

# Long-term Fairness with Bounded Worst-case Losses

Gabriel Balan and Dana Richards and Sean Luke

Department of Computer Science, George Mason University  
4400 University Drive, MSN 4A5, Fairfax, VA, 20030, USA  
{gbalan, richards, sean}@cs.gmu.edu

## Abstract

How does one repeatedly choose actions so as to be fairest to the multiple beneficiaries of those actions? We examine approaches to discovering sequences of actions for which the worst-off beneficiaries are treated maximally well, then secondarily the second-worst-off, and so on. We formulate the problem for the situation where the sequence of action choices continues forever; this problem may be reduced to a set of linear programs. We then extend the problem to situations where the game ends at some unknown finite time in the future. We demonstrate that an optimal solution is NP-hard, and present two good approximation algorithms.

## Introduction

Consider the problem of repeatedly assigning two AI professors to teach two classes offered by their department each semester. One class is much harder than the other one, so during any single semester any one-to-one assignment is unfair to one of the professors. One fair solution would be for the professors to teach both classes together. But assuming this requires more overall effort than teaching the classes separately, this solution would be inefficient over the long run.

There is of course a better solution. If the two professors instead took turns teaching the hard class, then in the long run their average utilities would be *more fair* than in the one-to-one assignments and *more efficient* than in the sharing assignment. This is the rough idea behind *long-term fairness*: repeated interactions offer opportunities for improved efficiency and fairness over the single interaction scenario. But in general the solutions will not be as simple as alternating assignments (e.g. suppose we extended the previous example with multiple assignments, involving many professors and classes).

The research presented here examines the following framework: there are a number of *beneficiaries* (e.g. professors), which receive different *rewards* from each of a finite set of *actions* (e.g. class assignments). The actions are chosen with replacement, and actions chosen early do not restrict what actions can be chosen later, or their rewards. This framework, borrowed from (Verbeeck, Parent,

and Nowé 2003), is very similar to the repeated normal-form game framework from game theory, except there is a single *decision maker* that chooses actions for the good of all beneficiaries.

We are ultimately interested in the game-theoretic multiple decision maker case. But it turns out that the single decision maker situation has been inadequately studied, and so our focus here is on that case.

We define a beneficiary's *utility* as the average of all rewards he received in the past. We will use the term *utility profile* to refer to the vector of utilities (one utility per beneficiary) that they derive from past actions.

In order to compare utility profiles, we use a fairness concept called *leximin*. Leximin is a total-order relation defined as follows: a utility profile  $U_1$  is preferable to another utility profile  $U_2$  if the beneficiary with the minimum utility in  $U_1$  is better-off than the minimum utility beneficiary in  $U_2$ ; ties are broken by comparing the utilities of the the second-worst beneficiary in each utility profile, then third-worst, and so on. From a procedural point of view, one should sort in increasing order the two utility profiles, and compare the sorted versions lexicographically. We choose leximin for our work because it is widely used in the literature, it is well understood, and it incorporates a certain measure of efficiency (a leximin-optimal profile is also Pareto-optimal), making it a good fit for long-term fairness.

In this paper we derive upper bounds for this notion of utility in infinitely repeated games and propose a particular approach to approximation solutions for more realistic cases (i.e. a finite sequence of actions).

## Infinite-length Games

An action may give high rewards to some beneficiaries and low rewards to others, so sticking to just one action might be unfair to some beneficiaries, and thus leximin-undesirable. However, if beneficiaries receive different rewards from different actions we may be able to improve all beneficiaries' reward averages by playing a combination of actions.

One approach to playing combinations is to use a periodic, repeated sequence of actions. In previous literature (Verbeeck, Parent, and Nowé 2003), distributed algorithms for discovering such sequences found suboptimal ones. We will show optimal solutions, albeit with non-distributed algorithms.

Periodic sequences are intuitively appealing, but even if one can find the periodic sequence with the best limit, that can still be suboptimal. We will show at the end of this section that in some problem instances with irrational coefficients there might exist infinite non-periodic sequences that achieve, at the limit, leximin-superior utility profiles to any utility profile achievable by a periodic sequence. Our algorithms are guaranteed to produce sequences converging to the optimal utility profile, and, whenever possible, those sequences are periodic.

