)}\|_2^2 \geq \|\mathbf{r}^{(t-1)}\|_2^2 + \|\mathbf{g}^{(t)}\|_2^2$, where $\mathbf{g}^{(t)} := \nabla u(\mathbf{x}^{(t)}) - \langle \nabla u(\mathbf{x}^{(t)}), \mathbf{x}^{(t)} \rangle \mathbf{1}$ is the instantaneous regret vector at round t .*

In particular,

$$\|\mathbf{r}^{(t)}\|_2^2 \geq \|\mathbf{r}^{(t-1)}\|_2^2 + \|\mathbf{g}^{(t)}\|_2^2 \geq \|\mathbf{r}^{(t-1)}\|_2^2 + \text{KKTGap}(\mathbf{x}^{(t)})^2$$

since $\|\mathbf{g}^{(t)}\|_2^2 \geq \text{KKTGap}(\mathbf{x}^{(t)})^2$. That is, not only is the ℓ_2 norm of the regret vector nondecreasing, but the increase is at least $\text{KKTGap}(\mathbf{x}^{(t)})^2$ at each round $t \in [T]$. Combining with [Lemma 3.7](#) yields the following.

Theorem 3.9 (Single simplex). *Let u be an L -smooth function in $\Delta(\mathcal{A}) \subset \mathbb{R}^m$ with range u_{range} and $R := \max\{2m, 9mL\}$. RM^+ requires at most $1 + \frac{(m(2u_{\text{range}} + R^2))^2}{\epsilon^4}$ rounds to reach an ϵ -KKT point.*

Simultaneous updates in symmetric potential games We now use [Theorem 3.9](#) to prove convergence of *simultaneous* RM^+ in *symmetric* potential games; our earlier result in [Theorem 3.4](#) shows convergence for arbitrary potential games but for the alternating version. The symmetry assumption here means that $\mathcal{A}_1 = \mathcal{A}_1 = \dots = \mathcal{A}_n = \mathcal{A}$ and $\mathbf{u}_1(\mathbf{x}_{-1}) = \mathbf{u}_2(\mathbf{x}_{-2}) = \dots = \mathbf{u}_n(\mathbf{x}_{-n})$ when $\mathbf{x}_1 = \mathbf{x}_2 = \dots = \mathbf{x}_n$. It is further assumed that all players initialize from the same strategy, so that the previous property implies that, inductively, it will be the case that $\mathbf{x}_1^{(t)} = \mathbf{x}_2^{(t)} = \dots = \mathbf{x}_n^{(t)}$ for all t under simultaneous updates because players observe exactly the same utility vector. A simple example of this is a two-player game with a common, symmetric payoff matrix $\mathbf{A} = \mathbf{A}^\top$. Then $\mathbf{u}_1(\mathbf{x}_2) = \mathbf{A}\mathbf{x}_2$ and $\mathbf{u}_2(\mathbf{x}_1) = \mathbf{A}\mathbf{x}_1$, so the previous assumption is satisfied.

Corollary 3.10. *In any symmetric potential game, simultaneous RM^+ converges to an ϵ -Nash equilibrium after $O_\epsilon(1/\epsilon^4)$ rounds. In particular, if convergence to the set of CCE happens at a rate of $T^{-(1-\alpha)}$, for some $\alpha \in [0, 1/2]$, the rate of convergence to Nash equilibria is no worse than $T^{-\frac{1-\alpha}{2}}$.*

Multiple simplices We now have the necessary tools to analyze the general case where we maximize u over a product of simplices. Similarly to [Theorem 3.4](#), we run alternating RM^+ , thinking of every individual simplex as being controlled by a single player; this is akin to coordinate descent.

Corollary 3.11. *If u is an L -smooth function in $\Delta(\mathcal{A}_1) \times \dots \times \Delta(\mathcal{A}_n)$ with range u_{range} , ϵ -lazy alternating RM^+ initialized at $\mathbf{r}_i^{(0)} = \max\{2\sqrt{m_i}, 9\sqrt{m_i}L\}\mathbf{1}$ for each player $i \in [n]$ requires at most $1 + \frac{4n^4 m^2 u_{\text{range}}^2}{\epsilon^4}$ rounds to reach an ϵ -KKT point of u .*

The proof follows directly from [Lemma 3.7](#) together with a telescopic summation. The non-lazy version of RM^+ admits a qualitatively similar bound, following [Theorem C.7](#). Furthermore, a similar bound holds even under simultaneous RM^+ ([Corollary C.11](#)), which follows by extending [Lemma 3.7](#) to multiple simplices ([Lemma C.10](#)).

One caveat of those results is that the regret vector of each player needs to be initialized at a specific threshold. Our next result addresses this limitation by analyzing the usual parameter-free and scale-invariant version of RM^+ , at the cost of introducing a worse dependence on $1/\epsilon$.

Theorem 3.12. *If u is an L -smooth function in $\Delta(\mathcal{A}_1) \times \dots \times \Delta(\mathcal{A}_n)$ with range u_{range} , ϵ -lazy alternating (or simultaneous) RM^+ requires at most $O_\epsilon(1/\epsilon^8)$ rounds to reach an ϵ -KKT point of u .*

4 EXPONENTIAL LOWER BOUNDS FOR REGRET MATCHING

In stark contrast, we show that RM , with or without alternation, can take exponentially many rounds to reach an approximate Nash equilibrium even in two-player identical-interest games. The underlying class of games is based on the one considered by [Panageas et al. \(2023a\)](#), who treated fictitious play. Specifically, for $m = 4, 6, \dots$ and $k \in \mathbb{N}$ we define the matrix $\mathbf{A}_{m,k}$ per the recursion

$$\mathbb{R}^{m \times m} \ni \mathbf{A}_{m,k} := \begin{bmatrix} k+1 & 0 & \dots & 0 & 0 \\ 0 & & & & k+4 \\ \vdots & & \mathbf{A}_{m-2,k+4} & & \vdots \\ 0 & & & & 0 \\ k+2 & 0 & \dots & 0 & k+3 \end{bmatrix}, \text{ where } \mathbf{A}_{2,k} := \begin{bmatrix} k+1 & 0 \\ k+2 & k+3 \end{bmatrix}.$$

(An illustrative example appears in [Section C.2](#).) For any even dimension m , we define $\mathbf{A} := \mathbf{A}_{m,0}$, with maximum entry $2m - 1$. Further, we define, for $1 \leq a_1 \leq m + 1$ and $1 \leq a_2 \leq m + 1$,

$$\mathbf{B}[a_1, a_2] := \begin{cases} \mathbf{A}[a_1, a_2] & \text{if } a_1 \leq m \text{ and } a_2 \leq m; \\ 1/2 & \text{if } (a_1 = m + 1 \text{ and } a_2 = 1) \text{ or } (a_1 = 1 \text{ and } a_2 = m + 1); \\ 0 & \text{otherwise.} \end{cases} \quad (6)$$

The action sets of the two players are $\mathcal{A}_1 = [m + 1] = \mathcal{A}_2$. We assume that RM is initialized at the pure strategy $(m + 1, m + 1)$; [Section C.2](#) shows how to adapt the lower bound when RM is initialized at the uniform random strategy ([Corollary C.16](#)), which is more common. We recall that one round includes one update from each player, which for now is assumed to be made in a simultaneous fashion. For a payoff $k \in \mathbb{N}$, we denote by $a_1(k), a_2(k) \in [m]$ the row and column index, respectively, corresponding to k in the matrix \mathbf{A} .

We begin by stating a basic invariance concerning the behavior of RM when executed on the game (6).

Property 4.1. *After the first round both players play the first action. Thereupon, either the players play with probability 1 ($a_1(k), a_2(k)$), or, when k is odd, only Player 1 (respectively, Player 2 when k is even) mixes between $a_1(k)$ and $a_1(k+1)$ (respectively, $a_2(k)$ and $a_2(k+1)$). If a row or a column stops being played, it will never be played henceforth. An action profile $(a_1(k+1), a_2(k+1))$ is played with positive probability only if $(a_1(k), a_2(k))$ was played at some previous round.*

We prove this property inductively in [Section C.2](#). We take it for granted in what follows.

In accordance with [Property 4.1](#), for $k \geq 2$, we define \underline{t}_k to be the first round in which the action profile corresponding to payoff k is played with positive probability and \bar{t}_k the last round before the action profile corresponding to payoff $k+1$ is played with positive probability. We then define $T_k := \bar{t}_k - \underline{t}_k + 1$ to be the number of rounds corresponding to the period $[\underline{t}_k, \bar{t}_k]$.

We also define $\mathcal{A}_1(k) := \{a_1(k') : 2m-1 \geq k' \geq k\}$ and $\mathcal{A}_2(k) := \{a_2(k') : 2m-1 \geq k' \geq k\}$. These are the rows and columns, respectively, that will be played after the action profile corresponding to k starts being played. The next crucial lemma shows that before an action becomes desirable, it will have accumulated very negative regret in the previous rounds.

Lemma 4.2. *For any even $k \geq 4$, let $\mathbf{r}_1^{(\bar{t}_k - 2)}[a_1]$ be the regret of Player 1 with respect to any action $a_1 \in \mathcal{A}_1(k)$. Then $\mathbf{r}_1^{(\bar{t}_k - 2)}[a_1] \leq -\sum_{l=2}^{k-2} (l-1)T_l$. Similarly, for any odd $k \geq 5$, if $\mathbf{r}_2^{(\bar{t}_k - 2)}[a_2]$ is the regret of Player 2 with respect to any action $a_2 \in \mathcal{A}_2(k)$, $\mathbf{r}_2^{(\bar{t}_k - 2)}[a_2] \leq -\sum_{l=2}^{k-2} (l-1)T_l$.*

At the same time, when an action has very negative regret, it will take a long time before that action gets played with positive probability, as formalized below.

Lemma 4.3. *For any even $k \geq 4$, $T_k \geq -\frac{1}{2}\mathbf{r}_2^{(\bar{t}_k - 1)}[a_2(k+1)]$. Similarly, for every odd $k \geq 5$, $T_k \geq -\frac{1}{2}\mathbf{r}_1^{(\bar{t}_k - 1)}[a_1(k+1)]$.*

By [Lemmas 4.2](#) and [4.3](#), it follows that $T_k \geq \sum_{l=2}^{k-1} \frac{l-1}{2}T_l$ for any $k \geq 4$. By the inductive basis, we know that $T_3 \geq 1$. As a result, $T_k \geq \frac{k-2}{2}T_{k-1} \geq \frac{k-2}{2} \frac{k-3}{2} \dots \frac{2}{2}T_3 \geq \frac{(k-2)!}{2^{k-3}}$ for all $k \geq 4$.

Moreover, it takes at least T_{2m-2} rounds to converge to an NE with approximation gap at most $1/2^{m+2}$ ([Lemma C.15](#)). We thus arrive at the following exponential lower bound.

Theorem 4.4. *Simultaneous RM requires $m^{\Omega(m)}$ rounds to converge to a $\frac{1}{2^m}$ -Nash equilibrium in two-player $m \times m$ identical-interest games.*

The same reasoning directly applies to alternating RM.

Corollary 4.5. *Alternating RM requires $m^{\Omega(m)}$ rounds to converge to a $\frac{1}{2^m}$ -Nash equilibrium in two-player $m \times m$ identical-interest games.*

5 FUTURE RESEARCH

Our paper sheds new light on the convergence properties of regret matching⁽⁺⁾ in constrained optimization problems in general, and potential games in particular. We showed that RM^+ is a sound and fast first-order optimizer; on the flip side, RM can be exponentially slow even in two-player identical-interest games.

Several interesting questions remain open. It would be interesting to understand whether RM^+ , with or without alternation, can experience $\Omega(\sqrt{T})$ regret in potential games; this is known to be the case in zero-sum games ([Farina et al., 2023](#)), but remains unclear for the class of potential games. In light of our results, any improvement over the \sqrt{T} barrier would automatically translate into a faster convergence rate to Nash equilibria. Moreover, we have worked exclusively in the full feedback setting; extending our results under stochastic or bandit feedback would be a natural next step. Finally, does RM asymptotically converge even under alternating updates?

REFERENCES

Anastasios N. Angelopoulos, Michael I. Jordan, and Ryan J. Tibshirani. Gradient equilibrium in online learning: Theory and applications. *arXiv:2501.08330*, 2025.

