

Algorithm Selection in Optimization and Application to Angry Birds

Shahaf S. Shperberg
 CS Department
 Ben Gurion University
 shperbsh@post.bgu.ac.il

Solomon Eyal Shimony
 CS Department
 Ben Gurion U., UMass Lowell
 shimony@cs.bgu.ac.il

Avinoam Yehezkel
 CS Department
 Ben Gurion University
 yehezavi@post.bgu.ac.il

Abstract

Consider the MaxScore algorithm selection problem: given some optimization problem instances, a set of algorithms that solve them, and a time limit, what is the optimal policy for selecting (algorithm, instance) runs so as to maximize the sum of solution qualities for all problem instances?

We analyze the computational complexity of restrictions of MaxScore (NP-hard), and provide a dynamic programming approximation algorithm. This algorithm, as well as new greedy algorithms, are evaluated empirically on data from agent runs on Angry Birds problem instances. Results show a significant improvement over a hyper-agent greedy scheme from related work.

1 Introduction

Algorithm selection is of significant interest to researchers in AI, and other fields where more than one algorithm is available to solve problems under computational resource constraints (Rice 1976; Huberman, Lukose, and Hogg 1997; Xu et al. 2008). This paper examines a variant of algorithm selection (“MaxScore”) where one needs to solve a *set* of optimization problems, with computational resource (assumed here to be time) limitation being over the entire set. This generalizes the standard setting handled in, e.g. SAT solver algorithm portfolios, where the time limit is separate for each individual problem instance.

Our original motivation for the MaxScore setting was combining multiple programs that compete in the AI Angry Birds (Copyright Rovio Entertainment) competition, on which we also base the empirical results of this paper. Angry birds is a physical simulation video game. In the Albirds competition, each agent program (or human) is presented with N previously unseen game levels (problem instances). The agents can select a level to play, where at each level the agent is presented with a screen-shot representing a physical simulation. The agent is supposed to kill off all the pigs with catapulted birds (shots), thereby completing (winning) the level. Points for completed levels are gained for destroying objects (pigs and blocks), and for using as few birds as possible. The agent may play any level as many times as desired, until its overall allocated time (typically 10 or 30

minutes) expires. Level score is the maximum achieved in all attempts, with total score being the sum of level scores (typically 4 or 8 levels in past competitions).

Of the numerous AI techniques used in AI birds agent programs, to-date none have achieved near-human performance; each program has strengths and weaknesses in different areas of the game. Rather than attempting to create a better AI for Angry Birds, our goal here is one of meta-reasoning: use a portfolio of existing programs to play better, an idea suggested originally by (Stephenson and Renz 2017) with promising initial results. The focus of this paper is on how to do this combination optimally given the available information, in a decision-theoretic sense, rather than on the aspect of learning to fit the program to the problem instance. In addition to the research interest of this meta-reasoning problem, such an optimization might have an impact on algorithm portfolio optimization in general.

Informally (see Section 2 for the formal definition) we are given a set of problem instances (levels), to each of which we can apply any of a set of given algorithms (agents). Each such application uses up an unknown amount of time, and results in a score for the level that can be observed after each algorithm run terminates. Our meta-reasoning problem is to find a policy for selecting which agent to apply to which level, at any point in time, such that the sum of scores for all the levels at timeout is maximized. Note that in order to make this policy optimization well defined, one must specify some prior knowledge about scores and runtimes. In this paper we assume that these are specified by random variables with known distribution models.

Following the formal problem statement (Section 2), we analyze the computational complexity of restricted versions of our score maximization problem (Section 3: NP-hard even with *known* independent scores and runtimes). An approximation algorithm for one simple case is proposed, as well as faster greedy heuristic-based algorithms (Section 3). Empirical evaluation on scores and runtimes gathered from actual agent program runs show that using greedy expected improvement was near-optimal in practice (Section 4), and much better than the greedy scheme based on just expected score from (Stephenson and Renz 2017).

We then briefly examine the case of *unknown* independent distributions. These we treat as a distribution over performance profiles, i.e. a distribution over score and runtime

distributions (Section 5), with parameters that have to be learned from previously observed problem instances. This results in a model with induced dependencies, which makes the meta-level decision problem harder, in addition to the above learning problem. A naive learning method based on features from (Stephenson and Renz 2017; Tziortziotis, Pagiannis, and Blekas 2016) is proposed, and an empirical evaluation shows that our proposed greedy algorithm (coupled with belief updating) still performs well despite the low-quality learning.

2 Formal Problem Statement

This section describes the formal metareasoning MaxScore problem. In order to make the statement as general as possible, we abstract away from Angry Birds programs and levels, and present this as a sequential decision problem under uncertainty.