In this section we will (1) prove that convergent sequences (periodic or non-periodic) are sufficient and (2) identify the class of sequences with the leximin-optimal limit-point. In the following sections we will propose additional requirements to impose on this class of sequences and then provide algorithms that produce sequences satisfying these requirements.

**Problem 1. Base Problem** *Let there be a set  $\mathbb{A}$  of  $n_a$  actions affecting a set  $\mathbb{B}$  of  $n_b$  beneficiaries through the reward function  $R : \mathbb{A} \times \mathbb{B} \rightarrow \mathbb{R}$ . Let  $\mathbb{S}$  be the set of infinite sequences of actions from  $\mathbb{A}$ :*

$$\mathbb{S} = \{S = \langle s_1, s_2, \dots \rangle \mid \forall i \in \mathbb{N} : s_i \in \mathbb{A}\}.$$

*Let  $\mathbb{U}$  be the set of all possible utility profiles (vectors) achievable from following any sequence for any number of time steps:*

$$\mathbb{U} = \{U \in \mathbb{R}^{n_b} \mid \exists t \in \mathbb{N}, \exists S \in \mathbb{S} : U_b = \frac{1}{t} \sum_{j=1}^t R(s_j, b)\}.$$

*The goal of the problem is to find  $U^* = \sup(\mathbb{U})$ , the supremum (least upper bound) with respect to leximin over the set of all achievable utility profiles.*

In order to solve this problem we focus on a subset of “well-behaved” sequences that achieve the entire set of possible utility profiles  $\mathbb{U}$ . We then provide a mapping of that subset into the  $n_a$ -dimensional simplex which can be searched efficiently using linear-programming-based algorithms from the literature.

To make things easier, we refer to the elements of  $\mathbb{A}$  as  $\{1 \dots n_a\}$  and the elements of  $\mathbb{B}$  as  $\{1 \dots n_b\}$ . Let  $\mathbb{S}' \subset \mathbb{S}$  be the set of all sequences  $S$  where the *proportions* of the  $n_a$  actions converge. Formally:

$$\mathbb{S}' = \{S \in \mathbb{S} \mid \forall j \in \mathbb{A} : \exists \lim_{t \rightarrow \infty} \frac{1}{t} k_j(S_t)\}$$

where  $S_t$  is the subsequence of  $S$  consisting of the first  $t$  elements and  $k$  is the *count* function (so  $k_j(S_t)$  is equal to the number of times action  $j$  is used in the first  $t$  positions of sequence  $S$ ). We denote with  $F(S)$  the vector of action proportions (or fractions) for  $S$  (i.e.  $F_j(S) = \lim_{t \rightarrow \infty} \frac{1}{t} k_j(S_t)$ ). Let  $U(S_t)$  be the vector of utilities achieved after following the first  $t$  steps of sequence  $S$  (i.e.  $U$ 's component for beneficiary  $b$  is  $U_b(S_t) = \frac{1}{t} \sum_{i=1}^t R(s_i, b)$ ). Because the action proportions converge, the sequence of utility vectors  $\langle U(S_1), U(S_2), \dots \rangle$  also converges component by component; we denote its limit point by the vector  $U(S)$ :

$$U_b(S) = \lim_{t \rightarrow \infty} U_b(S_t) = \sum_{j=1}^{n_a} F_j(S) R(s_j, b). \quad (1)$$

It can be shown that  $\mathbb{U}'$ , the set of utility profiles achievable by sequences in  $\mathbb{S}'$ , is equal to  $\mathbb{U}$ .<sup>‡</sup> Moreover, if  $\mathbb{U}''$  is the set of limit points of  $\mathbb{U}'$  ( $\mathbb{U}'' = \{U \mid \exists S \in \mathbb{S}' : U = U(S)\}$ ), then we can prove that:<sup>‡</sup>

$$\sup(\mathbb{U} = \mathbb{U}') = \sup(\mathbb{U}''). \quad (2)$$

Because  $U(S)$  depends only on  $F(S)$ , there is an obvious one-to-one correspondence between  $\mathbb{U}''$  and the  $n_a$ -dimensional unit simplex.