- 540 Robert Aumann. Subjectivity and correlation in randomized strategies. *Journal of Mathematical*
541 *Economics*, 1:67–96, 1974.
- 542 Yakov Babichenko and Aviad Rubinstein. Settling the complexity of Nash equilibrium in congestion
543 games. In *Proceedings of the Annual Symposium on Theory of Computing (STOC)*, 2021.
- 544 David Blackwell. An analog of the minmax theorem for vector payoffs. *Pacific Journal of Mathe-*
545 *matics*, 6:1–8, 1956.
- 546 Avrim Blum, Eyal Even-Dar, and Katrina Ligett. Routing without regret: On convergence to Nash
547 equilibria of regret-minimizing algorithms in routing games. In *Proceedings of the ACM Sympo-*
548 *sium on Principles of Distributed Computing*, 2006.
- 549 Michael Bowling, Neil Burch, Michael Johanson, and Oskari Tammelin. Heads-up limit hold'em
550 poker is solved. *Science*, 347(6218), January 2015.
- 551 Noam Brown and Tuomas Sandholm. Superhuman AI for heads-up no-limit poker: Libratus beats
552 top professionals. *Science*, pp. eaa01733, Dec. 2017.
- 553 Noam Brown and Tuomas Sandholm. Solving imperfect-information games via discounted regret
554 minimization. In *Conference on Artificial Intelligence (AAAI)*, 2019a.
- 555 Noam Brown and Tuomas Sandholm. Superhuman AI for multiplayer poker. *Science*, 365(6456):
556 885–890, 2019b.
- 557 Yang Cai, Gabriele Farina, Julien Grand-Clément, Christian Kroer, Chung-Wei Lee, Haipeng Luo,
558 and Weiqiang Zheng. Fast last-iterate convergence of learning in games requires forgetful algo-
559 rithms. In *Neural Information Processing Systems (NeurIPS)*, 2024.
- 560 Yang Cai, Gabriele Farina, Julien Grand-Clément, Christian Kroer, Chung-Wei Lee, Haipeng Luo,
561 and Weiqiang Zheng. Last-iterate convergence properties of regret-matching algorithms in games.
562 In *International Conference on Learning Representations (ICLR)*, 2025.
- 563 Volkan Cevher, Ashok Cutkosky, Ali Kavis, Georgios Piliouras, Stratis Skoulakis, and Luca Viano.
564 Alternation makes the adversary weaker in two-player games. In *Neural Information Processing*
565 *Systems (NeurIPS)*, 2023.
- 566 Darshan Chakrabarti, Julien Grand-Clément, and Christian Kroer. Extensive-form game solving via
567 blackwell approachability on treplexes. In *Neural Information Processing Systems (NeurIPS)*,
568 2024.
- 569 Xi Chen, Xiaotie Deng, and Shang-Hua Teng. Settling the complexity of computing two-player
570 Nash equilibria. *Journal of the ACM*, 2009.
- 571 Qiwen Cui, Zhihan Xiong, Maryam Fazel, and Simon S Du. Learning in congestion games with
572 bandit feedback. In *Neural Information Processing Systems (NeurIPS)*, 2022.
- 573 Constantinos Daskalakis, Paul Goldberg, and Christos Papadimitriou. The complexity of computing
574 a Nash equilibrium. *SIAM Journal on Computing*, 2008.
- 575 Aaron Defazio and Konstantin Mishchenko. Learning-rate-free learning by d-adaptation. In *Inter-*
576 *national Conference on Machine Learning (ICML)*, 2023.
- 577 Gabriele Farina, Christian Kroer, and Tuomas Sandholm. Faster game solving via predictive Black-
578 well approachability: Connecting regret matching and mirror descent. In *Conference on Artificial*
579 *Intelligence (AAAI)*, 2021.
- 580 Gabriele Farina, Julien Grand-Clément, Christian Kroer, Chung-Wei Lee, and Haipeng Luo. Regret
581 matching+(in) stability and fast convergence in games. In *Neural Information Processing Systems*
582 *(NeurIPS)*, 2023.
- 583 Yoav Freund and Robert Schapire. Adaptive game playing using multiplicative weights. *Games and*
584 *Economic Behavior*, 29:79–103, 1999.

- 594 Sergiu Hart and Andreu Mas-Colell. A simple adaptive procedure leading to correlated equilibrium.
595 *Econometrica*, 68:1127–1150, 2000.
596
- 597 Sergiu Hart and Andreu Mas-Colell. Regret-based continuous-time dynamics. *Games and Economic*
598 *Behavior*, 45(2):375–394, 2003.
- 599 Elad Hazan, Karan Singh, and Cyril Zhang. Efficient regret minimization in non-convex games. In
600 *International Conference on Machine Learning (ICML)*, 2017.
601
- 602 Amélie Héliou, Johanne Cohen, and Panayotis Mertikopoulos. Learning with bandit feedback in
603 potential games. In *Neural Information Processing Systems (NeurIPS)*, 2017.
604
- 605 Wassily Hoeffding and J. Wolfowitz. Distinguishability of sets of distributions. *The Annals of*
606 *Mathematical Statistics*, 29(3):700–718, 1958.
- 607 Maor Ivgi, Oliver Hinder, and Yair Carmon. Dog is SGD’s best friend: A parameter-free dynamic
608 step size schedule. In *International Conference on Machine Learning (ICML)*, 2023.
609
- 610 Robert Kleinberg, Georgios Piliouras, and Éva Tardos. Multiplicative updates outperform generic
611 no-regret learning in congestion games: extended abstract. In *Proceedings of the Annual Sympos-*
612 *ium on Theory of Computing (STOC)*, 2009.
- 613 Tai-Yu Ma and Philippe Gerber. Distributed regret matching algorithm for dynamic congestion
614 games with information provision. *Transportation Research Procedia*, 3:3–12, 2014.
615
- 616 Jason R. Marden, Gürdal Arslan, and Jeff S. Shamma. Regret based dynamics: convergence in
617 weakly acyclic games. In *Autonomous Agents and Multi-Agent Systems (AAMAS)*, 2007.
- 618 Linjian Meng, Youzhi Zhang, Zhenxing Ge, Tianpei Yang, and Yang Gao. Asynchronous
619 predictive counterfactual regret minimization⁺ algorithm in solving extensive-form games.
620 *arXiv:2503.12770*, 2025.
621
- 622 Dov Monderer and Lloyd S Shapley. Potential games. *Games and Economic Behavior*, 14(1):
623 124–143, 1996.
- 624 Matej Moravčík, Martin Schmid, Neil Burch, Viliam Lisý, Dustin Morrill, Nolan Bard, Trevor
625 Davis, Kevin Waugh, Michael Johanson, and Michael Bowling. Deepstack: Expert-level artificial
626 intelligence in heads-up no-limit poker. *Science*, May 2017.
627
- 628 H. Moulin and J.-P. Vial. Strategically zero-sum games: The class of games whose completely
629 mixed equilibria cannot be improved upon. *International Journal of Game Theory*, 7(3-4):201–
630 221, 1978.
- 631 Tianlong Nan, Shuvomoy Das Gupta, Garud Iyengar, and Christian Kroer. On the $O(1/T)$ conver-
632 gence of alternating gradient descent-ascent in bilinear games. *arXiv:2510.03855*, 2025.
633
- 634 John Nash. Equilibrium points in n-person games. *Proceedings of the National Academy of Sciences*,
635 36:48–49, 1950.
- 636 Francesco Orabona and Dávid Pál. Coin betting and parameter-free online learning. In *Neural*
637 *Information Processing Systems (NIPS)*, 2016.
638
- 639 Gerasimos Palaiopoulos, Ioannis Panageas, and Georgios Piliouras. Multiplicative weights update
640 with constant step-size in congestion games: Convergence, limit cycles and chaos. In *Neural*
641 *Information Processing Systems (NeurIPS)*, 2017.
- 642 Ioannis Panageas, Nikolas Patris, Stratis Skoulakis, and Volkan Cevher. Exponential lower bounds
643 for fictitious play in potential games. In *Neural Information Processing Systems (NeurIPS)*,
644 2023a.
645
- 646 Ioannis Panageas, Stratis Skoulakis, Luca Viano, Xiao Wang, and Volkan Cevher. Semi bandit
647 dynamics in congestion games: Convergence to Nash equilibrium and no-regret guarantees. In
International Conference on Machine Learning (ICML), 2023b.

648 Robert W Rosenthal. A class of games possessing pure-strategy nash equilibria. *International*
649 *Journal of Game Theory*, 2(1):65–67, 1973.

650 Oskari Tammelin. Solving large imperfect information games using CFR+. *arXiv:1407.5042*, 2014.

651 Emanuel Tewolde, Brian Hu Zhang, Ioannis Anagnostides, Tuomas Sandholm, and Vince Conitzer.
652 Decision making under imperfect recall: Algorithms and benchmarks. In *Uncertainty in Artificial*
653 *Intelligence (UAI)*, 2025.

654 Andre Wibisono, Molei Tao, and Georgios Piliouras. Alternating mirror descent for constrained
655 min-max games. In *Neural Information Processing Systems*, 2022.

656 Hang Xu, Kai Li, Haobo Fu, Qiang Fu, Junliang Xing, and Jian Cheng. Dynamic discounted coun-
657 terfactual regret minimization. In *International Conference on Learning Representations (ICLR)*,
658 2024a.

659 Hang Xu, Kai Li, Bingyun Liu, Haobo Fu, Qiang Fu, Junliang Xing, and Jian Cheng. Minimizing
660 weighted counterfactual regret with optimistic online mirror descent. *arXiv:2404.13891*, 2024b.

661 Brian Hu Zhang and Tuomas Sandholm. General search techniques without common knowl-
662 edge for imperfect-information games, and application to superhuman fog of war chess.
663 *arXiv:2506.01242*, 2025.

664 Brian Hu Zhang, Ioannis Anagnostides, and Tuomas Sandholm. Scale-invariant regret matching
665 and online learning with optimal convergence: Bridging theory and practice in zero-sum games.
666 *arXiv:2510.04407*, 2025.

667 Naifeng Zhang, Stephen McAleer, and Tuomas Sandholm. Faster game solving via hyperparameter
668 schedules. *arXiv:2404.09097*, 2024.

669 Martin Zinkevich, Michael Bowling, Michael Johanson, and Carmelo Piccione. Regret minimization
670 in games with incomplete information. In *Neural Information Processing Systems (NIPS)*, 2007.

671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

702 A FURTHER RELATED WORK

703
704
705 Much of the existing research on regret matching revolves around zero-sum games. Many variants
706 have been proposed over the years to speed up its convergence (Xu et al., 2024b; Cai et al., 2025;
707 Chakrabarti et al., 2024; Meng et al., 2025; Farina et al., 2021; Tammelin, 2014; Brown & Sandholm,
708 2019a). Some notable variations that have considerably improved performance are *predictive* RM
709 and RM^+ (Farina et al., 2021), which rely on predicting the next utility, and *discounted* RM and
710 RM^+ (Brown & Sandholm, 2019a; Zhang et al., 2024; Xu et al., 2024a), where one dynamically
711 discounts the accumulated regret; in a similar vein, our work shows that a discounted variant of
712 RM^+ achieves a better convergence upper bound than RM^+ in our setting (Corollary 3.5). Moreover,
713 alternation is known to speed up performance, at least in zero-sum games (Tammelin, 2014), and
714 has been the subject of much recent research (Wibisono et al., 2022; Cevher et al., 2023; Nan et al.,
715 2025). It must be stressed that the focus of all that prior work was on zero-sum games. Constrained
716 optimization is a fundamentally different problem. For one, in zero-sum games, it is only the average
717 strategy of RM and RM^+ that converges, not the last iterate (Farina et al., 2023).

718 The recent paper of Tewolde et al. (2025) demonstrated that the regret matching family is a
719 formidable first-order optimizer in constrained optimization problems. In particular, their focus was
720 on (single-player) imperfect-recall problems, which are tantamount to general polynomial optimiza-
721 tion problems over a product of simplices. Interestingly, many of the trends observed in zero-sum
722 games are actually reversed in constrained optimization. For example, the predictive versions of
723 RM and RM^+ generally performed worse than their non-predictive counterparts. One trend that did
724 persist was the superiority of RM^+ over RM. It is also worth mentioning an earlier work by Ma &
725 Gerber (2014) that also reported fast empirical convergence in a certain class of congestion games.
726 Yet, there was hitherto no theoretical understanding of those algorithms in this setting. The main
727 precursors of our work are the paper of Hart & Mas-Colell (2003), which established asymptotic
728 convergence in discrete time but for a somewhat artificial variant of regret matching, and the paper
729 of Marden et al. (2007), which analyzed asymptotically a certain variant of regret matching that
730 aggressively discounts the regrets (cf. Corollary 3.5).

731 An interesting result that sheds light on RM and RM^+ is by Farina et al. (2021), who showed that RM
732 can be obtained by running *follow the regularized leader* (FTRL) in a certain lifted space, whereas
733 RM^+ can be obtained through *mirror descent* (MD) in the same space; this is despite the fact that, un-
734 like FTRL and MD, RM and RM^+ are both parameter free. On a related note, Cai et al. (2024) showed
735 that only forgetful algorithms—closer to MD than to FTRL—can attain fast last-iterate convergence.
736 Our exponential separation of RM and RM^+ echoes their finding, although in a different setting and
737 class of algorithms.

738 Zooming out of the RM family, understanding the convergence of no-regret dynamics in potential
739 games has been a popular research topic (Kleinberg et al., 2009; Héliou et al., 2017; Palaiopanos
740 et al., 2017; Panageas et al., 2023b; Cui et al., 2022; Blum et al., 2006). Our research also relates
741 to parameter-free optimization; for example, we refer to Ivgi et al. (2023); Orabona & Pál (2016);
742 Defazio & Mishchenko (2023) and references therein.

743 B FURTHER BACKGROUND

744
745 **Coarse correlated equilibria** For completeness, we provide the definition of a coarse correlated
746 equilibrium (Moulin & Vial, 1978), which is a relaxation of correlated equilibria (Aumann, 1974).
747 The key connection that relates to our results is that if all players in a normal-form game have
748 sublinear regret, the average correlated distribution of play converges to the set of coarse correlated
749 equilibria. In particular, the rate of convergence is driven by the maximum of the players’ regrets
750 (Proposition B.2).

751 **Definition B.1** (Coarse correlated equilibrium). Consider an n -player game in normal form. A
752 correlated distribution $\mu \in \Delta(\mathcal{A}_1 \times \dots \times \mathcal{A}_n)$ is an ϵ -*coarse correlated equilibrium* (CCE) if for
753 any player $i \in [n]$ and deviation $a'_i \in \mathcal{A}_i$,

$$754 \mathbb{E}_{(a_1, \dots, a_n) \sim \mu} u_i(a_1, \dots, a_n) \geq \mathbb{E}_{(a_1, \dots, a_n) \sim \mu} u_i(a'_i, a_{-i}) - \epsilon. 755$$

Proposition B.2. *If each player $i \in [n]$ observes the sequence of utilities $(\mathbf{u}_i(\mathbf{x}_{-i}^{(t)}))_{t=1}^T$, the average correlated distribution of play is an ϵ -CCE with $\epsilon \leq \frac{1}{T} \max_{1 \leq i \leq n} \text{Reg}_i^{(T)}$, where $\text{Reg}_i^{(T)}$ is the regret of the i th player.*

This connection holds for simultaneous updates. It is unclear if and how it can be extended under alternating updates. For the special case of potential games and RM^+ , which is our main focus here, we are indeed able to establish convergence to the set of CCEs even under alternating updates by bounding the path length of the players' strategies (Remark C.9).