A MaxScore problem is a 4-tuple (I, A, T, P) : I is a set of problem instances to be optimized (game levels in Angry Birds); A is a set of algorithms (or agent programs); T is a time limit; and P is a known distribution model over problems in I , agents in A , that describes the (non-negative) score $S(a, l, i)$ achieved by agent $a \in A$ and the (positive valued) runtime $T(a, l, i)$ of a when applied to problem instance $l \in I$ during decision-making round i of the problem-solving task (or game play). Distribution P (also known as a *performance profile* (Zilberstein and Russell 1996)) can be defined in various ways. We consider the following cases for P :

1. P is deterministic.
2. P is a known distribution with independent variables.
3. P is a distribution with dependent variables (some of which are unobservable).

A policy π is a mapping from process histories to actions. The process is to select, at each round $i \geq 1$, an agent program $a(i)$ to apply to a problem instance (level) $l(i)$, given the past observations, according to π . The results $S(a(i), l(i), i)$ and $T(a(i), l(i), i)$ are observed after the selection at round i . Then i is incremented as we go to the next round. The process stops when we reach the time limit, i.e. at the first k such that:

$$\sum_{i=1}^k T(a(i), l(i), i) > T$$

The score of the process is the sum of maximal achieved scores for each problem instance, i.e.

$$S = \sum_{l \in I} \max_{i=1}^{k-1} S(a(i), l(i), i) \delta(l, l(i))$$

where $\delta(i, j)$ is the Kronecker delta (1 if $i = j$, 0 otherwise). The problem is: find a policy (mapping from process history, or alternately belief state and round number, to (agent, problem instance) pairs) that maximizes the expected value of S . This is the *sequential (adaptive)* decision making version of the problem.

We also consider, for computational complexity analysis, simpler *linear* settings of MaxScore, where the decision on

which agents to run on which problem instances (and in which order) is decided only once, in advance. A policy in this setting is simply a sequence of (agent, problem instance) pairs. The linear setting is the same as the *batch* setting (also called *non-adaptive*) commonly used in the research literature (Shperberg and Shimony 2017), except that if the runtimes are not both known and deterministic, one must specify the ordering so as to have a well-defined policy (some of the agents may not get to run at all, due to runtimes uncertainty).

Performance Models

Assuming that the distribution model is Markov, the MaxScore problem is a POMDP with states defined by the current maximum scores vector cR and the play time elapsed. As the number of rounds is not known in advance, we define the problem as an indefinite horizon POMDP with terminal states being those where the sum of runtimes exceeds T ¹. The transition probabilities in this POMDP are trivially (and deterministically) defined given the score and runtime achieved in the current round (which in turn are defined by distribution P). The reward function is 0 except for transitions into terminal states. In general, POMDPs are intractable (PSPACE complete even if the belief space is finite). The actual complexity of MaxScore depends on the setting (sequential vs. linear/batch), on the performance profiles distribution model P , and on the size of sets I and A .

A major issue is the performance profile (distribution model) P . Typically, exactly what scores and runtimes to expect is unknown, except by running the programs on the problem instances, which is too late to make the needed decision. However, we can run the programs on similar instances, collect statistics and learn a prediction model given instance features. Related work involved learning to predict the *expected score* of an agent in an unseen level (Stephenson and Renz 2017). However, as argued below, such information is insufficient for optimal choice: one may need to predict the *whole score distribution* (or equivalently, the expected improvement over each possible current score).

If the agent programs are effectively memoryless, i.e. attempt to solve the level from scratch each time they encounter it, then the order of the observed scores and runtimes is irrelevant. This behavior is reasonable for Angry Birds, as the game is effectively stochastic. Additionally (unlike search problems in most search domains), even if an agent knows the optimal play, it must still wait for the simulation to run its course, which usually takes on the order of one minute of real time per level attempt. Finally, the correlation coefficient between the actual score and time results measured over a few dozen instances was very small (≈ -0.015). Although this does not preclude a more complicated dependency between them, modeling score and time as independent is a reasonable approximation. We thus consider the agent scores and runtimes for a problem instance

¹Because the timeout T is known, one could use a finite horizon POMDP with T time-slices, but this would necessitate many dummy transitions, for time points where an agent is running and no decision is to be made, which is inefficient.

as i.i.d samples drawn from the distributions $P_S(a, l)$ and $P_T(a, l)$, respectively. If these distributions are known, the MaxScore problem is in fact an MDP, analyzed below.

3 Analysis: Known IID Case

We examine the computational complexity of some settings of the MaxScore problem. We begin with the fully deterministic case (scores and runtimes known in advance), and proceed to the independent case.

Complexity: Restricted Versions

We show that the MaxScore problem is NP hard even in the following extremely restricted cases:

1. Independent score distributions, deterministic runtimes, and only a single problem instance ($|I| = 1$).
2. Deterministic scores and runtimes, with $|A| = 1$ (but with $|I|$ unbounded).