Once reformulated as an optimization problem over a compact and convex set, the problem becomes readily solvable with the algorithm proposed in (Potters and Tijs 1992). In our particular case, the algorithm consists of solving  $O(n_a^2)$  linear programs (LPs), but approximate results based on a floating-point fixed-length representation might be required to make sure the complexity of this algorithm does not dominate that of the algorithms we propose here. The algorithm produces  $U^*$ , the unique leximin-optimal utility vector, and  $F^*$ , a point (not necessarily unique) inside the simplex such that:  $\forall b \in \mathbb{B} : \sum_{i=1}^{n_a} R(b, i) \times F_i^* = U_b^*$ . Since  $U^*$  is the leximin-maximal utility vector in  $\mathbb{U}''$  (by the algorithm's guarantee), and it is unique, it implies that  $U^* = \sup(\mathbb{U}'')$ , and, by Equation 2:

$$U^* = \sup(\mathbb{U}). \quad (3)$$

We revisit the existence of a periodic sequence for a given optimal utility profile solution  $U^*$ . If such a periodic policy exists, let  $c_j$  denote the number of times of action  $j$  appears in one period (we call  $c_j$  the *multiplicity* of action  $j$ ). Then by normalizing the vector  $c = [c_1 \dots c_{n_a}]$  one should obtain a valid  $F^*$  vector. If all rewards are rational then there must exist an  $F^*$  with rational components, and hence a periodic sequence.<sup>‡</sup> However, irrational rewards could mean there is no  $F^*$  with rational coefficients, in which case there is no optimal periodic sequence.

For example, if both professors get a reward of 0 from teaching the hard class, but teaching the easy class gives one professor a reward of 1 and the other a reward of  $\sqrt{2}$ , then  $U^* = [\frac{\sqrt{2}}{1+\sqrt{2}}, \frac{\sqrt{2}}{1+\sqrt{2}}]$  and there is a unique  $F^* = [\frac{\sqrt{2}}{1+\sqrt{2}}, \frac{1}{1+\sqrt{2}}]$ . We make the observation that depending on the values in  $F^*$ , the multiplicities could be arbitrarily large. Thus, one can see the non-periodic sequences due to  $F^*$  having (some) irrational components as special cases of periodic sequences where the period length is infinite.

## Finite-length Games

Solving the **Base Problem** provides  $U^*$ , the least upper bound over the set of achievable utility profiles. If the game is guaranteed to last forever, one could be satisfied with the goal of finding a sequence that achieves  $U^*$  at the limit. But in more practical applications, one must consider the implications of having the process end after a finite number of steps, or more generally, having the length of the sequence of actions drawn from some probability distribution.

<sup>‡</sup>Due to limited space we omit all proofs and mathematical derivations. They may be found in (Balan, Richards, and Luke 2008), as well as an extended discussion.

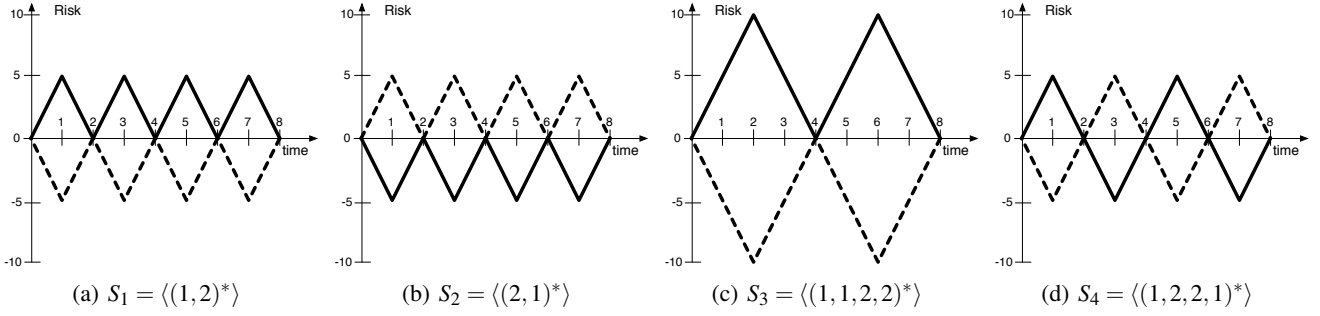


Figure 1: Risks as functions of time for the two beneficiaries in Example 1 (the solid line for the first and the dotted line for the second) as produced by various periodic sequences.