Other notions of regret Section 2.2 introduced the usual notion of regret used in online linear optimization. For a constrained optimization problem with respect to a differentiable function u , we have $\text{Reg}^{(T)} = \max_{\mathbf{x}' \in \mathcal{X}} \sum_{t=1}^T \langle \mathbf{x}' - \mathbf{x}^{(t)}, \nabla_{\mathbf{x}} u(\mathbf{x}^{(t)}) \rangle$. This is a linearized version of $\max_{\mathbf{x}' \in \mathcal{X}} \sum_{t=1}^T (u(\mathbf{x}') - u(\mathbf{x}^{(t)}))$. Minimizing the latter notion is computationally intractable unless one places restrictive assumptions on u . Hazan et al. (2017) (cf. Angelopoulos et al., 2025) introduced the notion of "local regret," and showed that having sublinear local regret implies that a randomly selected iterate will be an approximate stationary point; this is in stark contrast to regret as defined in Section 2.2 (Proposition 3.2).

Discounting Next, we spell out regret matching⁺ with discounting (DRM^+ ; Algorithm 3). The only difference from RM^+ is that the regret vector is multiplied by a discounting coefficient $\alpha^{(t)} \in (0, 1]$ in every round $t \in [T]$ (Algorithm 3); the special case where $\alpha^{(t)} = 1$ for all $t \in [T]$ is RM^+ .

Algorithm 3: Regret matching⁺ with discounting (DRM^+)

- 1 **Input:** discounting coefficients $(\alpha^{(1)}, \dots, \alpha^{(T)}) \in (0, 1]^T$;
 - 2 Initialize cumulative regrets $\mathbf{r}^{(0)} := \mathbf{0}$;
 - 3 Initialize strategy $\mathbf{x}^{(1)} \in \Delta(\mathcal{A})$;
 - 4 **for** $t = 1, \dots, T$ **do**
 - 5 Set $\boldsymbol{\theta}^{(t)} \leftarrow \mathbf{r}^{(t-1)}$;
 - 6 **if** $\boldsymbol{\theta}^{(t)} \neq \mathbf{0}$ **then**
 - 7 Compute $\mathbf{x}^{(t)} \leftarrow \boldsymbol{\theta}^{(t)} / \|\boldsymbol{\theta}^{(t)}\|_1$;
 - 8 **else**
 - 9 $\mathbf{x}^{(t)} \leftarrow \mathbf{x}^{(t-1)}$;
 - 10 Output strategy $\mathbf{x}^{(t)} \in \Delta(\mathcal{A})$;
 - 11 Observe utility $\mathbf{u}^{(t)} \in \mathbb{R}^{\mathcal{A}}$;
 - 12 $\mathbf{r}^{(t)} \leftarrow \alpha^{(t)} [\mathbf{r}^{(t-1)} + \mathbf{u}^{(t)} - \langle \mathbf{x}^{(t)}, \mathbf{u}^{(t)} \rangle \mathbf{1}]^+$;
-

C OMITTED PROOFS

This section provides the proofs missing from the main body. We begin by stating a simple lemma that bounds the regret of RM^+ , implying Proposition 2.3; we will then adapt it to account for discounting per Algorithm 3.

Lemma C.1 (Regret vector upper bound). *For any time $t \in [T]$, RM^+ guarantees $\|\mathbf{r}^{(t)}\|_2^2 \leq \|\mathbf{r}^{(t-1)}\|_2^2 + \|\mathbf{g}^{(t)}\|_2^2$, where $\mathbf{g}^{(t)} := \mathbf{u}^{(t)} - \langle \mathbf{x}^{(t)}, \mathbf{u}^{(t)} \rangle$ is the instantaneous regret at time t .*

Proof. By definition of RM^+ , $\langle \mathbf{r}^{(t-1)}, \mathbf{g}^{(t)} \rangle = \langle \mathbf{x}^{(t)}, \mathbf{g}^{(t)} \rangle = 0$ since $\mathbf{x}^{(t)} \propto \mathbf{r}^{(t-1)}$. Thus,

$$\|\mathbf{r}^{(t)}\|_2^2 = \|\mathbf{r}^{(t-1)} + \mathbf{g}^{(t)}\|_2^2 \leq \|\mathbf{r}^{(t-1)}\|_2^2 + \|\mathbf{g}^{(t)}\|_2^2,$$

by orthogonality. □

As a result, the telescopic summation yields $\|\mathbf{r}^{(T)}\|_2^2 \leq \sum_{t=1}^T \|\mathbf{g}^{(t)}\|_2^2 \leq mT$ since $\|\mathbf{g}^{(t)}\|_\infty \leq 1$ (by the assumption that the range of the utilities is bounded by 1). It is worth noting that a similar proof works for RM . We now adapt Lemma C.1 for DRM^+ .

Lemma C.2. *For any time $t \in [T]$, DRM^+ guarantees $\|\mathbf{r}^{(t)}\|_2^2 \leq (\alpha^{(t)})^2 (\|\mathbf{r}^{(t-1)}\|_2^2 + \|\mathbf{g}^{(t)}\|_2^2)$.*

810 *Proof.* As before, $\langle \mathbf{r}^{(t-1)}, \mathbf{g}^{(t)} \rangle = \langle \mathbf{x}^{(t)}, \mathbf{g}^{(t)} \rangle = 0$ since $\mathbf{x}^{(t)} \propto \mathbf{r}^{(t-1)}$. Thus,
 811 $\|\mathbf{r}^{(t)}\|_2^2 = (\alpha^{(t)})^2 \|\mathbf{r}^{(t-1)} + \mathbf{g}^{(t)}\|_2^2 \leq (\alpha^{(t)})^2 \|\mathbf{r}^{(t-1)} + \mathbf{g}^{(t)}\|_2^2 \leq (\alpha^{(t)})^2 (\|\mathbf{r}^{(t-1)}\|_2^2 + \|\mathbf{g}^{(t)}\|_2^2)$.
 812 \square

813
 814
 815 A direct consequence is the following bound on the regret vector.

816 **Corollary C.3.** *For any time $t \in [T]$, DRM^+ guarantees*

$$817 \|\mathbf{r}^{(t)}\|_2^2 \leq (\alpha^{(t)})^2 \|\mathbf{g}^{(t)}\|_2^2 + (\alpha^{(t)} \alpha^{(t-1)})^2 \|\mathbf{g}^{(t-1)}\|_2^2 + \dots + \left(\prod_{\tau=1}^t \alpha^{(\tau)} \right)^2 \|\mathbf{g}^{(1)}\|_2^2. \quad (7)$$

818
 819
 820 *In particular, if $\alpha^{(t)} = 1 - \gamma$ for some constant $\gamma \in (0, 1)$, it follows that $\|\mathbf{r}^{(T)}\|_2 \leq \sqrt{m/\gamma}$.*

821
 822 *Proof.* The first part of the claim follows by unfolding the bound of [Lemma C.2](#). For the second
 823 part, using (7) and the fact that $\|\mathbf{g}^{(t)}\|_2^2 \leq m$ for any t , we have

$$824 \|\mathbf{r}^{(t)}\|_2^2 \leq m \left((1 - \gamma)^2 + (1 - \gamma)^4 + \dots + (1 - \gamma)^{2t} \right) \leq m \frac{(1 - \gamma)^2}{1 - (1 - \gamma)^2} \leq m \frac{1}{\gamma}.$$

825
 826
 827 \square

828 C.1 PROOFS FROM SECTION 3

829
 830 We continue with the proofs from [Section 3](#). We first establish that RM^+ enjoys a one-step improve-
 831 ment property when the utility is linear.

832 **Lemma 3.3** (One-step improvement for RM^+). *For any $\mathbf{r} \in \mathbb{R}_{\geq 0}^{\mathcal{A}}$ and $\mathbf{u} \in \mathbb{R}^{\mathcal{A}}$, we define $\mathbf{x} :=$
 833 $\mathbf{r}/\|\mathbf{r}\|_1$; if $\mathbf{r} = \mathbf{0}$, $\mathbf{x} \in \Delta(\mathcal{A})$ can be arbitrary. If $\mathbf{r}' := [\mathbf{r} + \mathbf{u} - \langle \mathbf{x}, \mathbf{u} \rangle \mathbf{1}]^+ \neq \mathbf{0}$ and $\mathbf{x}' := \mathbf{r}'/\|\mathbf{r}'\|_1$,
 834 then*

$$835 \langle \mathbf{x}' - \mathbf{x}, \mathbf{u} \rangle \geq \frac{1}{\|\mathbf{r}'\|_1} \left(\max_{a \in \mathcal{A}} \mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle \right)^2 = \frac{1}{\|\mathbf{r}'\|_1} \text{BRGap}(\mathbf{x}, \mathbf{u})^2. \quad (4)$$

836
 837 *If $\mathbf{r}' = \mathbf{0}$, then $\langle \mathbf{x}, \mathbf{u} \rangle = \langle \mathbf{x}', \mathbf{u} \rangle \geq \max_{a \in \mathcal{A}} \mathbf{u}[a]$.*

838
 839 *Proof.* First, if $\mathbf{r}' = \mathbf{0}$, it follows that $\mathbf{r} + \mathbf{u} - \langle \mathbf{x}, \mathbf{u} \rangle \mathbf{1} \leq \mathbf{0}$, where the inequality is to be taken
 840 coordinate-wise. Since $\mathbf{r} \geq \mathbf{0}$, we have $\langle \mathbf{x}, \mathbf{u} \rangle \geq \mathbf{u}[a]$ for all $a \in \mathcal{A}$, as claimed.

841
 842 We now assume $\mathbf{r}' \neq \mathbf{0}$. If $\mathbf{r} = \mathbf{0}$, we have $\mathbf{r}' = [\mathbf{u} - \langle \mathbf{x}, \mathbf{u} \rangle \mathbf{1}]^+$. (4) can then be equivalently
 843 expressed as

$$844 \sum_{a \in \mathcal{A}} \mathbf{r}'[a] (\mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle) \geq \left(\max_{a \in \mathcal{A}} \mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle \right)^2,$$

845
 846 which holds since $\mathbf{r}' = [\mathbf{u} - \langle \mathbf{x}, \mathbf{u} \rangle \mathbf{1}]^+$. So we can assume $\mathbf{r} \neq \mathbf{0}$. We define $\delta := \mathbf{r}' - \mathbf{r}$. (4) can
 847 be expressed as

$$848 \frac{\sum_{a \in \mathcal{A}} (\mathbf{r}[a] + \delta[a]) \mathbf{u}[a]}{\sum_{a' \in \mathcal{A}} (\mathbf{r}[a'] + \delta[a'])} \geq \frac{\sum_{a \in \mathcal{A}} \mathbf{r}[a] \mathbf{u}[a]}{\sum_{a' \in \mathcal{A}} \mathbf{r}[a']} + \frac{(\max_{a \in \mathcal{A}} \mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle)^2}{\sum_{a' \in \mathcal{A}} (\mathbf{r}[a'] + \delta[a'])}.$$

849
 850 Equivalently,

$$851 \sum_{a' \in \mathcal{A}} \mathbf{r}[a'] \sum_{a \in \mathcal{A}} (\mathbf{r}[a] + \delta[a]) \mathbf{u}[a] \geq \sum_{a \in \mathcal{A}} \mathbf{r}[a] \sum_{a' \in \mathcal{A}} (\mathbf{r}[a'] + \delta[a']) \mathbf{u}[a]$$

$$852 + \sum_{a' \in \mathcal{A}} \mathbf{r}[a'] \left(\max_{a \in \mathcal{A}} \mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle \right)^2.$$

853
 854 This in turn is equivalent to

$$855 \sum_{a' \in \mathcal{A}} \mathbf{r}[a'] \sum_{a \in \mathcal{A}} \delta[a] \mathbf{u}[a] \geq \sum_{a \in \mathcal{A}} \mathbf{r}[a] \sum_{a' \in \mathcal{A}} \delta[a'] \mathbf{u}[a] + \sum_{a' \in \mathcal{A}} \mathbf{r}[a'] \left(\max_{a \in \mathcal{A}} \mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle \right)^2$$

$$856 = \sum_{a' \in \mathcal{A}} \delta[a'] \sum_{a \in \mathcal{A}} \mathbf{r}[a] \langle \mathbf{x}, \mathbf{u} \rangle + \sum_{a' \in \mathcal{A}} \mathbf{r}[a'] \left(\max_{a \in \mathcal{A}} \mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle \right)^2.$$

Rearranging,

$$\sum_{a' \in \mathcal{A}} r[a'] \sum_{a \in \mathcal{A}} \delta[a](\mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle) \geq \sum_{a' \in \mathcal{A}} r[a'] \left(\max_{a \in \mathcal{A}} \mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle \right)^2.$$

Now, for any $a \in \mathcal{A}$ such that $\mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle \geq 0$, it follows that $\delta[a] = \mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle \geq 0$; on the other hand, for $a \in \mathcal{A}$ such that $\mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle < 0$, we have $\delta[a] \leq 0$. That is, $\delta[a](\mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle) \geq 0$, and the claim follows. \square

We will now show that [Lemma 3.3](#) is, in a certain sense, tight. We consider a simple linear maximization over the simplex. If the regret vector of RM^+ can be initialized arbitrarily, as is the premise in [Lemma 3.3](#), we make the following observation.