We begin with the latter instance, as proving NP-hardness here is immediate through a straightforward mapping from the Knapsack problem. Simply map Knapsack item values to scores, item weights to runtimes, and the weight limit to T . Note that as the scores and runtimes are deterministic, in this case there is no difference between a linear setting, a batch setting, and a sequential setting of the problem.

With only one problem instance, we need to be more careful, but still get (see appendix for proof):

Theorem 1. *The linear setting of the MaxScore problem with independent score distributions, deterministic runtimes, and $|I| = 1$, is NP-hard.*

We believe that the complexity of *sequential* setting with the same restrictions is at least as hard as the linear setting, but have not proved it. Also note that the *linear* setting MaxScore problem with $|I| = 1$ is non-trivial even if we further restrict it to *unit runtimes*. For example, using a natural greedy scheme that picks the agent with the best expected score can be suboptimal. Consider having agents A, B, C, with time limit $T=2$. Suppose A always scores 100, B scores 101 with probability 0.99 and 0 otherwise, and C scores 200 with probability 0.001, and 99 with probability 0.999. A greedy scheme would first pick A, as it has the certain value 100, higher than the expected scores of B and C. In fact the optimal policy is to pick B and C (expected score just over 101, whereas anything containing A achieves less than 101). The computational complexity of this setting of MaxScore is, as far as we know, an open problem.

Approximation Algorithms

Using dynamic programming schemes it is possible to achieve a pseudo polynomial algorithm for the case of (known distribution) independent scores and runtimes, and $|I|$ bounded by a constant, using the following scheme. (The assumption that $|I|$ is bounded by a constant is reasonable for e.g. ABirds competitions, where $|I|$ is 4 or 8.) Additionally, we are assuming that score items (e.g. score for killing a pig in Angry Birds) and runtimes are integer valued, and that the time span and score items have a unary representation in

the input. The following dynamic programming value determination scheme (a variant of the Bellman equation) computes the optimal policy, and has a time complexity linear in the time span and the maximum score, and exponential in $|I|$.

Let $OPT(rT, cR)$ be the optimal solution value to the MaxScore problem with rT remaining time, and current maximum score vector $cR = \langle cR_1, \dots, cR_m \rangle$. The value determination recursive equation for $OPT(rT, cR)$ appears in Figure 1, where $R' = \langle cR_1, \dots, \max\{cR_l, r\}, \dots, cR_m \rangle$, and $sp(D)$ is the support of distribution D . The value of $OPT(T, \langle 0, \dots, 0 \rangle)$ is that of the optimal policy at the initial state.

If the score distributions are continuous, or have too many values, we can round them into bins, achieving a $(1 - \epsilon)$ -approximation to the optimal policy. Likewise discretizing the runtime distributions is possible, but here near-optimality is not guaranteed. Although the dynamic programming approximation scheme can be computed in pseudo-polynomial time, it is still too computationally demanding to be practical. We would thus like to use a greedy scheme in practice, and the one that comes to mind immediately is to use the agent that has the best expected score, as essentially done in (Stephenson and Renz 2017). A slightly better scheme is to take into account the runtime, and use the ratio of expected score over expected runtime. However, it is easy to show that these schemes are far from optimal (as verified by the empirical results).

For example, suppose we have only one problem instance and two agents. We have already achieved a score of 10,000, and have time for exactly one more run. Suppose the first agent always scores 10,000. The second agent scores 100,000 with probability 0.05, and otherwise fails and scores 0, thus its expected score is 5,000. The above greedy schemes would select the first agent and always get 10,000, while the optimal policy would obviously select the second agent, to possibly end up with 100,000 (with expected final score 10,500). An improved greedy scheme instead looks at the *expected improvement* to the score over the current score, i.e. the value:

$$E[S(a, l, i) - \max_{a \in A} \max_{j=1}^{i-1} S(a, l, j)] \quad (1)$$

rather than just the expected score $E[S(a, l, i)]$. (Consider $S(a, l, j)$ to be 0 if agent a was not selected in round j to be run on problem instance l .) Although the improved greedy scheme is suboptimal even for unit runtimes, in practice it does well (Section 4).

4 Experiments: Known IID

Quality/runtime tradeoffs for the independent model were examined for score and runtime distributions based on runs of ABirds algorithms on original Angry Birds game levels. Each algorithm was run 10 times on each level to obtain the empirical distributions, which were then treated as if they were the true distributions. Levels that caused issues with the agent's vision module were filtered out.