Consider the professor-assignment example from the introduction: the first action corresponds to the assignment where the first professor teaches the easy class, the second action corresponds to the assignment where the second professor teaches the easy class, and the third action is the assignment where they teach the classes together:

	Beneficiaries	
	$b_1$	$b_2$
<b>Example 1.</b>		
Actions	1	2
	2	3
	10	0
	0	10
	1	1

The utility least upper bound for Example 1 is  $U^* = [5, 5]$ , which can only be achieved through  $F^* = [0.5, 0.5, 0]$ , and which coincides with the optimal solution if the game lasts for an even number of steps: use actions 1 and 2 in equal proportions. However, one should use action 3 if the game lasts for only one round. This shows that the optimal sequence of decisions depends on the duration of the game, and so if the duration is not known in advance, some sort of tradeoff might be required.

One could choose an action stochastically, using the values in  $F^*$  as probabilities. This approach produces the leximin-optimal *expected rewards* (equal to  $U^*$ ). This sort of guarantee is usually referred to in the economics literature as *ex-ante* fairness (ex-ante means “beforehand” in latin). If all beneficiaries are risk-neutral, a mechanism choosing repeatedly from the  $F^*$  distribution promises leximin-optimal *expected utilities* before the process starts, but not necessarily leximin-optimal actual utilities when the game ends. For Example 1 this approach would flip a coin between actions 1 and 2; although the expected rewards (and utilities) are equal to  $U^*$ , if the game ends after two rounds, there is a 50% chance that one beneficiary accumulated 10 reward units and the other none. Moreover, if the game continues, the probability that the difference between the two will shrink is equal to the probability that the difference will grow. Thus, the weakness of this ex-ante approach would be its lack of concern with paying reparations for unfairness resulting from past randomness.

One can address the weaknesses of the previous method with a deterministic approach, aiming to find a sequence which leximin-optimizes beneficiaries’ *expected utilities* weighted by the probabilities in the distribution of game

lengths. Unlike the previous method, this could use action 3 in Example 1, provided there is a high enough probability that the game will end after one step.

Both previous methods assume *risk-neutral* beneficiaries. Given that *leximin* is a *risk-averse* fairness concept (“no one is left behind”), it makes sense to instead focus on a *risk-averse* solution approach: minimize the largest amount a beneficiary risks losing due to the game ending prematurely. We define the *risk* of beneficiary  $k$  at time step  $t$  while executing sequence  $S$  as the difference between the first  $t$  steps, and the amount he was entitled to, which is  $t \times U_k^*$ .

$$Risk_b(S, t) = \sum_{i=1}^t R(s_i, b) - t \times U_b^*. \quad (4)$$

We use the term *worst risk* (WR) of a sequence  $S$  to mean a lower bound on all risk values, regardless of beneficiary (i.e.  $WR \leq Risk_b(S, t)$ ,  $\forall t \in \mathbb{N}$  and  $\forall b \in \mathbb{B}$ ). The worst risk is a negative value, and its absolute value is an upper bound on the largest amount that any beneficiary can lose.

This approach is not meant to replace expectation optimization; we are simply adding another tool to the arsenal available for this problem. It is reasonable to expect the computational complexity of the expectation-optimization approach to grow with the maximum game length, and so our proposed solution is particularly suitable to scenarios with very large game durations (our approach is insensitive to the distribution of game durations). Furthermore, lower-bounding the worst risk will always make the utilities converge to  $U^*$  as  $t \rightarrow \infty$ .<sup>‡</sup>

For convenience, we define for each action  $j$  a vector  $X_j = [X_{j,1} \dots X_{j,n_b}]$ , where  $X_{j,b} = R(j, b) - U_b^*$ , the amount action  $j$  changes the risk of beneficiary  $b$ .

$$Risk_b(S, t) = \sum_{i=1}^t [R(s_i, b) - U_b^*] = \sum_{i=1}^t X_{s_i, b} \quad (5)$$

$$\forall b \in \mathbb{B} : \sum_{j=1}^{n_a} X_{j,b} \times F_j^* = 0 \quad (6)$$