Lemma C.4. *Consider a utility vector $\mathbf{u} \in \mathbb{R}^A$ and some initial regret vector $\mathbb{R}_{\geq 0}^A \ni \mathbf{r}^{(1)} \neq \mathbf{0}$. If $\mathbf{x}^{(1)} = \mathbf{r}^{(1)} / \|\mathbf{r}^{(1)}\|_1$ and $\epsilon = \max_{a \in \mathcal{A}} \mathbf{u}[a] - \langle \mathbf{x}^{(1)}, \mathbf{u} \rangle$ is the initial best-response gap, it takes at least $\|\mathbf{r}^{(1)}\|_1 / 2\epsilon$ iterations for RM^+ to reach a point $\mathbf{x}^{(t)}$ with best-response gap at most $\epsilon/2$.*

Indeed, we consider the two-dimensional problem in which $\mathbf{u} = (1 - \epsilon, 1)$ and $\mathbf{r}^{(1)} = (\|\mathbf{r}^{(1)}\|_1, 0)$. To incur a best-response gap of at most $\epsilon/2$, the player needs to allot a probability mass of at least $1/2$ to the second action. In the meantime, the decrement of the first coordinate of $\mathbf{r}^{(t)}$ will be at most ϵ while the increment of the second coordinate of $\mathbf{r}^{(t)}$ will be at most ϵ . But, for the algorithm to terminate, it must be the case that the second coordinate of $\mathbf{r}^{(t)}$ is at least as large as the first coordinate of $\mathbf{r}^{(t)}$, leading to [Lemma C.4](#).

Now, given that $\langle \mathbf{x}^{(t)} - \mathbf{x}^{(1)}, \mathbf{u} \rangle \leq \epsilon$, [Lemma C.4](#) matches the bound obtained for this problem through [Lemma 3.3](#) in the regime where $\|\mathbf{r}^{(1)}\|_1$ is at least as large as $1/\epsilon$ (so that the norm of $\mathbf{r}^{(t)}$ is within a constant factor of that of $\mathbf{r}^{(1)}$, by [Lemma C.1](#)).

Taking this argument a step further, if we have a regret bound of the form $\|\mathbf{r}^{(t)}\|_1 \leq \|\mathbf{r}\|_1 = O_t(1)$, which holds when the utility is fixed, [Lemma 3.3](#) implies that $2\|\mathbf{r}\|_1 + 4\|\mathbf{r}\|_1 + \dots + \|\mathbf{r}\|_1/\epsilon = O_\epsilon(1/\epsilon)$ iterations suffice for RM^+ to have a best-response gap at most ϵ when facing a fixed utility; this follows by applying [Lemma 3.3](#) first for all iterations in which the best-response gap is at least $1/2$, then for all iterations in which it is at least $1/4$, and so forth.

Analysis of non-lazy alternating RM^+ In the main body, we used [Lemma 3.3](#) to argue that *lazy* alternating RM^+ converges in potential games ([Theorem 3.4](#)). One caveat of the lazy version of alternation is that it requires knowing the desired precision ϵ ahead of time. We will now extend the analysis to encompass the usual version of alternation, albeit at the cost of introducing a further dependence on the iteration bound.

To begin with, we first state a direct refinement of [Lemma 3.3](#) that we rely on; this refinement will also be useful in the more general setting of constrained optimization.

Lemma C.5 (Refinement of [Lemma 3.3](#)). *Under the preconditions of [Lemma 3.3](#),*

$$\langle \mathbf{x}' - \mathbf{x}, \mathbf{u} \rangle \geq \frac{1}{\|\mathbf{r}'\|_1} \|\mathbf{r} - \mathbf{r}'\|_2^2. \quad (8)$$

In particular, (8) implies (4) since $\|\mathbf{r} - \mathbf{r}'\|_2^2 \geq \|\mathbf{r} - \mathbf{r}'\|_\infty^2 \geq (\max_{a \in \mathcal{A}} \mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle)^2$, by definition of \mathbf{r}' . The proof of [Lemma C.5](#) is identical to that of [Lemma 3.3](#).

The next elementary lemma shows that, so long as the norm of the regret vector is not too small, closeness in regrets implies closeness in strategies.

Lemma C.6. *For $\mathbb{R}_{\geq 0}^A \ni \mathbf{r}, \mathbf{r}' \neq \mathbf{0}$, let $\mathbf{x} := \mathbf{r} / \|\mathbf{r}\|_1$ and $\mathbf{x}' := \mathbf{r}' / \|\mathbf{r}'\|_1$. Then*

$$\|\mathbf{x} - \mathbf{x}'\|_1 \leq \|\mathbf{r} - \mathbf{r}'\|_1 \left(\frac{1}{\|\mathbf{r}\|_1} + \frac{1}{\|\mathbf{r}'\|_1} \right).$$

Proof. The term $\mathbf{x}[a] - \mathbf{x}'[a]$ can be expressed, for any $a \in \mathcal{A}$, as

$$\begin{aligned} \frac{\mathbf{r}[a]}{\sum_{a' \in \mathcal{A}} \mathbf{r}[a']} - \frac{\mathbf{r}'[a]}{\sum_{a' \in \mathcal{A}} \mathbf{r}'[a']} &= \frac{\sum_{a' \in \mathcal{A}} (\mathbf{r}[a] \mathbf{r}'[a'] - \mathbf{r}'[a] \mathbf{r}[a'])}{\|\mathbf{r}\|_1 \|\mathbf{r}'\|_1} \\ &= \frac{\sum_{a' \in \mathcal{A}} (\mathbf{r}[a] (\mathbf{r}'[a'] - \mathbf{r}[a']) + \mathbf{r}[a'] (\mathbf{r}[a] - \mathbf{r}'[a]))}{\|\mathbf{r}\|_1 \|\mathbf{r}'\|_1}, \end{aligned}$$

and the claim follows. \square

With those two helper lemmas in hand, we are now ready to analyze (non-lazy) alternating RM^+ in potential games.

Theorem C.7. *In any potential game with utilities in $[0, 1]$, alternating RM^+ requires at most $2 \lceil \frac{625m^4 n^4 \Phi_{\text{range}}^2}{\delta^4 \epsilon^4} \rceil$ rounds to converge to an ϵ -Nash equilibrium, where $\delta_i := \text{BRGap}_i(\mathbf{x}_i^{(1)}, \mathbf{u}_i^{(1)}) > 0$ and $\delta := \min_{1 \leq i \leq n} \delta_i$.*

The assumption that $\delta_i > 0$ is without any loss in the following sense. The analysis requires that $\|\mathbf{r}_i^{(t)}\|_2 > 0$ for all $i \in [n]$ and any sufficiently large t . By Lemma 3.8, it suffices if a player incurs a positive best-response gap at some time. In the contrary case, if a player always has zero best-response gap, it will always play the same strategy (by definition of RM^+ in Algorithm 2), thereby reducing to a potential game with $n - 1$ players. We further remark that, in accordance with Theorem 3.4, the bound in Theorem C.7 can be similarly parameterized in terms of the maximum regret incurred by a player. A more pedantic point about Theorem C.7 is that utilities are taken to be in $[0, 1]$, while in the rest of the paper we only assume that the range is bounded by 1; this innocuous assumption is used in Claim C.8 below.

Proof of Theorem C.7. By Lemma C.5, the telescopic summation after T rounds yields

$$\Phi_{\text{range}} \geq \sum_{t=1}^T \sum_{i=1}^n \frac{1}{\|\mathbf{r}_i^{(t)}\|_1} \|\mathbf{r}_i^{(t)} - \mathbf{r}_i^{(t-1)}\|_2^2 \geq \sum_{t=1}^T \frac{1}{\max_i \|\mathbf{r}_i^{(t)}\|_1} \sum_{i=1}^n \|\mathbf{r}_i^{(t)} - \mathbf{r}_i^{(t-1)}\|_2^2. \quad (9)$$

Since $\|\mathbf{r}_i^{(t)}\|_1 \leq m\sqrt{t}$ for any player $i \in [n]$, it follows that $\sum_{t=1}^T \sum_{i=1}^n \|\mathbf{r}_i^{(t+1)} - \mathbf{r}_i^{(t)}\|_2^2 \leq m\Phi_{\text{range}}\sqrt{T}$. As a result, after at most $2 \lceil m^2 \Phi_{\text{range}}^2 / \epsilon^4 \rceil$ rounds, there will be a time t such that $\sum_{i=1}^n \|\mathbf{r}_i^{(t)} - \mathbf{r}_i^{(t-1)}\|_2^2 + \sum_{i=1}^n \|\mathbf{r}_i^{(t+1)} - \mathbf{r}_i^{(t)}\|_2^2 \leq \epsilon^2$, which in turn implies $\|\mathbf{r}_i^{(t)} - \mathbf{r}_i^{(t-1)}\|_2, \|\mathbf{r}_i^{(t+1)} - \mathbf{r}_i^{(t)}\|_2 \leq \epsilon$ for all $i \in [n]$. Now, by assumption, we know that $\|\mathbf{r}_i^{(1)}\|_2 \geq \delta_i > 0$ for all $i \in [n]$. By the monotonicity property of the regret vector (Lemma 3.8, proven in the sequel), it follows that $\|\mathbf{r}_i^{(t)}\|_2 \geq \delta_i > 0$ for all $t \in [T + 1]$ and $i \in [n]$. Consequently, applying Lemma C.6, we have $\|\mathbf{x}_i^{(t+1)} - \mathbf{x}_i^{(t)}\|_1 \leq 2\|\mathbf{r}_i^{(t+1)} - \mathbf{r}_i^{(t)}\|_1 / \delta_i \leq 2\sqrt{m}\epsilon / \delta_i$ since $\|\cdot\|_1 \leq \sqrt{m}\|\cdot\|_2$. Next, we will make use of the following simple claim.

Claim C.8. *Consider a normal-form game with utilities in $[0, 1]$. For any two joint strategies $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ and $(\mathbf{x}'_1, \dots, \mathbf{x}'_n)$, it holds that $\|\mathbf{u}_i(\mathbf{x}_{-i}) - \mathbf{u}_i(\mathbf{x}'_{-i})\|_\infty \leq \sum_{i' \neq i} \|\mathbf{x}_{i'} - \mathbf{x}'_{i'}\|_1$ for any player $i \in [n]$.*

Proof. We fix a player $i \in [n]$. We have

$$\begin{aligned} \|\mathbf{u}_i(\mathbf{x}_{-i}) - \mathbf{u}_i(\mathbf{x}'_{-i})\|_\infty &= \left| \sum_{\mathbf{a} \in \mathcal{A}_1 \times \dots \times \mathcal{A}_n} u_i(\cdot, \mathbf{a}_{-i}) \left(\prod_{i' \neq i} \mathbf{x}_{i'}[a_{i'}] - \prod_{i' \neq i} \mathbf{x}'_{i'}[a_{i'}] \right) \right| \\ &\leq \left| \sum_{\mathbf{a} \in \mathcal{A}_1 \times \dots \times \mathcal{A}_n} \left(\prod_{i' \neq i} \mathbf{x}_{i'}[a_{i'}] - \prod_{i' \neq i} \mathbf{x}'_{i'}[a_{i'}] \right) \right| \\ &\leq \sum_{i' \neq i} \|\mathbf{x}_{i'} - \mathbf{x}'_{i'}\|_1, \end{aligned} \quad (10)$$

$$\leq \sum_{i' \neq i} \|\mathbf{x}_{i'} - \mathbf{x}'_{i'}\|_1, \quad (11)$$

where (10) uses triangle inequality together with the assumption that $|u_i(\cdot)| \leq 1$, and (11) uses a bound on the total variation distance of a product distribution in terms of the sum of the total variation distances of its marginals (Hoeffding & Wolfowitz, 1958). \square

Using this lemma, we now observe that $\|\mathbf{u}_i^{(t)} - \mathbf{u}_i(\mathbf{x}_{-i}^{(t)})\|_\infty \leq \sum_{i' < i} \|\mathbf{x}_{i'}^{(t+1)} - \mathbf{x}_{i'}^{(t)}\|_1 \leq \sum_{i' \neq i} \|\mathbf{x}_{i'}^{(t+1)} - \mathbf{x}_{i'}^{(t)}\|_1 \leq 2n\sqrt{m}\epsilon/\delta$ for each $i \in [n]$, where $\delta = \min_{1 \leq i \leq n} \delta_i$. Furthermore, in view of the fact that $\|\mathbf{r}_i^{(t)} - \mathbf{r}_i^{(t-1)}\|_2 \leq \epsilon$, it follows that $\|\mathbf{u}_i^{(t)}\|_\infty - \langle \mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)} \rangle \leq \epsilon$. Thus,

$$\begin{aligned} \|\mathbf{u}_i(\mathbf{x}_{-i}^{(t)})\|_\infty - \langle \mathbf{x}_i^{(t)}, \mathbf{u}_i(\mathbf{x}_{-i}^{(t)}) \rangle &= \|\mathbf{u}_i^{(t)}\|_\infty - \langle \mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)} \rangle \\ &\quad + \|\mathbf{u}_i(\mathbf{x}_{-i}^{(t)})\|_\infty - \|\mathbf{u}_i^{(t)}\|_\infty + \langle \mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)} - \mathbf{u}_i(\mathbf{x}_{-i}^{(t)}) \rangle \\ &\leq \|\mathbf{u}_i^{(t)}\|_\infty - \langle \mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)} \rangle + 2\|\mathbf{u}_i^{(t)} - \mathbf{u}_i(\mathbf{x}_{-i}^{(t)})\|_\infty \quad (12) \\ &\leq \epsilon \left(1 + \frac{4n\sqrt{m}}{\delta}\right) \leq \epsilon \left(\frac{5n\sqrt{m}}{\delta}\right). \end{aligned}$$

where (12) follows from the fact that $\langle \mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)} - \mathbf{u}_i(\mathbf{x}_{-i}^{(t)}) \rangle \leq \|\mathbf{x}_i^{(t)}\|_1 \|\mathbf{u}_i^{(t)} - \mathbf{u}_i(\mathbf{x}_{-i}^{(t)})\|_\infty \leq \|\mathbf{u}_i^{(t)} - \mathbf{u}_i(\mathbf{x}_{-i}^{(t)})\|_\infty$ since $\|\mathbf{x}_i^{(t)}\|_1 = 1$. We conclude that $(\mathbf{x}_1^{(t)}, \dots, \mathbf{x}_n^{(t)})$ is an $\epsilon(5n\sqrt{m}/\delta)$ -Nash equilibrium; rescaling ϵ leads to the claim. \square