We applied the meta-agent to the open source versions of the following five existing agents: Naive, AngryBER, ihsev,

$$OPT(rT, cR) = \max_{\substack{a \in A \\ l \in I}} \left(\sum_{\substack{r' \in \\ sp(P_S(a,l))}} \sum_{\substack{l' \in \\ sp(P_T(a,l)) \\ \wedge t \leq rT}} OPT(rT - t, R') \times P_T(a, l)[t] \times P_S(a, l)[r] + \sum_{\substack{t \in \\ 1 \\ \wedge t > rT}}^m \sum_{i=1}^m cR_i \times P_T(a, l)[t] \right)$$

Figure 1: Optimal Solution to the MaxScore problem

Eagle’s Wing and planA, which competed in past AIBirds competitions. Note that these versions are not necessarily identical to the versions submitted for the competition. All tests were conducted on Windows 10 using a machine with Intel(R) i7-4700HQ 2.40GHz processor and 12 GB RAM. The evaluation process was performed using different numbers of levels and time budgets as described below:

```

1 for 1 to 50 do
2   for Every number of levels and time budget T
   configuration in {2, 3, 4} × {200, 400, 600, 800, 1000}
   do
3     Draw random levels uniformly from the level
     pool.
4     for 1 to 10,000 do
5       while Time budget was not exceeded do
6         Compute policy and choose a move
         using collected statistics as true
         distribution.
7         Execute the selected move (agent and
         level) by drawing score and time
         according to the statistics.

```

We also evaluated some of the algorithms in AIBirds competition settings: 8 levels with $T = 1800$ seconds.

The following optimization algorithms were compared:

1. **Dynamic Programing (optimal)**: compute the recurrence relation in Figure 1, using memoization of the results from all recursive calls.
2. **Binned Dynamic Programing(X, Y)**: same as the optimal solution, with scores rounded up to the next multiple of X , and times rounded to the next multiple of Y . We used this algorithmic scheme with $X \in \{1, 1000, 10000\}$ and $Y \in \{1, 10, 25\}$.
3. **Score Greedy**: choose the agent and level that maximize the expected score.
4. **Rate Greedy**: choose the agent and level that maximize the expected score divided by the expected time given that the time is less than or equal to the remaining time, multiplied by the probability that the time is less or equal to the remaining time.
5. **Improved Score/Rate Greedy**: same as the score/rate greedy except for considering expected score/rate improvement instead of expected score.
6. **Round Robin Score Greedy**: the algorithm used by the hyper-agent from (Stephenson and Renz 2017), which selects a level using round-robin, and chooses the agent that maximizes the expected score for that level, preferring agents not selected in previous attempts.

7. **Play Single Agent(A)**: select a level using round-robin, always with agent A . This scheme was evaluated for each of the 4 possible agents.
8. **Random**: draw a pair of level and agent uniformly.

The results appear in Figure 2. The scores in the plot are normalized to the highest score for each setting. For clarity, we show only a subset of the algorithms. In Play Single Agent policies, we show only the maximum value among them. The score greedy and improved score greedy were dominated by the rate greedy and improved rate greedy respectively, and are not shown. Finally, we showed the Binned Dynamic Programing(10000,10) as a sole representative of its category, since it achieved the best balance between runtime and score. The optimal policy did not always result in the best score, as the process is stochastic and thus exhibits measurement noise. The results indicate that the optimal policy is indeed the best in terms of scores. However, the binned version and the improved rate greedy achieved near-optimal results. When considering runtime and space usage, the optimal solution could not solve instances greater than 4 levels with $T = 400$ time budget. The binned version was able to solve all instances, with a maximal overhead of 13.7 seconds and a maximal memory usage of 52 MB across all instances with 4 levels or less. However, it required an average of 335 seconds and 270 MB of memory to handle the full competition setting (8 levels with 1800 seconds time budget). All the other algorithms ran in negligible time (several milliseconds) and memory. This makes the improved rate greedy the best algorithm in terms of balance between score and resources. The large gap in scores between schemes that selected instances either randomly or in round-robin fashion and those that attempted to optimize instance selection (improved greedy and dynamic programming) suggest that the multi-instance setting is very different from the single instance setting, and that a good instance selection scheme is crucial. The apparent increase in the gap as the number of instances grows further supports this observation (Figure 2). Note that in the experiments we did not include the optimization runtime in the total available runtime T ; but in a competition it must be. Also, the simulation runtime of the agents averaged roughly 90 seconds, so 1800 seconds consists of about 20 rounds (agent runs).

5 Unknown Distributions

A major point of competitions like AIBirds (as well as other competitions, such as IPC, SAT-solving, etc.) is that they are done with *previously unseen* problem instances, so the score and runtime distributions are unknown. The latter issue then becomes a learning problem, which can be modeled by treating the agent performance quality as hidden random