Let us illustrate the concept of risk using Example 1. The  $X$  vectors are as follows:  $X_1 = [5, -5]$ ,  $X_2 = [-5, 5]$ , and  $X_3 = [-4, -4]$ . We will compare the following periodic ac-

tion sequences (see Figure 1):

$$\begin{aligned} S_1 &= \langle 1, 2, 1, 2, \dots \rangle \\ S_2 &= \langle 2, 1, 2, 1, \dots \rangle \\ S_3 &= \langle 1, 1, 2, 2, 1, 1, 2, 2, \dots \rangle \\ S_4 &= \langle 1, 2, 2, 1, 1, 2, 2, 1, \dots \rangle \end{aligned}$$

Sequence  $S_1$  makes the first beneficiary's risks equal to 5 on odd steps and 0 on even steps, while the second beneficiary's risks are 0 on odd steps and  $-5$  on even ones. Therefore the second beneficiary risks coming up 5 units short if the game stops after an odd number of steps and breaks even otherwise. Sequence  $S_2$  has identical effects as  $S_1$ , but to different beneficiaries, so the two are equally good. Using sequence  $S_3$ , the second beneficiary risks losing as much as 10 units, so we prefer  $S_1$  (or  $S_2$ ) to  $S_3$ .

Note that using sequence  $S_1$  makes the second beneficiary take all the "bad" (negative) risks, while the first beneficiary takes only "good" (positive) risks. Although sequence  $S_1$  is as fair as possible with respect to average utilities, one cannot help but notice a "second degree" unfairness. Instead, one might consider the goal of optimizing cumulative risks (e.g. sequence  $S_4$  alternates which beneficiary is taking negative risks). This leads to the following paradox: although intuitively one prefers  $S_4$  to  $S_1$  or  $S_2$ , the worst risks for all beneficiaries during sequences  $S_1$  or  $S_2$  (i.e.  $-5$  and  $0$ ) are leximin superior to the worst risks during  $S_4$  (i.e.  $-5$  and  $-5$ ). This suggests that one might consider optimizing only the "min" of the worst risks, instead of leximin optimizing all beneficiaries' worst risks.

**Problem 2. Worst Risk Maximization Problem** Given the setup in the **Base Problem**, find an infinite sequence of actions  $S$  with a maximal worst risk.

**Theorem 1.** Problem 2 is intractable<sup>‡</sup> (since even finding a sequence of a given finite length is NP-hard).

## Algorithms

The most intuitive solution to optimizing worst risk(s) is to try to keep all risks as large as possible at all times, i.e. greedily choose the action that leximin optimizes risks during the next step (choose the action with the best immediate effects). We will refer to this strategy as **GR** (Greedy leximin-optimizing next-step Risks). One weakness of this approach is that it assumes the world ends at the next step. In Example 1 it will always choose action 3, leading to arbitrarily bad risks for both beneficiaries if the game lasts long enough. Note that  $F_3^* = 0$ , meaning action 3 should not be played at all. One can easily extend the greedy heuristic to ignore the "unusable" actions, but that is not enough to prevent risks from getting arbitrarily bad. Consider Example 2, with  $U^* = [240, 240, 240]$  and unique  $F^* = [\frac{4}{15}, \frac{6}{15}, \frac{5}{15}]$ . For this problem **GR** will choose action 1 repeatedly, leading to arbitrarily bad risks for the first and last beneficiaries. The reason is that action 2 and 3 hurt one of those beneficiaries more than action 1 hurts any of them, and **GR** lacks the look-ahead to see the benefits of using actions 2 and 3.

		Beneficiaries			
		b1	b2	b3	
<b>Example 2.</b>	Actions	1	-20	30	-20
		2	60	30	-40
		3	-56	-60	64

The algorithms we propose in this paper are based on the following observation. All beneficiaries' risks are zero at time  $t$  if the number of times each action  $j$  was used up to time  $t$  are all proportional to the corresponding components of some  $F^*$  vector (i.e.  $k_j(S_t) = F_j^* \times t, \forall j$ ). Ideally all  $k_j(S_t) = F_j^* \times t$  at all times, but since the  $k_j$  functions only take integer values, that is not always possible. It is intuitive that the risks cannot get very bad at any time  $t'$  if the relative counts ( $k_j/t'$ ) for how much each action was used up to time  $t'$  are close enough to the corresponding values in  $F^*$ . There are many ways one can construct sequences  $S$  to keep the  $k_j(S_t)$  values close enough to the  $F_j^* \times t$  values; we will present two simple methods, and derive risk upper and lower bounds for both of them.