Remark C.9 (Convergence to CCE under alternating updates). The folk connection linking no-regret learning and coarse correlated equilibria (Proposition B.2) is predicated on the dynamics being executed in a simultaneous fashion. We observe that, in potential games, even alternating RM^+ converges to the set of CCEs in the following sense. Lemma C.6 together with (9) imply that $\sum_{t=1}^T \sum_{i=1}^n \|\mathbf{x}_i^{(t+1)} - \mathbf{x}_i^{(t)}\|_1^2 = O_T(\sqrt{T})$. The Cauchy-Schwarz inequality in turn yields

$$\sum_{t=1}^T \sum_{i=1}^n \|\mathbf{x}_i^{(t+1)} - \mathbf{x}_i^{(t)}\|_1 \leq \sqrt{\sum_{t=1}^T \left(\sum_{i=1}^n \|\mathbf{x}_i^{(t+1)} - \mathbf{x}_i^{(t)}\|_1 \right)^2 \sum_{t=1}^T 1^2} = O_T(T^{3/4}). \quad (13)$$

Now, by the fact that RM^+ has the no-regret property (Proposition 2.3), $\max_{\mathbf{x}'_i \in \Delta(\mathcal{A}_i)} \sum_{t=1}^T \langle \mathbf{x}'_i - \mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)} \rangle = O_T(\sqrt{T})$. Further, by Claim C.8 and (13), $\sum_{t=1}^T \|\mathbf{u}_i(\mathbf{x}_{-i}^{(t)}) - \mathbf{u}_i^{(t)}\|_\infty \leq \sum_{t=1}^T \sum_{i' \neq i} \|\mathbf{x}_{i'}^{(t+1)} - \mathbf{x}_{i'}^{(t)}\|_1 = O_T(T^{3/4})$. As a result, we conclude that, for any player $i \in [n]$,

$$\max_{\mathbf{x}'_i \in \Delta(\mathcal{A}_i)} \sum_{t=1}^T \langle \mathbf{x}'_i - \mathbf{x}_i^{(t)}, \mathbf{u}_i(\mathbf{x}_{-i}^{(t)}) \rangle = O_T(T^{3/4}).$$

This implies that the correlated distribution $\frac{1}{T} \sum_{t=1}^T \otimes_{i=1}^n \mathbf{x}_i^{(t)}$ is an ϵ -CCE for some $\epsilon = O_T(T^{-1/4})$ (Proposition B.2); here, $\otimes_{i=1}^n \mathbf{x}_i^{(t)}$ denotes the product distribution induced by $(\mathbf{x}_1^{(t)}, \dots, \mathbf{x}_n^{(t)})$.

Moving on, a further implication of Lemma C.4 is that having a large regret vector can slow down RM^+ . Employing discounting, in the form of DRM^+ (Algorithm 3), addresses this deficiency in potential games, as we prove below. It is worth stressing that while discounting is sensible when employing alternating RM^+ in potential games, this is not the case in more general constrained optimization problems; there, having a regret vector with a *small* norm can impede convergence (cf. Lemma 3.7).

Corollary 3.5. *In any potential game, ϵ -lazy alternating DRM^+ with discount factor $1 - \gamma \in (0, 1)$ requires at most $1 + \frac{m\Phi_{\text{range}}}{\epsilon^2\sqrt{\gamma}}$ rounds to converge to an ϵ -Nash equilibrium.*

Proof. Lemma 3.3 can be directly adjusted for DRM^+ : the updated strategy $\mathbf{x}' = \mathbf{r}'/\|\mathbf{r}'\|_1$, where $\mathbf{r}' = \alpha[\mathbf{r} + \mathbf{u} - \langle \mathbf{x}, \mathbf{u} \rangle \mathbf{1}]^+ \neq \mathbf{0}$, remains the same since we only rescaled the regret vector by α . That is, for any $\mathbf{u} \in \mathbb{R}^{\mathcal{A}}$,

$$\langle \mathbf{x}' - \mathbf{x}, \mathbf{u} \rangle \geq \frac{1}{\|\mathbf{r}'\|_1} \left(\max_{a \in \mathcal{A}} \mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle \right)^2.$$

The key difference compared to Theorem 3.4 is that we now maintain the invariance $\|\mathbf{r}_i^{(t)}\|_2 \leq \sqrt{m/\gamma}$ for all $i \in [n]$ and $t \in [T]$ (Corollary C.3). Consequently,

$$\Phi_{\text{range}} \geq \sum_{t=1}^T \sum_{i=1}^n \frac{1}{\|\mathbf{r}_i^{(t)}\|_1} \text{BRGap}_i(\mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)})^2 \mathbb{1}\{\text{BRGap}_i(\mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)}) > \epsilon\}.$$

If in every round $t \in [T]$ there is a player $i \in [n]$ such that $\text{BRGap}_i(\mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)}) > \epsilon$, we have $\Phi_{\text{range}} \geq \sum_{t=1}^T (\sqrt{\gamma}/m)\epsilon^2$, so $T \leq \frac{m\Phi_{\text{range}}}{\epsilon^2\sqrt{\gamma}}$. This concludes the proof. \square

Unlike RM^+ and DRM^+ , RM only has a *conditional* one-step improvement because the regret vector can have negative coordinates.

Lemma 3.6. For any $\mathbf{r} \in \mathbb{R}_{\geq 0}^A$ and $\mathbf{u} \in \mathbb{R}^A$, we define $\mathbf{x} := \boldsymbol{\theta}/\|\boldsymbol{\theta}\|_1$, where $\boldsymbol{\theta} := \max(\mathbf{r}, \mathbf{0})$; if $\boldsymbol{\theta} = \mathbf{0}$, $\mathbf{x} \in \Delta(\mathcal{A})$ can be arbitrary. If $\mathbf{r}' := \mathbf{r} + \mathbf{u} - \langle \mathbf{x}, \mathbf{u} \rangle \mathbf{1}$ and $\mathbf{x}' := \boldsymbol{\theta}'/\|\boldsymbol{\theta}'\|_1$, where $\boldsymbol{\theta}' = \max(\mathbf{r}', \mathbf{0}) \neq \mathbf{0}$, then $\langle \mathbf{x}' - \mathbf{x}, \mathbf{u} \rangle \geq \frac{1}{\|\boldsymbol{\theta}'\|_1} \|\boldsymbol{\theta}' - \boldsymbol{\theta}\|_2^2 \geq \frac{1}{\|\boldsymbol{\theta}'\|_1} (\max_{a \in \mathcal{A}} \mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle)^2 \mathbb{1}\{\mathbf{r}[a] \geq 0\}$, where $a \in \arg \max_{a' \in \mathcal{A}} \mathbf{u}[a']$. If $\boldsymbol{\theta}' = \mathbf{0}$, then $\langle \mathbf{x}, \mathbf{u} \rangle = \langle \mathbf{x}', \mathbf{u} \rangle \geq \mathbf{u}[a]$.

Proof. We define $\boldsymbol{\delta} := \boldsymbol{\theta}' - \boldsymbol{\theta}$. Following the proof of Lemma 3.3, it suffices to show that

$$\sum_{a' \in \mathcal{A}} \boldsymbol{\theta}[a'] \sum_{a \in \mathcal{A}} \boldsymbol{\delta}[a] (\mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle) \geq \sum_{a' \in \mathcal{A}} \boldsymbol{\theta}[a'] \left(\max_{a \in \mathcal{A}} \mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle \right)^2 \mathbb{1}\{\mathbf{r}[a] \geq 0\}. \quad (14)$$

For an action $a \in \mathcal{A}$, we consider the following cases.

- If $\mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle \geq 0$,
 - if $\mathbf{r}[a] \geq 0$, we have $\boldsymbol{\delta}[a] = \mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle$.
 - If $\mathbf{r}[a] < 0$, it follows that $\boldsymbol{\delta}[a] \geq 0$; in particular, $\boldsymbol{\delta}[a] = 0$ if $\mathbf{r}'[a] \leq 0$ and $\boldsymbol{\delta}[a] > 0$ otherwise. As a result, $\boldsymbol{\delta}[a](\mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle) \geq 0$.
- If $\mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle < 0$,
 - if $\mathbf{r}[a] \leq 0$, we have $\boldsymbol{\delta}[a] = 0$ since $\boldsymbol{\theta}[a] = 0 = \boldsymbol{\theta}'[a]$.
 - if $\mathbf{r}[a] > 0$, it follows that $\boldsymbol{\delta}[a] < 0$. Again, we have $\boldsymbol{\delta}[a](\mathbf{u}[a] - \langle \mathbf{x}, \mathbf{u} \rangle) \geq 0$.

Combining those items, (14) follows. \square

We now turn to the more general setting of constrained optimization. First, combining Lemmas C.5 and C.6 that were proven earlier, we formally show that RM^+ improves the value of the underlying function provided that the norm of the regret vector is not too small.

Lemma 3.7. Let u be an L -smooth function over $\Delta(\mathcal{A})$. For any $\mathbf{r} \in \mathbb{R}_{\geq 0}^A$ with $\mathbf{r} \neq \mathbf{0}$, we define $\mathbf{x} := \mathbf{r}/\|\mathbf{r}\|_1$. Further, let $\mathbf{r}' := [\mathbf{r} + \nabla u(\mathbf{x}) - \langle \mathbf{x}, \nabla u(\mathbf{x}) \rangle \mathbf{1}]^+ \neq \mathbf{0}$ and $\mathbf{x}' := \mathbf{r}'/\|\mathbf{r}'\|_1$. If $\|\mathbf{r}'\|_2 \geq \max\{2m, 9mL\}$, then

$$u(\mathbf{x}') - u(\mathbf{x}) \geq \frac{1}{2\|\mathbf{r}'\|_1} \left(\max_{\mathbf{x}^* \in \Delta(\mathcal{A})} \langle \mathbf{x}^* - \mathbf{x}, \nabla u(\mathbf{x}) \rangle \right)^2.$$

Proof. Using the quadratic bound for u , we have

$$\begin{aligned} u(\mathbf{x}') - u(\mathbf{x}) &\geq \langle \nabla u(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle - \frac{L}{2} \|\mathbf{x} - \mathbf{x}'\|_2^2 \\ &\geq \frac{1}{\|\mathbf{r}'\|_1} \|\mathbf{r} - \mathbf{r}'\|_2^2 - \frac{L}{2} \|\mathbf{r} - \mathbf{r}'\|_1^2 \left(\frac{1}{\|\mathbf{r}\|_1} + \frac{1}{\|\mathbf{r}'\|_1} \right)^2 \end{aligned} \quad (15)$$

$$\geq \frac{1}{\|\mathbf{r}'\|_1} \|\mathbf{r} - \mathbf{r}'\|_2^2 - \frac{9mL}{2\|\mathbf{r}'\|_1^2} \|\mathbf{r} - \mathbf{r}'\|_2^2 \quad (16)$$

$$\geq \frac{1}{2\|\mathbf{r}'\|_1} \|\mathbf{r} - \mathbf{r}'\|_2^2, \quad (17)$$

where (15) uses the one-step improvement property (Lemma C.5) applied for $\mathbf{u} := \nabla u(\mathbf{x})$ together with Lemma C.6; (16) follows from the fact that $\|\mathbf{r}\|_1 \geq \|\mathbf{r}'\|_1 - m \geq \frac{1}{2}\|\mathbf{r}'\|_1$ since $\|\mathbf{r}'\|_1 \geq \|\mathbf{r}'\|_2 \geq 2m$ and $|\langle \mathbf{x} - \mathbf{x}', \nabla u(\mathbf{x}) \rangle| \leq 1$ for all $\mathbf{x}' \in \Delta(\mathcal{A})$ (per our normalization assumption); and (17) follows from the assumption that $\|\mathbf{r}'\|_1 \geq 9mL$. \square

To make use of [Lemma 3.7](#), we next establish that the ℓ_2 norm of the regret vector under RM^+ is nondecreasing.