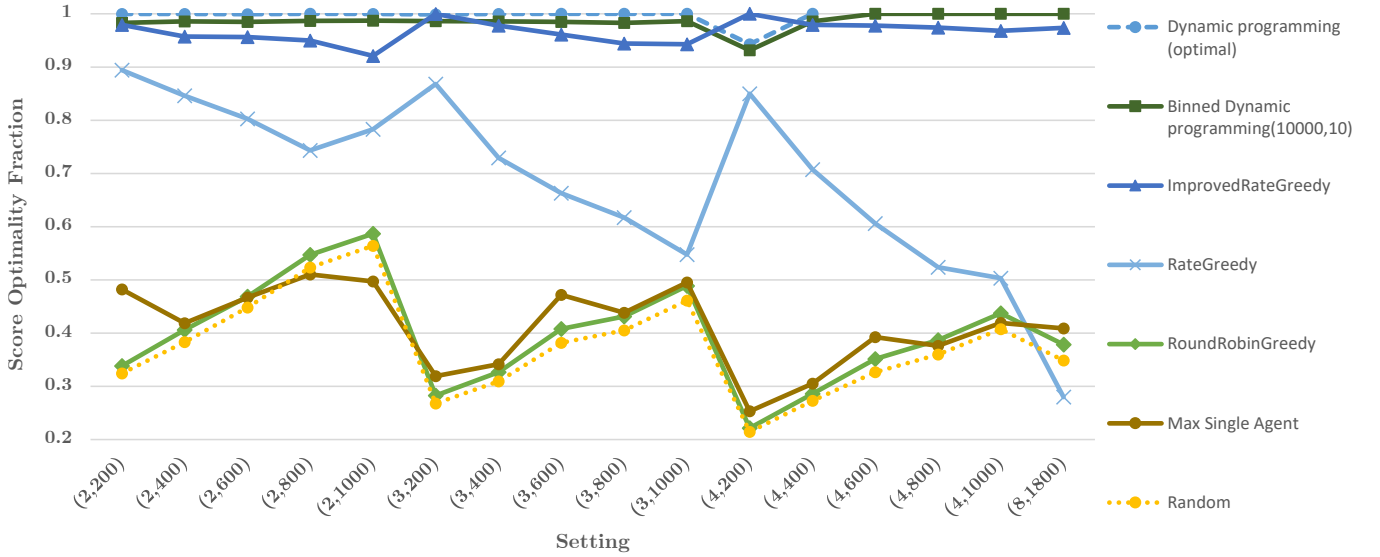


Figure 2: Optimization Algorithms Score Evaluation

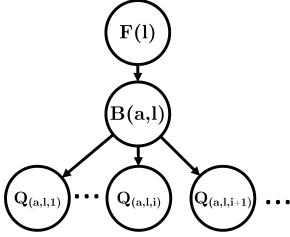


Figure 3: Dependency model (Bayes net fragment)

variables, with some assumed prior distributions based on observing similar problem instances. We adopt a naive learning scheme (described briefly below for independent scores and runtimes for simplicity, as implemented for the Albirds meta-agent).

Unknown IID Score and Runtime

In our naive learning scheme, we are assuming that agent performance profiles of a previously unseen problem instance are (almost) equal to their performance profile on some problem instance(s) for which performance statistics were already collected. Hence, we are essentially taking a case-based reasoning (CBR) approach to predicting the performance profile for an unknown instance. However, we do not assume knowledge of *which* of the previously seen performance statistics fits the current instance.

Therefore, if we need to predict an agent performance $Q(a, l, i)$ for round i of the current problem instance l (which we model as the score and runtime iid $P_S(a, l)$ and $P_T(a, l)$) it makes sense to condition on the (unknown) agent performance profile $Q(a, l)$ for the current problem instance. Essentially, what we need is a mapping from features to agent performance (i.e. *distributions* over score and runtime distributions), as we do in this paper.

The mapping we adopt here is simply a smoothed version of the performance profiles of the K most similar instances. (Smoothing is done in by assigning scores and runtimes into bins.) That is, for every problem instance l in the training set, and every agent program a , collect and store the score and runtime distributions estimate as $Q(a, l)$ indexed by the feature values (vector $F(l)$) for problem instance l .

When a new problem instance l is encountered online, compute its feature vector values $F(l)$, and find the K most similar instances l_1, \dots, l_K according to some appropriately defined similarity measure $s(l_i, l)$. Now we assume that the agent performance for instance l has a distribution over performance profiles, and that it is equal to the (smoothed version) of its performance profile for some instance l_i , with probability proportional to $s(l_i, l)$. That is, denote by $B(a, l)$ a K -valued random variable, with integer values denoting the respective $Q(a, l_i)$ profile. Then we have: $P(B(a, l) = i) = \frac{s(l_i, l)}{\sum_{j=1}^K s(l_j, l)}$.

The performance profiles describe both score and runtime distributions, and additionally we assumed that these are drawn i.i.d. given the value of $B(a, l)$. The distribution model topology is summarized in Figure 3, for each problem instance l (shown for one agent program). We have an observable feature vector variable $F(l)$. Belief updating for this conditionally i.i.d. model is straightforward, as this is a naive Bayes model.