Our methods completely ignore the actions  $j$  which are not "used" in  $F^*$  (i.e.  $F_j^* = 0$ ). Let  $n_a^*$  be the number of actions "used" in  $F^*$ ; for simplicity, we rename the actions (reorder the dimensions of the simplex) such that all actions used in  $F^*$  come before the other ones:  $\forall j \in \{1 \dots n_a^*\} : F_j^* > 0$  and  $\forall j \in \{n_a^* + 1 \dots n_a\} : F_j^* = 0$ .

**Method 1** Intuitively, this method chooses an action that, up to that point, has been used the least relative to how often the action should have been used. Consequently, an action  $j$  can be chosen at time  $t + 1$  if it has the smallest ratio  $\frac{k_j(D_t)}{F_j^* \times t}$ , where  $D$  is the sequence of decisions produced by this method (so  $D_t$  contains all its past decisions). Without affecting the decision process, one can eliminate  $t$  from the denominator. We formally describe this method with the following notation:

$$d_t \in \left\{ j \leq n_a^* \mid \forall i \leq n_a^* : \frac{k_j(D_{t-1})}{F_j^*} \leq \frac{k_i(D_{t-1})}{F_i^*} \right\} \quad (7)$$

where all functions  $k_j$  are extended such that  $k_j(D_0) = 0, \forall j \in \{1 \dots n_b\}$ . Note that multiple actions could tie for the minimum.

**Theorem 2.** Regardless of how the first method breaks ties, the following holds  $\forall b \in \mathbb{B}$  and  $\forall t \in \mathbb{N}^{\ddagger}$ :

$$\sum_{i=1}^{n_a^*} \min(X_{i,b}, 0) \leq Risk_b(D, t) \leq \sum_{i=1}^{n_a^*} \max(X_{i,b}, 0). \quad (8)$$

**Method 2** While the previous approach helps less-often chosen actions catch up to the others, the next approach chooses actions that — if used — will get the least ahead of the others. The sequence of decisions  $D'$  for this method is described formally as follows:

$$d'(t) \in \left\{ j \leq n_a^* \mid \forall i \leq n_a^* : \frac{k_j(D'_{t-1}) + 1}{F_j^*} \leq \frac{k_i(D'_{t-1}) + 1}{F_i^*} \right\} \quad (9)$$

**Theorem 3.** *Regardless of how the second method breaks ties, the following holds for  $\forall b \in \mathbb{B}$  and  $\forall t \in \mathbb{N}$ :*<sup>‡</sup>

$$-\sum_{i=1}^{n_a^*} \max(X_{i,b}, 0) \leq \text{Risk}_b(D', t) \leq -\sum_{i=1}^{n_a^*} \min(X_{i,b}, 0). \quad (10)$$

### Discussion

We proposed two methods for choosing actions such that the risks are always lower-bounded. While **GR** leximin-optimizes next-step risks, our methods *greedily* optimize actions' usage *frequencies* (relative counts) relative to some optimal configuration  $F^*$ ; we will call our methods **GF**.

We denote with  $LB_1$  and  $LB_2$  the lower bounds on risks guaranteed by the two **GF** methods respectively.

$$LB_1 = \min_b \sum_{i=1}^{n_a^*} \min(X_{i,b}, 0) \quad (11)$$

$$LB_2 = \min_b \sum_{i=1}^{n_a^*} -\max(X_{i,b}, 0). \quad (12)$$

There is an obvious way to unify the two: compute  $LB_1$  and  $LB_2$  ahead of time, then use the most promising method. Thus our best worst risk guarantee is:

$$WR \geq \max(LB_1, LB_2). \quad (13)$$

The computational complexity of our algorithms is  $O(\lg n_a^*)$  per time step. This is because one can use a heap to store actions'  $k_j[+1]/F_j^*$  scores, since our algorithms change a single action's score at each time step.