Lemma 3.8. *For any t , RM^+ guarantees $\|\mathbf{r}^{(t)}\|_2^2 \geq \|\mathbf{r}^{(t-1)}\|_2^2 + \|[\mathbf{g}^{(t)}]^+\|_2^2$, where $\mathbf{g}^{(t)} := \nabla u(\mathbf{x}^{(t)}) - \langle \nabla u(\mathbf{x}^{(t)}), \mathbf{x}^{(t)} \rangle \mathbf{1}$ is the instantaneous regret vector at round t .*

Proof. We have $\mathbf{r}^{(t)} - \mathbf{r}^{(t-1)} = \max(\mathbf{g}^{(t)}, -\mathbf{r}^{(t-1)})$ (element-wise), so

$$\|\mathbf{r}^{(t)} - \mathbf{r}^{(t-1)}\|_2 = \|\max(\mathbf{g}^{(t)}, -\mathbf{r}^{(t-1)})\|_2 \geq \|[\mathbf{g}^{(t)}]^+\|_2.$$

Further, $\langle \mathbf{r}^{(t-1)}, \mathbf{r}^{(t)} - \mathbf{r}^{(t-1)} \rangle = \langle \mathbf{r}^{(t-1)}, \max(\mathbf{g}^{(t)}, -\mathbf{r}^{(t-1)}) \rangle \geq \langle \mathbf{r}^{(t-1)}, \mathbf{g}^{(t)} \rangle = 0$, where we used the fact that $\mathbf{r}^t \geq \mathbf{0}$, coordinate-wise. Therefore,

$$\begin{aligned} \|\mathbf{r}^{(t)}\|_2^2 &= \|\mathbf{r}^{(t)} - \mathbf{r}^{(t-1)} + \mathbf{r}^{(t-1)}\|_2^2 = \|\mathbf{r}^{(t)} - \mathbf{r}^{(t-1)}\|_2^2 + \|\mathbf{r}^{(t-1)}\|_2^2 + 2\langle \mathbf{r}^{(t-1)}, \mathbf{r}^{(t)} - \mathbf{r}^{(t-1)} \rangle \\ &\geq \|\mathbf{r}^{(t-1)}\|_2^2 + \|[\mathbf{g}^{(t)}]^+\|_2^2, \end{aligned}$$

as claimed. \square

Armed with [Lemmas 3.7](#) and [3.8](#), we can now prove [Theorem 3.9](#).

Theorem 3.9 (Single simplex). *Let u be an L -smooth function in $\Delta(\mathcal{A}) \subset \mathbb{R}^m$ with range u_{range} and $R := \max\{2m, 9mL\}$. RM^+ requires at most $1 + \frac{(m(2u_{\text{range}} + R^2))^2}{\epsilon^4}$ rounds to reach an ϵ -KKT point.*

Proof. Let $t_c \in [T]$ be the largest t such that $\|\mathbf{r}^{(t)}\|_2 < \max\{2m, 9mL\} = R$. By [Lemma 3.8](#),

$$\|\mathbf{r}^{(t)}\|_2^2 \geq \|\mathbf{r}^{(t-1)}\|_2^2 + \|[\mathbf{g}^{(t)}]^+\|_2^2 \geq \|\mathbf{r}^{(t-1)}\|_2^2 + \text{KKTGap}(\mathbf{x}^{(t)})^2;$$

so,

$$\sum_{t=1}^{t_c} \text{KKTGap}(\mathbf{x}^{(t)})^2 \leq \sum_{t=1}^{t_c} (\|\mathbf{r}^{(t)}\|_2^2 - \|\mathbf{r}^{(t-1)}\|_2^2) = \|\mathbf{r}^{(t_c)}\|_2^2 \leq R^2. \quad (18)$$

Further, for any $t \geq t_c + 1$, we have $\|\mathbf{r}^{(t)}\|_2 \geq R$ since the ℓ_2 norm of the regret vector is nondecreasing ([Lemma 3.8](#)) and $\|\mathbf{r}^{(t_c+1)}\|_2 \geq R$. Thus, by [Lemma 3.7](#),

$$\sum_{t=t_c+1}^T \frac{1}{2\|\mathbf{r}^{(t)}\|_1} \text{KKTGap}(\mathbf{x}^{(t)})^2 \leq u(\mathbf{x}^{(T+1)}) - u(\mathbf{x}^{(t_c+1)}) \leq u_{\text{range}}. \quad (19)$$

Combining (18) and (19), together with the fact that $\|\mathbf{r}^{(t)}\|_1 \leq m\sqrt{t}$,

$$\sum_{t=1}^T \frac{1}{m\sqrt{t}} \text{KKTGap}(\mathbf{x}^{(t)})^2 \leq 2u_{\text{range}} + R^2.$$

This leads to the claim. \square

Next, we use [Theorem 3.9](#) to establish that even *simultaneous* RM^+ converges in symmetric potential games, as long as the players have the same initialization.

Corollary 3.10. *In any symmetric potential game, simultaneous RM^+ converges to an ϵ -Nash equilibrium after $O_\epsilon(1/\epsilon^4)$ rounds. In particular, if convergence to the set of CCE happens at a rate of $T^{-(1-\alpha)}$, for some $\alpha \in [0, 1/2]$, the rate of convergence to Nash equilibria is no worse than $T^{-\frac{1-\alpha}{2}}$.*

Proof. We will argue that simultaneous RM^+ , under the same initialization, in a symmetric game is tantamount to running RM^+ with respect to the function $\Delta(\mathcal{A}) \ni \mathbf{x} \mapsto \Phi(\mathbf{x}, \dots, \mathbf{x})$, where $\mathcal{A}_1 = \dots = \mathcal{A}_n = \mathcal{A}$. We first claim that for any $a \in \mathcal{A}$,

$$\frac{\partial \Phi(\mathbf{x}, \dots, \mathbf{x})}{\partial \mathbf{x}[a]} = \sum_{i=1}^n \frac{\partial \Phi(\mathbf{x}_1, \dots, \mathbf{x}_n)}{\partial \mathbf{x}_i[a]} \Big|_{\mathbf{x}_1 = \dots = \mathbf{x}_n = \mathbf{x}}. \quad (20)$$

By definition, $\Phi(\mathbf{x}_1, \dots, \mathbf{x}_n)$ is multilinear. Let us consider a monomial of Φ of the form

$$m_{\mathbf{a}}(\mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_{i=1}^n \mathbf{x}_i[a_i]$$

for some joint action $\mathbf{a} = (a_1, \dots, a_n) \in \mathcal{A}_1 \times \dots \times \mathcal{A}_n$. Then

$$\begin{aligned} \frac{\partial m_{\mathbf{a}}(\mathbf{x}, \dots, \mathbf{x})}{\partial \mathbf{x}[a]} &= \frac{\partial}{\partial \mathbf{x}[a]} \left(\prod_{a' \in \mathcal{A}'} (\mathbf{x}[a'])^{d(a')} \right) \\ &= \begin{cases} d(a) (\mathbf{x}[a])^{d(a)-1} \prod_{a' \in \mathcal{A}' \setminus \{a\}} (\mathbf{x}[a'])^{d(a')} & \text{if } a \in \mathcal{A}' \text{ and} \\ 0 & \text{otherwise,} \end{cases} \end{aligned} \quad (21)$$

where $\mathcal{A}' = \mathcal{A}'(\mathbf{a}) = \{a' \in \mathcal{A} : \exists i \in [n] \text{ such that } a' = a_i\}$ and degrees $d(a') = |\{i \in [n] : a_i = a'\}| \geq 1$ for $a' \in \mathcal{A}'$. Further,

$$\frac{\partial m_{\mathbf{a}}(\mathbf{x}_1, \dots, \mathbf{x}_n)}{\partial \mathbf{x}_i[a]} = \begin{cases} \prod_{i' \neq i} \mathbf{x}_{i'}[a_{i'}] & \text{if } a = a_i \\ 0 & \text{otherwise.} \end{cases}$$

In particular,

$$\left. \frac{\partial m_{\mathbf{a}}(\mathbf{x}_1, \dots, \mathbf{x}_n)}{\partial \mathbf{x}_i[a]} \right|_{\mathbf{x}_1 = \dots = \mathbf{x}_n} = \mathbb{1}\{a = a_i\} \mathbf{x}[a]^{d(a)-1} \prod_{a' \in \mathcal{A}' \setminus \{a\}} (\mathbf{x}[a'])^{d(a')}.$$

As a result,

$$\sum_{i=1}^n \left. \frac{\partial m_{\mathbf{a}}(\mathbf{x}_1, \dots, \mathbf{x}_n)}{\partial \mathbf{x}_i[a]} \right|_{\mathbf{x}_1 = \dots = \mathbf{x}_n} = \mathbb{1}\{a \in \mathcal{A}'\} d(a) \mathbf{x}[a]^{d(a)-1} \prod_{a' \in \mathcal{A}' \setminus \{a\}} (\mathbf{x}[a'])^{d(a')},$$

which matches the expression in (21). That is, we have shown that for any $\mathbf{a} \in \mathcal{A}_1 \times \dots \times \mathcal{A}_n$,

$$\frac{\partial m_{\mathbf{a}}(\mathbf{x}, \dots, \mathbf{x})}{\partial \mathbf{x}[a]} = \sum_{i=1}^n \left. \frac{\partial m_{\mathbf{a}}(\mathbf{x}_1, \dots, \mathbf{x}_n)}{\partial \mathbf{x}_i[a]} \right|_{\mathbf{x}_1 = \dots = \mathbf{x}_n}.$$

Since $\Phi(\mathbf{x}_1, \dots, \mathbf{x}_n) = \sum_{\mathbf{a} \in \mathcal{A}_1 \times \dots \times \mathcal{A}_n} \Phi(\mathbf{a}) m_{\mathbf{a}}(\mathbf{x}_1, \dots, \mathbf{x}_n)$, (20) follows by linearity. Moreover, the symmetry assumption concerning the potential Φ tells us that

$$\left. \frac{\partial \Phi(\mathbf{x}_1, \dots, \mathbf{x}_n)}{\partial \mathbf{x}_i[a]} \right|_{\mathbf{x}_1 = \dots = \mathbf{x}_n} = \left. \frac{\partial \Phi(\mathbf{x}_1, \dots, \mathbf{x}_n)}{\partial \mathbf{x}_{i'}[a]} \right|_{\mathbf{x}_1 = \dots = \mathbf{x}_n}$$

for any $i, i' \in [n]$. Combining with (20), we have that for any $a \in \mathcal{A}$,

$$\frac{\partial \Phi(\mathbf{x}, \dots, \mathbf{x})}{\partial \mathbf{x}[a]} = n \left. \frac{\partial \Phi(\mathbf{x}_1, \dots, \mathbf{x}_n)}{\partial \mathbf{x}_i[a]} \right|_{\mathbf{x}_1 = \dots = \mathbf{x}_n = \mathbf{x}} \quad \forall i \in [n].$$

Given that RM^+ is scale invariant, we conclude that simultaneous RM^+ on the potential game is equivalent to running RM^+ on the function $\Delta(\mathcal{A}) \ni \mathbf{x} \mapsto \Phi(\mathbf{x}, \dots, \mathbf{x})$, and the claim follows from [Theorem 3.9](#). \square

We now turn to the analysis of alternating RM^+ in constrained optimization problems over multiple probability simplices. We clarify that when the regret vector is initialized at $\mathbb{R}_{\geq 0}^{\mathcal{A}_i} \ni \mathbf{r}_i^{(0)} \neq \mathbf{0}$, as assumed below, the first strategy must be defined consistently as $\mathbf{x}_i^{(1)} \propto \mathbf{r}_i^{(0)}$.

Corollary 3.11. *If u is an L -smooth function in $\Delta(\mathcal{A}_1) \times \dots \times \Delta(\mathcal{A}_n)$ with range u_{range} , ϵ -lazy alternating RM^+ initialized at $\mathbf{r}_i^{(0)} = \max\{2\sqrt{m_i}, 9\sqrt{m_i}L\} \mathbf{1}$ for each player $i \in [n]$ requires at most $1 + \frac{4n^4 m^2 u_{\text{range}}^2}{\epsilon^4}$ rounds to reach an ϵ -KKT point of u .*

Proof. By Lemma 3.8, it follows that $\|\mathbf{r}_i^{(t)}\|_2 \geq \|\mathbf{r}_i^{(0)}\|_2 = \max\{2m_i, 9m_iL\}$. Following the proof of Lemma 3.7, we have

$$u(\mathbf{x}_{i' \leq i}^{(t+1)}, \mathbf{x}_{i' > i}^{(t)}) - u(\mathbf{x}_{i' < i}^{(t+1)}, \mathbf{x}_{i' \geq i}^{(t)}) \geq \frac{1}{2\|\mathbf{r}_i^{(t)}\|_1} \left(\text{BRGap}_i(\mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)}) \right)^2 \mathbb{1}\{\text{BRGap}_i(\mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)}) > \epsilon\}.$$

As a result, the telescopic summation yields

$$\sum_{t=1}^T \sum_{i=1}^n \frac{1}{2\|\mathbf{r}_i^{(t)}\|_1} \left(\text{BRGap}_i(\mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)}) \right)^2 \mathbb{1}\{\text{BRGap}_i(\mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)}) > \epsilon\} \leq u_{\text{range}},$$

Since $\|\mathbf{r}_i^{(t)}\|_1 \leq m\sqrt{t}$ for all $i \in [n]$ and $t \in [T]$, it will take at most $1 + 4m^2u_{\text{range}}^2/\epsilon^4$ rounds to converge to a point in which all players have at most an ϵ best-response gap, which in turn implies that the KKT gap is at most $n\epsilon$. Rescaling ϵ concludes the proof. \square

To analyze simultaneous RM^+ , we adapt Lemma 3.7 as follows.