In the conditionally i.i.d. model, since $B(a, l)$ is unobservable, but its current distribution (thus belief state) changes given new observations of $T(a, l, i)$ and $S(a, l, i)$, we now have a POMDP that we cannot hope solve optimally, especially in real time. Instead, we can solve an MDP where the $T(a, l, i)$ and $S(a, l, i)$ are assumed to be i.i.d. as before, but based on the current belief state of $B(a, l)$.

That is, we can do the belief updating given the new observed scores and runtimes, but in the policy computation

act as if future updates are not observed. Then one can re-compute the MDP policy after each observation and belief update. However, the MDP solution was also quite computationally intensive, and re-computation makes it even more so. As the improved greedy scheme performed almost as well as the MDP solution w.r.t. optimality, we no longer considered using the MDP solution, for practical reasons.

Obtaining a reasonable $s(l, l')$ is a learning problem, which was tackled by normalizing all feature values to $[0, 1]$, and taking the inverse of the Euclidean distance as the similarity. Despite the naive nature of our method for defining prior probability of unobserved levels, our algorithms, provided with such priors, showed a major improvement over existing methods according to the results presented below.

Experiments: Unknown IID

We normalized each level’s score by $maxScore_l$, an upper bound on achievable score in each level, that can be computed using the features. We used the following subset of features, described in (Stephenson and Renz 2017; Tziortziotis, Papagiannis, and Blekas 2016): #Blocks, targetWidth, targetHeight, closestObjDist, farthestObjDist, density, #Objects, iceObjects, woodObjects, stoneObjects, #Pigs, helmetPigs, noHelmetPigs, #Birds, #RedBirds, #YellowBirds, #BlueBirds, #BlackBirds, #WhiteBirds, varietyOfBirds, feasibleObjects, feasiblePigs, roundObjectsNotPigs, icedTerritory, woodenTerritory, stonedTerritory, averagePigsInBlocks, blocksWithPigs and #TNTs.

We tested the optimization algorithms using the same process as in Section 4, where the algorithms had to rely on the predicted distribution based on the naive learning scheme. We incorporated the resulting distribution in the following algorithms: (1) the improved rate greedy defined in Section 4, using the distribution of distributions with K neighbors (denoted by $IRG(K)$) without belief updating; (2) a binned version of the improved rate greedy with Bayesian belief updating, using a zero value bin and 10 additional bins uniform in $(i \cdot maxScore_l, i + 0.1 \cdot maxScore_l]$, $0 \leq i \leq 9$ (denoted as BIRG). We used $K = 128$ in the learning process for both algorithms as a default value. We also tested IRG with other K values, specified in parenthesis.

Table 1 shows the solution quality achieved by the different algorithms, relative to the quality achieved by an improved rate greedy algorithm acting on *known* distributions (denoted "omniscient"). The projected standard deviation σ' of the results was at most ± 0.013 , small enough to maintain the performance ordering between the algorithms as presented below.²

In most cases, IRG (using *unknown* distributions) shows a major improvement over the baseline methods: choosing agent and level at random; choosing the post-facto best performing agent (denoted MSA, with the first letter of agent achieved that score in parenthesis); and the round robin greedy algorithm (denoted as RRG) with *known* expectation. BIRG (using *unknown* distributions), further improved the

²The projection σ' was based on the standard deviation σ measured over averages of sets of 35 random runs. To project to the 1000 runs per instance actually used we have $\sigma' = \frac{\sigma}{\sqrt{\frac{1000}{35}}}$.

results, achieving an average solution quality of 0.75 despite using a naive learning scheme. Note that when BIRG had sufficient time to perform updates (600 seconds and above) it achieved the best score out of all tested algorithms. We also tracked the improvement due to Bayesian updates by comparing the Wasserstein distance (also known as "earth mover’s distance" (EMD)) between the predicted and true distributions. As expected, results improved as the above EMD decreased, which also explains the improved results for BIRG in the longer sequences. With Bayesian updating, the EMD has improved from the beginning to the end of each test setting, with an average improvement of 17.2% (not shown).

In IRG performance seems to improve with increased K , up to a point where this trend peters off and even reverses. We believe that including too few cases is insufficient to predict the performance profile well, but too many cases leads to overfitting. This phenomenon does not occur in BIRG due to the Bayesian updating which quickly disregards irrelevant cases, and improves monotonically with K .

Experiments: Angry Birds Game

After experimenting using data collected from the game, we implemented a full version of the meta-agent, which interfaces with the AIBirds server and plays the actual game. Our meta-agent implementation starts by collecting information on all levels using the provided vision module. Based on this, the meta-agent constructs an objects-tree for each level, extracts features from the objects-trees, and predicts a performance profile for each level and agent pair. Then, the meta-agent applies the BIRG scheme to select a pair of agent and a level to play. The meta-agent sends the selection to the server and observes the results of the run. The observations are used for belief updating. The select-and-play process repeats until the time limit is reached.