**Eliminating Unnecessary Actions** If the vector  $F^*$  is not unique, the particular choice of  $F^*$  influences both the time complexity (through  $n_a^*$ ), and the worst risk. Each beneficiary's worst risk is bounded below by the sum of his negative  $X$  values (or the negative sum of his positive  $X$  values), so eliminating an action can only improve the worst risk.

**Lemma 1.** *There must always exist  $F^*$  such that  $n_a^* \leq n_b$ .*<sup>‡</sup>

Based on this result, one can eliminate at least  $\max(n_a - n_b, 0)$  actions in a preprocessing phase, thus reducing the per-step complexity to  $O(\lg(\min(n_a, n_b)))$ . There exists a simple algorithm for this task based on a particular constructive proof for Carathéodory's theorem (e.g. (Florenzano and Le Van 2001)), using Gauss elimination. The complexity of the Gauss elimination is  $O(n_a n_b^2)$ , and there are at most  $n_a - 1$  iterations, so the entire pre-process of eliminating actions can be done in  $O(n_a^2 n_b^2)$ .

**Breaking Ties with GR** The lower-bounds on worst risks (Equations 11 and 12) make no assumptions on how ties are broken, which means they assume the ties are broken in the worst possible way. Finding the best way to break ties is NP-hard (see the proof of Theorem 1), so a greedy approach will have to do. The most obvious solution would be to use **GR** to break ties for **GF** (for a time complexity of  $O(n_b \times n_a^*)$  per time step).

### Related Work

The game-theoretic work in (Bhaskar 2000; Lau and Mui 2005) is concerned with finding Nash equilibria that result in alternating joint-actions (also referred to as turn-taking), but these results were tailored for specific classes of 2-by-2 games.

The starting point in our research on long-term fairness was the work in (Verbeeck, Parent, and Nowé 2003) on "periodic policies." Their reward model comes in the form of a normal-form game, but the players are actually cooperative learning agents (rather than self-interested). The process consists of the learners playing selfishly to discover a pure Nash equilibrium, while being interrupted periodically to compare accumulated rewards. The player gaining the most (in the current Nash equilibrium and overall) has its action put off-limits until the others catch up. Alternatively (Verbeeck et al. 2006), after the learners discover all pure Nash equilibria, they create a periodic policy consisting of those joint actions with the fairest outcomes. In the authors' examples the players have only two actions: a highly lucrative one and a social one they fall back on while waiting for the other(s) to catch up. Both of these greedy algorithms could lead to utility profiles arbitrarily worse than the optimum.

The **Worst Risk Maximization Problem** is actually a variant of the **Compact Vector Summation Problem** (Sevast'janov 1991): given a finite set of vectors  $X_1, \dots, X_n \in \mathbb{R}^m$  such that  $\sum_{i=1}^n X_i = 0$ , one must find a permutation  $\pi$  of  $\{1, 2, \dots, n\}$  that minimizes  $\max_{1 \leq k \leq n} \|\sum_{i=1}^k X_{\pi_i}\|$ . In the earliest such work  $\|\cdot\|$  was the Euclidian norm, so the problem consisted of ordering the vectors such that the path resulting from adding vectors one by one stays inside a minimum-radius  $m$ -dimensional circle centered at the origin. Later research has focused on results general enough to accommodate arbitrary norms (intuitively a norm is a function associating a "size" value to every vector).

We submit that the algorithm in (Sevast'janov 1991) is the most relevant algorithm for a comparison with our algorithms: it has the best guarantee and time complexity we are aware of and it accommodates "asymmetric norms." The last item is relevant because it is more meaningful to compare against algorithms trying to optimize the same function, and the function we try to optimize can be rearranged as an asymmetric norm but not a norm. Maximizing the worst risk is equivalent to minimizing is the largest absolute value of any negative coordinate of any partial sum of  $X$  vectors. This function is an asymmetric norm, but not a norm, since it satisfies the triangle inequality ( $\|y+z\| \leq \|y\| + \|z\|$ ) and positive definiteness ( $\|y\| = 0 \Rightarrow y = 0$ ), but it only satisfies the scalability condition ( $\|ky\| = |k| \times \|y\|$ ) for *positive scaling* factors (Borodin 2001).