Lemma C.10. *Let u be an L -smooth function over $\Delta(\mathcal{A}_1) \times \dots \times \Delta(\mathcal{A}_n)$. For any $i \in [n]$ and $\mathbf{r}_i \in \mathbb{R}_{\geq 0}^{\mathcal{A}_i}$ with $\mathbf{r}_i \neq \mathbf{0}$, we define $\mathbf{x}_i := \mathbf{r}_i / \|\mathbf{r}_i\|_1$. Further, let $\mathbf{r}'_i := [\mathbf{r}_i + \nabla_{\mathbf{x}_i} u(\mathbf{x}) - \langle \mathbf{x}_i, \nabla_{\mathbf{x}_i} u(\mathbf{x}) \rangle \mathbf{1}]^+ \neq \mathbf{0}$ and $\mathbf{x}'_i := \mathbf{r}'_i / \|\mathbf{r}'_i\|_1$. If $\|\mathbf{r}'_i\|_2 \geq \max\{2m_i, 9m_iL\}$ for any $i \in [n]$, then*

$$u(\mathbf{x}') - u(\mathbf{x}) \geq \frac{1}{2 \max_{1 \leq i \leq n} \|\mathbf{r}'_i\|_1} \sum_{i=1}^n \left(\max_{\mathbf{x}'_i \in \Delta(\mathcal{A}_i)} \langle \mathbf{x}'_i - \mathbf{x}_i, \nabla_{\mathbf{x}_i} u(\mathbf{x}) \rangle \right)^2.$$

Proof. Using the quadratic bound for u , we have

$$\begin{aligned} u(\mathbf{x}') - u(\mathbf{x}) &\geq \langle \nabla u(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle - \frac{L}{2} \|\mathbf{x} - \mathbf{x}'\|_2^2 \\ &= \sum_{i=1}^n \left(\langle \nabla_{\mathbf{x}_i} u(\mathbf{x}), \mathbf{x}'_i - \mathbf{x}_i \rangle - \frac{L}{2} \|\mathbf{x}_i - \mathbf{x}'_i\|_2^2 \right) \\ &\geq \sum_{i=1}^n \left(\frac{1}{\|\mathbf{r}'_i\|_1} \|\mathbf{r}_i - \mathbf{r}'_i\|_2^2 - \frac{L}{2} \|\mathbf{r}_i - \mathbf{r}'_i\|_1^2 \left(\frac{1}{\|\mathbf{r}_i\|_1} + \frac{1}{\|\mathbf{r}'_i\|_1} \right)^2 \right) \end{aligned} \quad (22)$$

$$\geq \sum_{i=1}^n \left(\frac{1}{\|\mathbf{r}'_i\|_1} \|\mathbf{r}_i - \mathbf{r}'_i\|_2^2 - \frac{9m_iL}{2\|\mathbf{r}'_i\|_1^2} \|\mathbf{r}_i - \mathbf{r}'_i\|_2^2 \right) \quad (23)$$

$$\geq \frac{1}{2} \sum_{i=1}^n \left(\frac{1}{\|\mathbf{r}'_i\|_1} \|\mathbf{r}_i - \mathbf{r}'_i\|_2^2 \right), \quad (24)$$

where (22) uses the one-step improvement property (Lemma C.5) applied for each player $i \in [n]$ and $\mathbf{u}_i := \nabla_{\mathbf{x}_i} u(\mathbf{x})$ together with Lemma C.6; (23) follows from the fact that $\|\mathbf{r}_i\|_1 \geq \|\mathbf{r}'_i\|_1 - m_i \geq \frac{1}{2}\|\mathbf{r}'_i\|_1$ since $\|\mathbf{r}'_i\|_1 \geq \|\mathbf{r}'_i\|_2 \geq 2m_i$ and $|\langle \mathbf{x}_i - \mathbf{x}'_i, \nabla_{\mathbf{x}_i} u(\mathbf{x}) \rangle| \leq 1$ for all $\mathbf{x}'_i \in \Delta(\mathcal{A}_i)$ (per our normalization assumption); and (24) follows from the assumption that $\|\mathbf{r}'_i\|_1 \geq 9m_iL$ for each $i \in [n]$. \square

Similarly to Corollary 3.11, we can now show the following concerning simultaneous RM^+ .

Corollary C.11. *If u is an L -smooth function in $\Delta(\mathcal{A}_1) \times \dots \times \Delta(\mathcal{A}_n)$ with range u_{range} , simultaneous RM^+ initialized at $\mathbf{r}_i^{(0)} = \max\{2\sqrt{m_i}, 9\sqrt{m_i}L\}\mathbf{1}$ for each player $i \in [n]$ requires at most $1 + \frac{4n^2m^2u_{\text{range}}^2}{\epsilon^4}$ rounds to reach an ϵ -KKT point of u .*

Proof. By Lemma 3.8, it follows that $\|\mathbf{r}_i^{(t)}\|_2 \geq \|\mathbf{r}_i^{(0)}\|_2 = \max\{2m_i, 9m_iL\}$. By Lemma C.10,

$$u(\mathbf{x}^{(t+1)}) - u(\mathbf{x}^{(t)}) \geq \frac{1}{2n} \frac{1}{\max_{1 \leq i \leq n} \|\mathbf{r}_i^{(t)}\|_1} \left(\text{KKTGap}(\mathbf{x}^{(t)}) \right)^2,$$

and the claim follows. \square

Next, we show how to extend the analysis even when the regret vector is initialized at zero, at the cost of incurring a worse dependence on $1/\epsilon$.

Theorem 3.12. *If u is an L -smooth function in $\Delta(\mathcal{A}_1) \times \dots \times \Delta(\mathcal{A}_n)$ with range u_{range} , ϵ -lazy alternating (or simultaneous) RM^+ requires at most $O_\epsilon(1/\epsilon^8)$ rounds to reach an ϵ -KKT point of u .*

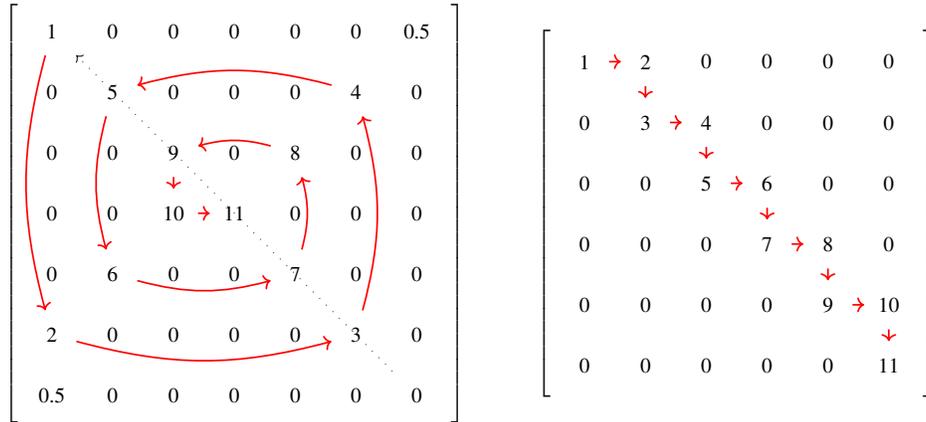
Proof. We analyze ϵ -lazy simultaneous RM^+ ; the alternating version can be treated similarly. If $\|\mathbf{r}_i^{(t)}\|_2 \geq \max\{2m_i, 9m_iL\}$ for all players with best-response gap at least ϵ , we have that $u(\mathbf{x}^{(t+1)}) - u(\mathbf{x}^{(t)})$ is lower bounded by

$$\frac{1}{2 \max_i \|\mathbf{r}_i^{(t)}\|_1} \sum_{i=1}^n \left(\text{BRGap}_i(\mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)}) \right)^2 \mathbb{1}\{\text{BRGap}_i(\mathbf{x}_i^{(t)}, \mathbf{u}_i^{(t)}) > \epsilon\}; \quad (25)$$

this follows similarly to [Lemma C.10](#) since players with best-response gap below ϵ do not update their strategies (under lazy updates). Now, by [Lemma 3.8](#), it follows that the total number of rounds in which there is a player i with at least an ϵ best-response gap and $\|\mathbf{r}_i^{(t)}\|_2 < \max\{2m_i, 9m_iL\} = R_i$ is at most $\sum_{i=1}^n (1 + R_i^2/\epsilon^2)$. We will now bound the number of rounds T' it takes to encounter two such rounds. We observe that, between such rounds, we only update players that have at least an ϵ best-response gap and $\|\mathbf{r}_i^{(t)}\|_2 \geq R_i$. By (25), this can continue for at most $1 + 2u_{\text{range}}m\sqrt{T}/\epsilon^2$ rounds, where T is the total number of rounds it takes for all players to have a best-response gap of at most ϵ . That is, $T' \leq 1 + 2u_{\text{range}}m\sqrt{T}/\epsilon^2$. Further, $T \leq T' + T' \sum_{i=1}^n (1 + R_i^2/\epsilon^2)$, and the claim follows by solving in terms of T and rescaling ϵ . \square

C.2 PROOFS FROM SECTION 4

We conclude with the proofs from [Section 4](#). Below (left figure), we provide an illustrative example of the matrix \mathbf{B} , defined earlier in (6), for $m = 6$. The plots in [Figure 1](#) are obtained by running simultaneous RM (left) and alternating RM^+ (right) on this exact game. It is worth pointing out that one could also use a staircase—instead of a spiral—pattern as shown in the right figure.



Our main goal is to prove the following invariance.

Property 4.1. *After the first round both players play the first action. Thereupon, either the players play with probability 1 ($a_1(k), a_2(k)$), or, when k is odd, only Player 1 (respectively, Player 2 when k is even) mixes between $a_1(k)$ and $a_1(k+1)$ (respectively, $a_2(k)$ and $a_2(k+1)$). If a row or a column stops being played, it will never be played henceforth. An action profile $(a_1(k+1), a_2(k+1))$ is played with positive probability only if $(a_1(k), a_2(k))$ was played at some previous round.*

It is possible to check the claim for $k = 1, 2, 3, 4$ by executing RM for a certain number of rounds. In particular, we find that $T_3 \geq 5$ and $T_4 \geq 20$. We proceed by induction in k . Suppose that it holds for all payoffs $1, \dots, \kappa$. We will show that it holds for $\kappa + 1$.

Lemma C.12. *For any even $\kappa + 2 \geq k \geq 4$, let $\mathbf{r}_1^{(\overline{t_{k-2}})}[a_1]$ be the regret of Player 1 with respect to any action $a_1 \in \mathcal{A}_1(k)$. Then $\mathbf{r}_1^{(\overline{t_{k-2}})}[a_1] \leq -\sum_{l=2}^{k-2} (l-1)T_l$. Similarly, for any odd $\kappa + 2 \geq$*

1296 $k \geq 5$, if $\mathbf{r}_2^{(\overline{t_{k-2}})}[a_2]$ is the regret of Player 2 with respect to any action $a_2 \in \mathcal{A}_2(k)$, $\mathbf{r}_2^{(\overline{t_{k-2}})}[a_2] \leq$
 1297 $-\sum_{l=2}^{k-2} (l-1)T_l$.
 1298

1299 *Proof.* Let $a_1 \in \mathcal{A}_1(k)$ and $l \in [k-2]$ for an even k . Playing $a_1 \in \mathcal{A}_1(k)$ during $[t_l, \bar{t}_l]$ gives Player
 1300 1 a utility of 0; this follows from the fact that for any column $a_2 \in \{a_2(1), a_2(3), \dots, a_2(k-3)\} =$
 1301 $\{a_2(1), a_2(2), a_2(3), \dots, a_2(k-3), a_2(k-2)\}$, it holds that $\mathbf{A}[a_1(k), a_2] = 0$, by construction of
 1302 \mathbf{A} . At the same time, Player 1 actually got a utility of at least $l-1$ for each round in $[t_l, \bar{t}_l]$. This
 1303 means that every time Player 1 updates its regret vector within the time period $[t_l, \bar{t}_l]$, the regret of
 1304 a_1 decreases by at least $l-1$. The same reasoning applies for Player 2 when k is odd. \square
 1305

1306 **Lemma C.13.** For any even $\kappa \geq k \geq 4$, $T_k \geq -\frac{1}{2}\mathbf{r}_2^{(\overline{t_{k-1}})}[a_2(k+1)]$. Similarly, for every odd
 1307 $\kappa \geq k \geq 5$, $T_k \geq -\frac{1}{2}\mathbf{r}_1^{(\overline{t_{k-1}})}[a_1(k+1)]$.
 1308

1309 *Proof.* We fix an even $k \geq 4$. T_k is at least as large as the number of rounds it takes for $a_2(k+1)$
 1310 to have nonnegative regret, starting from $\mathbf{r}_2^{(\overline{t_{k-1}})}[a_2(k+1)]$. But in every round in $[t_k, \bar{t}_k]$ the regret
 1311 of $a_2(k+1)$ can increase additively by at most 2. This holds because the utility of Player 2 for
 1312 playing $a_2(k+1)$ is larger than the utility obtained in each round in $[t_k, \bar{t}_k]$ by at most 2. So,
 1313 $T_k \geq \lceil -\frac{1}{2}\mathbf{r}_2^{(\overline{t_{k-1}})}[a_2(k+1)] \rceil \geq -\frac{1}{2}\mathbf{r}_2^{(\overline{t_{k-1}})}[a_2(k+1)]$. The same reasoning applies when k is
 1314 odd. \square
 1315

1316 The following upper bound on the regret is crude, but will suffice for our purposes.
 1317

1318 **Lemma C.14** (Regret upper bound). For any even $\kappa \geq k \geq 4$, $\|\mathbf{r}_1^{(\overline{t_k})}\|^+ \leq 2\|\mathbf{r}_1^{(\overline{t_{k-2}})}\|^+ + \infty +$
 1319 $2 \leq \frac{5}{3}2^{k/2}$ since $\|\mathbf{r}_1^{(\overline{t_2})}\|_\infty \leq \frac{4}{3}$. Similarly, for any odd $k \geq 5$, $\|\mathbf{r}_2^{(\overline{t_k})}\|^+ \leq 2\|\mathbf{r}_2^{(\overline{t_{k-2}})}\|^+ + \infty + 2 \leq$
 1320 $\frac{5}{3}2^{(k-1)/2}$ since $\|\mathbf{r}_2^{(\overline{t_3})}\|_\infty \leq \frac{4}{3}$.
 1321