Our evaluation was based on past competition levels (between 2014 and 2016), a total of 72 levels (8 at a time with a 30 minutes time budget). The results are shown in table 2. The improved-greedy based meta-agent achieved an average score of 441,752, compared to PlanA (357,468), ihsev (321,280), AngryBER (303,166) and Eagle’s Wing (323,099). Note that the above 4 agents were the ones actually used by the meta-agent, and all of them contributed to its score. The hyper-agent of (Stephenson and Renz 2017) achieved an impressive average score of 424,740; However, the authors, which are the organizers of the AI-Birds competition, had access to a total of 8 agents as opposed to our 4 open-sourced agents. We strongly believe that using the full set of 8 agents would have further improved the performance of our algorithm, as in auxiliary experiments (not shown) our scheme was relatively robust to adding agents (including dummy, useless agents). This belief is further supported by Stephenson and Renz (2017) which showed that each individual agent actually contributed to the score of the hyper-agent, meaning that the agents to which we had no access were actually *not* useless.

Algorithm	Number Of Levels															AVG
	2					3					4					
	200	400	600	800	1000	200	400	600	800	1000	200	400	600	800	1000	
IRG	0.71	0.72	0.70	0.75	0.78	0.66	0.71	0.63	0.67	0.69	0.55	0.63	0.62	0.63	0.67	0.67
BIRG	0.69	0.76	0.81	0.84	0.85	0.67	0.75	0.73	0.76	0.83	0.54	0.68	0.77	0.75	0.79	0.75
IRG (1)	0.60	0.60	0.64	0.63	0.59	0.57	0.54	0.52	0.52	0.54	0.54	0.53	0.46	0.44	0.47	0.55
IRG (50)	0.70	0.73	0.74	0.75	0.78	0.72	0.73	0.65	0.70	0.72	0.57	0.69	0.64	0.70	0.71	0.70
IRG (100)	0.70	0.70	0.71	0.75	0.76	0.68	0.69	0.64	0.67	0.70	0.59	0.65	0.62	0.63	0.68	0.68
Random	0.38	0.50	0.56	0.64	0.68	0.31	0.40	0.45	0.52	0.59	0.24	0.34	0.36	0.47	0.50	0.46
RRG	0.45	0.56	0.63	0.71	0.75	0.32	0.48	0.52	0.60	0.66	0.28	0.40	0.42	0.55	0.57	0.53
MSA	0.53(p)	0.80(e)	0.66(p)	0.72(i)	0.75(i)	0.48(p)	0.57(i)	0.60(p)	0.66(p)	0.65(i)	0.38(e)	0.60(p)	0.61(p)	0.64(p)	0.63(i)	0.62

Table 1: Solution qualities as a fraction of the solution quality obtained by the "omnicient" improved rate greedy scheme

Round	Naive	PlanA	ihsev	A-BER	E-Wing	Hyper	Meta
Q 2014	187,180	314,540	109,920	188,710	282,110	332,270	418,210
S 2014	400,980	541,220	439,520	429,680	442,800	524,400	602,050
F 2014	209,130	193,110	257,410	90,110	250,970	338,330	231,480
Q 2015	68,020	316,850	163,790	367,000	346,760	351,300	372,260
S 2015	145,910	288,870	166,750	143,270	299,220	375,670	369,130
F 2015	131,660	452,860	458,030	392,420	191,970	483,610	490,050
Q 2016	251,080	313,440	444,560	310,600	252,100	336,840	426,570
S 2016	436,870	372,330	562,820	445,030	420,170	610,280	593,020
F 2016	390,050	423,990	288,720	361,670	421,790	469,960	471,380
Average	246,764	357,468	321,280	303,166	323,099	424,740	441,572

Table 2: Actual game results using past competitions setting

6 Discussion

In this paper we defined the MaxScore optimization problem, analyzed its computational complexity (NP-hard even under extreme restrictions), and suggested approximation algorithms for known independent distributions. In practice, based on empirical evaluation on AI birds, it turns out that a greedy algorithm based on expected improvement is near-optimal. Despite the latter having no theoretical guarantees, it currently seems to be the only viable alternative for real-time computation. Applying these results to unknown distributions requires learning performance profiles given problem instance features. A naive learning scheme applicable to the AIbirds application was proposed. This results in imperfect predicted distributions, which degrades the meta-reasoning results. Nevertheless, the greedy algorithm is still the better option, especially if the distribution model is updated using scores and runtimes observed during the run.