**Guarantee Comparison** Sevast'janov's algorithm guarantees no risk will be worse than  $-M(n_a^* - 1 + \frac{1}{n_a^*})$ , where  $M = \max_{1 \leq k \leq n} \|X_k\|$  (i.e.  $-M$  is the smallest coordinate of any  $X$  vector size). We submit that  $LB_1 \geq -M(n_a^* - 1)$ . This results from Equation 11, replacing every negative  $X_{i,b}$  with  $-M$ , and noticing that for any given beneficiary  $b$  at most

$n_a^* - 1$  of his  $X$  values can be negative.<sup>‡</sup> Therefore our worst risk,  $WR \geq LB_1 \geq -M(n_a^* - 1) = -M(\min(n_a, n_b) - 1) > -M(n_b - 1 + \frac{1}{n_b})$ , the guarantee of Sevast'janov's algorithm.

**Complexity Comparison** Sevast'janov's algorithm has a complexity of  $O(n^2m^2)$ , because it picks a vector  $n - m$  times, and each such operation has a complexity of  $O(km^2)$ , where  $k$  is the number of alternatives ( $k = n, \dots, n - m$ ). An iteration in Sevast'janov's algorithm has the same complexity as an iteration in our preprocessing phase (they are both based on Gauss elimination). We also submit that the number of iterations in Sevast'janov's algorithm is at least the number of iteration in our preprocessing phase (each action has at least multiplicity 1). Therefore the time complexity of our preprocessing phase cannot be larger than that of Sevastianov's algorithm. More importantly, even if that algorithm were extended to benefit from our preprocessing phase and to explicitly deal with multiplicities (i.e.  $k = n_a^*$ ), its complexity would still be  $O(n_a^*n_b^2)$  per time step which is higher than the complexity of our algorithms (even when breaking ties with **GR**).

In summary, by eliminating unnecessary actions and only keeping track of multiplicities, we are able to offer worst case guarantees that are never worse than (and sometimes much better than) Sevast'janov's.

## Conclusions

In this paper we studied the problem of achieving certain long-term fairness guarantees in a simple repeated-game setup: (1) all beneficiaries are entitled to their socially-optimal utilities and (2) no matter when the game ends all beneficiaries are guaranteed to have received close to what they were entitled to that point. We proved that finding an optimal solution with respect to the second guarantee is NP-hard and proposed two efficient approximation algorithms.

## Acknowledgements

We thank Octav Olteanu, Joey Harrison, Zoran Duric, Alexei Samsonovich, and Alex Brodsky for their help.

## References

- Balan, G.; Richards, D.; and Luke, S. 2008. Long-term fairness with bounded worst-case losses. Technical Report GMU-CS-TR-2008-2, Department of Computer Science, George Mason University.
- Bhaskar, V. 2000. Egalitarianism and efficiency in repeated symmetric games. *Games and Economic Behavior* 32(2):247–262.
- Borodin, P. A. 2001. The Banach-Mazur theorem for spaces with asymmetric norm. *Mathematical Notes* 69:298–305.
- Florenzano, M., and Le Van, C. 2001. *Finite Dimensional Convexity and Optimization*. Springer.
- Lau, S.-H. P., and Mui, V.-L. 2005. Using turn taking to mitigate conflict and coordination problems in the battle of the sexes game. Technical Report 1129, Hong Kong Institute of Economics and Business Strategy.

Potters, J. A. M., and Tijs, S. H. 1992. The nucleolus of a matrix game and other nucleoli. *Mathematics of Operations Research* 17(1):164–174.

Sevast'janov, S. 1991. On the compact summation of vectors. *Diskretnaya Matematika* 3(3):66–72.

Verbeeck, K.; Nowé, A.; Parent, J.; and Tuyls, K. 2006. Exploring selfish reinforcement learning in repeated games with stochastic rewards. *Journal of Autonomous Agents and Multi-Agent Systems* 14(3):239–269.

Verbeeck, K.; Parent, J.; and Nowé, A. 2003. Homo equalis reinforcement learning agents for load balancing. In *Innovative Concepts for Agent-Based Systems: 1<sup>st</sup> International Workshop on Radical Agent Concepts*, volume 2564 of *Lecture Notes in Computer Science*, 81–91.