1322 *Proof.* The fact that $\|\mathbf{r}_1^{(\overline{t_2})}\|_\infty, \|\mathbf{r}_2^{(\overline{t_3})}\|_\infty \leq \frac{4}{3}$ can be shown as part of the basis of the induction. We
 1323 make the argument for an even k . From round \underline{t}_k until Player 1 plays $a_1(k)$ with probability 1, the
 1324 regret of $a_1(k)$ increases by $k - (k\mathbf{x}_1^{(t)}[a_1(k)] + (k-1)\mathbf{x}_1^{(t)}[a_1(k-2)]) = \mathbf{x}_1^{(t)}[a_1(k-2)]$ and the
 1325 regret of $a_1(k-2)$ increases by $k-1 - (k\mathbf{x}_1^{(t)}[a_1(k)] + (k-1)\mathbf{x}_1^{(t)}[a_1(k-2)]) = -1 + \mathbf{x}_1^{(t)}[a_1(k-2)]$;
 1326 that is, it decreases by $1 - \mathbf{x}_1^{(t)}[a_1(k-2)]$. Let t' be the first round for which $\mathbf{r}^{(t')}[a_1(k)] \geq$
 1327 $\mathbf{r}^{(t')}[a_1(k-2)]$. It holds that $\mathbf{r}^{(t')}[a_1(k)] \leq \|\mathbf{r}_1^{(\overline{t_{k-2}})}\|^+ + \infty + 1$ since the regret of $a_1(k)$ is increasing
 1328 by at most 1 in each round and $\mathbf{r}^{(t')}[a_1(k-2)] \leq \|\mathbf{r}_1^{(\overline{t_{k-2}})}\|^+ + \infty$. From then onward, the regret of
 1329 $a_1(k)$ is increasing by at most 1/2 while the regret of $a_1(k-2)$ is decreasing by at least 1/2. Thus,
 1330 it will take at most $\lceil 2|\mathbf{r}^{(t')}[a_1(k-2)]| \rceil \leq 2|\mathbf{r}^{(t')}[a_1(k-2)]| + 1 \leq 2\|\mathbf{r}_1^{(\overline{t_{k-2}})}\|^+ + \infty + 1$ rounds
 1331 for the regret of $a_1(k-2)$ to be nonpositive. During that time, the regret of $a_1(k)$ can increase by
 1332 at most $\|\mathbf{r}_1^{(\overline{t_{k-2}})}\|^+ + \infty + 1$. \square
 1333

1334 *Proof of Property 4.1.* If κ is odd, it suffices to prove that in every round in which Player 1 mixes
 1335 between $a_1(\kappa)$ and $a_1(\kappa+1)$, Player 2 plays $a_2(\kappa) = a_2(\kappa+1)$ with probability 1. Similarly, if κ is
 1336 even, it suffices to prove that in every round in which Player 2 mixes between $a_2(\kappa)$ and $a_2(\kappa+1)$,
 1337 Player 1 plays $a_1(\kappa) = a_1(\kappa+1)$ with probability 1. Let us analyze the case where κ is even; the
 1338 odd case is similar. When Player 2 starts mixing more and more to $a_2(\kappa+1)$, it makes the row
 1339 $a_1(\kappa+2)$ more attractive for Player 2. By Lemmas C.12 and C.13,
 1340

$$1341 \mathbf{r}_1^{(\overline{t_{\kappa}})}[a_1(\kappa+2)] \leq -\frac{\kappa-1}{2}T_\kappa - \frac{\kappa-2}{2}T_{\kappa-1} \leq -\frac{(\kappa-1)!}{2^{\kappa-2}}T_4 - \frac{(\kappa-2)!}{2^{\kappa-3}}T_3. \quad (26)$$

1342 At the same time, Lemma C.14 implies that Player 2 is mixing between $a_2(\kappa)$ and $a_2(\kappa+1)$ for at
 1343 most $3\|\mathbf{r}_2^{(\overline{t_{\kappa-1}})}\|^+ + \infty + 2 \leq 5 \cdot 2^{(\kappa-1)/2} + 2$ rounds. To see this, we observe that it takes at most
 1344 $\lceil \|\mathbf{r}_2^{(\overline{t_{\kappa-1}})}\|^+ + \infty \rceil$ rounds for the action $a_2(\kappa+1)$ to be played with at least the same probability as
 1345 $a_2(\kappa)$, which in turn holds because the regret of $a_2(\kappa+1)$ increases by $\mathbf{x}_2^{(t)}[a_2(\kappa)]$ while the regret
 1346 of $a_2(\kappa)$ decreases by $1 - \mathbf{x}_2^{(t)}[a_2(\kappa)]$. From then on, the regret of $a_2(\kappa)$ decreases by at least 1/2
 1347

in each round, so it takes at most $\lceil 2\lceil \|\mathbf{r}_2^{(\overline{t_{\kappa-1}})}\| + \|\infty\| \rceil$ rounds for it to be nonpositive. We now claim that, by (26), action $a_1(\kappa + 2)$ is never played during those rounds. The reason is that since $T_3 \geq 5$ and $T_4 \geq 20$ (by our inductive basis),

$$\mathbf{r}_1^{(\overline{t_{\kappa}})}[a_1(\kappa + 2)] \leq -20 \frac{(\kappa - 1)!}{2^{\kappa-2}} - 5 \frac{(\kappa - 2)!}{2^{\kappa-3}}$$

and in each round the regret of $a_1(\kappa + 2)$ can only increase additively by 2. Since

$$\frac{1}{2} \left(20 \frac{(\kappa - 1)!}{2^{\kappa-2}} + 5 \frac{(\kappa - 2)!}{2^{\kappa-3}} \right) > 5 \cdot 2^{(\kappa-1)/2} + 2 \quad \forall \kappa \geq 4,$$

the inductive step follows. \square

The next lemma shows that, under the invariance of [Property 4.1](#), the only way to reach an approximate Nash equilibrium is to start playing the actions corresponding to $2m - 1$, which is the maximum payoff in the matrix.

Lemma C.15. *Consider any strategy profile $(\mathbf{x}_1, \mathbf{x}_2)$ such that Player 1 only assigns positive probability to actions in $\{a_1(k), a_1(k + 1)\}$ and Player 2 only assigns positive probability to actions in $\{a_2(k), a_2(k + 1)\}$, where $k + 1 < 2m - 1$. Then either Player 1 or Player 2 has a deviation benefit of at least $1/k + 2 - \gamma$ for any $\gamma > 0$.*

Proof. By construction of the game, either $a_1(k) = a_1(k + 1)$ or $a_2(k) = a_2(k + 1)$. We can assume that $a_1(k) = a_1(k + 1)$; the argument when $a_2(k) = a_2(k + 1)$ is symmetric. Let p be the probability Player 2 places at $a_2(k + 1)$ and $1 - p$ at $a_2(k)$. Suppose that the deviation benefit of each player is at most ϵ . The utility of Player 2 under the current strategy profile is $k(1 - p) + (k + 1)p = k + p$, while deviating to $a_2(k + 1)$ gives $k + 1$. So, $p > 1 - \epsilon$. Given that $k + 1 < 2m - 1$, Player 1 can deviate to $a_1(k + 2)$ to obtain a utility of $p(k + 2) \leq k + p + \epsilon$. Combining with the fact that $p \geq 1 - \epsilon$, this implies $\epsilon \geq 1/k + 2$. \square

Uniform initialization So far, our lower bound assumes that one can initialize RM arbitrarily. A more common initialization prescribes randomizing uniformly at random. In what follows, we point out that the same argument works by considering a different common payoff matrix; namely,

$$\tilde{\mathbf{B}}[a_1, a_2] := \begin{cases} \mathbf{A}[a_1, a_2] & \text{if } 1 \leq a_1 \leq m \text{ and } 1 \leq a_2 \leq m; \\ 1 - \frac{1}{m} & \text{if } a_1 = 1 \text{ and } m + 1 \leq a_2 \leq 2m; \\ 1 - \frac{3}{m} & \text{if } m + 1 \leq a_1 \leq 2m \text{ and } a_2 = 1; \\ -\frac{1}{m} \sum_{a'_2=1}^m \mathbf{A}[a_1, a'_2] & \text{if } 2 \leq a_1 \leq m \text{ and } m + 1 \leq a_2 \leq 2m; \\ -\frac{1}{m} \sum_{a'_1=1}^m \mathbf{A}[a'_1, a_2] & \text{if } m + 1 \leq a_1 \leq 2m \text{ and } 2 \leq a_2 \leq m; \\ \frac{1}{m^2} \sum_{a'_1=1}^m \sum_{a'_2=1}^m \mathbf{A}[a'_1, a'_2] - \frac{2}{m} & \text{if } m + 1 \leq a_1 \leq 2m \text{ and } m + 1 \leq a_2 \leq 2m. \end{cases}$$

By assumption, $\mathbf{x}_1^{(1)}$ and $\mathbf{x}_2^{(1)}$ are the uniform distributions over $[2m]$. We observe that

$$\begin{aligned} \sum_{a_1=1}^{2m} \sum_{a_2=1}^{2m} \tilde{\mathbf{B}}[a_1, a_2] &= \sum_{a_1=1}^m \sum_{a_2=1}^m \mathbf{A}[a_1, a_2] + m \left(1 - \frac{1}{m}\right) + m \left(1 - \frac{3}{m}\right) \\ &\quad - \sum_{a_1=2}^m \sum_{a_2=1}^m \mathbf{A}[a_1, a_2] - \sum_{a_1=1}^m \sum_{a_2=2}^m \mathbf{A}[a_1, a_2] + \sum_{a_1=1}^m \sum_{a_2=1}^m \mathbf{A}[a_1, a_2] - 2m \\ &= 2 \sum_{a_1=1}^m \sum_{a_2=1}^m \mathbf{A}[a_1, a_2] - \sum_{a_1=2}^m \sum_{a_2=1}^m \mathbf{A}[a_1, a_2] - 1 - \sum_{a_1=1}^m \sum_{a_2=2}^m \mathbf{A}[a_1, a_2] - 3 = 0 \end{aligned}$$

since $1 = \sum_{a_2=1}^m \mathbf{A}[a_1 = 1, a_2]$ and $3 = \sum_{a_1=1}^m \mathbf{A}[a_1, a_2 = 1]$. In turn, this implies that $\mathbf{x}_1^{(1)} \tilde{\mathbf{B}} \mathbf{x}_2^{(1)} = \frac{1}{4m^2} \sum_{a_1=1}^{2m} \sum_{a_2=1}^{2m} \tilde{\mathbf{B}}[a_1, a_2] = 0$. Let $\mathbf{u}_1^{(1)} = \tilde{\mathbf{B}} \mathbf{x}_2^{(1)}$ and $\mathbf{u}_2^{(1)} = \tilde{\mathbf{B}}^\top \mathbf{x}_1^{(1)}$. We have

$$\mathbf{u}_1^{(1)}[a_1] = \begin{cases} \frac{1}{2m} \sum_{a_2=1}^m \mathbf{A}[a_1 = 1, a_2] + \frac{1}{2m} m \left(1 - \frac{1}{m}\right) = \frac{1}{2} & \text{if } a_1 = 1; \\ \frac{1}{2m} \sum_{a_2=1}^m \mathbf{A}[a_1, a_2] - \frac{1}{2m} m \frac{1}{m} \sum_{a_2=1}^m \mathbf{A}[a_1, a_2] = 0 & \text{if } 2 \leq a_1 \leq m; \\ \frac{1}{2m} (1 - 2) = -\frac{1}{2m} & \text{if } m + 1 \leq a_1 \leq 2m. \end{cases}$$

1404 Since $\langle \mathbf{x}_1^{(1)}, \mathbf{u}_1^{(1)} \rangle = 0$,

1405

1406

1407

1408

1409

$$\mathbf{r}_1^{(1)}[a_1] = \begin{cases} \frac{1}{2} & \text{if } a_1 = 1; \\ 0 & \text{if } 2 \leq a_1 \leq m; \\ -\frac{1}{2m} & \text{if } m + 1 \leq a_1 \leq 2m. \end{cases}$$

1410

Similarly, we have

1411

1412

1413

1414

$$\mathbf{r}_2^{(1)}[a_2] = \begin{cases} \frac{1}{2} & \text{if } a_2 = 1; \\ 0 & \text{if } 2 \leq a_2 \leq m; \\ -\frac{1}{2m} & \text{if } m + 1 \leq a_2 \leq 2m. \end{cases}$$

1415

1416

1417

1418

1419

1420

1421

1422

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

1433

1434

1435

1436

1437

1438

1439

1440

1441

1442

1443

1444

1445

1446

1447

1448

1449

1450

1451

1452

1453

1454

1455

1456

1457

As a result, the second round sees both players playing the first action with probability 1. In view of the fact that $\mathbf{A}[a_1 = 1, a_2] < 1$ when $1 \leq a_2 \leq m$ and $\mathbf{A}[a_1, a_2 = 1] < 1$ when $1 \leq a_1 \leq m$, we immediately revert to the previous analysis concerning simultaneous RM executed on the game described in (6). Thus, we arrive at the following lower bound (by a change of variables $m' \leftarrow 2m$).

Corollary C.16. *Simultaneous RM requires $m^{\Omega(m)}$ rounds to converge to a $\frac{1}{m+1}$ -Nash equilibrium in two-player $m \times m$ identical-interest games. This holds even when both players initialize at the uniform random strategy.*