The MaxScore problem is closely related to *algorithm selection*, as originally defined by Rice in 1976 (Rice 1976). *Algorithm portfolios* (Gomes and Selman 2001; Huberman, Lukose, and Hogg 1997) are a natural and popular extension of the idea of algorithm selection. Such techniques are based on minimizing risk in economics. This approach defines a collection of algorithms (a portfolio) and establish a resource allocation to the algorithms in the portfolio in order to solve a given problem instance (instead of choosing a single algorithm for a given problem instance). This field has been studied extensively in the last decades, including works on different computational settings (parallel, sequential or in-between), many applications with outstanding results (Xu et al. 2008; Hoos, Lindauer, and Schaub 2014; Kadioglu et al. 2010; 2011) and even meta-level techniques for choosing a selector (Lindauer et al. 2015). Most common settings of algorithm portfolios focus on finding a solution to a single given problem instance. Our setting generalizes the meta-level decision problem solved in algorithm portfolios

to choosing which problem instance to work on, as well as selecting algorithms to use at any given time. A paper on dynamic restart policies (Kautz et al. 2002) proposes an optimal restart scheme in a decision-theoretic sense, similar to that defined in our paper, and with a scheme for learning a runtime distribution. Since our setting allows non-binary scores, where it is important to get a good score on a problem instance, rather than just solve it, the optimization scheme used in the restart policies paper is not directly applicable here. The scheme they use to learn runtime distributions may be applicable to our setting, but must be extended to predict score distributions as well before it can be used here. Maximizing the *number* of instances solved is also mentioned in (Kautz et al. 2002), but their instances are drawn randomly and independently, so there seems to be no allowance for the capability of choosing to return to a previously run instance as in our setting, in addition to there being no notion of instance score in (Kautz et al. 2002).

The MaxScore problem is also loosely related to multi-armed bandit (MAB) problems (Auer, Cesa-Bianchi, and Fischer 2002). Much of the related work on MABs does not assume a known distribution, or even a distribution over distributions as done in this paper. Rather, bounds on regret are analyzed, both asymptotic and finite. However, the fact that the reward in MaxScore is the maximum rather than the sum makes it unclear how such techniques might carry over. Additionally, in the motivating application of AIbirds, the number of rounds is small, further complicating such attempts. In fact, if we tried to apply an MAB scheme directly, we would get a random selection of problem instance, against which we did compare in the Angry Birds domain (random did poorly, as expected).

A significant part of the research on algorithm portfolios and multi-armed bandits focuses on learning issues. E.g. in (Kotthoff 2016), the focus is on analyzing problem features and applying different varieties of machine learning techniques in order to find scheduling policies for the portfolios. In this paper we achieved good results despite using a rather naive learning scheme to obtain a mapping from features to score and runtime distributions. Note that the relative performances in Table 1 still leave much room for improvement by better predicting the distributions: these performance figures are still well below the 1.0 value obtained by the an "omnicient" rate-greedy scheme that has access to the true distributions. Introducing better learning schemes for better prediction of the distributions should thus result in better performance.

Another issue for future work is learning the distribution models with time-score and inter-round dependencies, thus extending MaxScore solutions to more general settings of algorithm portfolios over optimization problems. Fully testing such generalized scenarios would require changing the rules of the competitions to maximizing total score over a global time limit, rather than the current setting where the time limits are per-instance.

Acknowledgements

Supported by ISF grant 417/13 and the BGU Frankel Center. Meta-agent implemented by: Lior Schachter, Dor Bareket, Ori Zviran.

Appendix: Proof of Theorem 1

Theorem 1. *The linear setting of the MaxScore problem with independent score distributions, deterministic runtimes, and $|I| = 1$, is NP-hard.*

Proof: by reduction from the optimization version of knapsack ((Garey and Johnson 1979), problem number [MP9]), re-stated below. Given a set of items $\mathcal{S} = \{s_1, \dots, s_n\}$, each with a positive integer weight w_i and a positive integer value v_i , a weight limit W , find a sub-multiset S of \mathcal{S} with a maximal total value, subject to: total weight of S at most W .

In the reduction, each agent represents an item in the Knapsack problem, where $P_T(a_i, l) = [1 : w_i]$ and $P_S(a_i, l) = [\varepsilon : v_i, 1 - \varepsilon : 0]$. As this is a simple one-to-one mapping, we abuse the notation and treat the agents as if they are actually the respective elements from \mathcal{S} in the Knapsack problem. In the MaxScore problem, let:

$$T = W, \quad H = \max_{s_i \in \mathcal{S}} v_i, \quad M = \frac{W}{\min_{s_i \in \mathcal{S}} w_i}, \quad \varepsilon = \frac{1}{M^2 H + 1}$$

Let S be a candidate solution to the MaxScore problem, with $m = |S| \leq n$. Assume w.l.o.g. that $S = \{s_1, \dots, s_m\}$ and that the items are sorted in non-descending order of values v_i . Denote by $P(S)$ expected value from selecting the items in the sequence S as a policy. Then:

$$P(S) = \sum_{i=1}^m v_i \varepsilon (1 - \varepsilon)^{m-i} \leq \sum_{i=1}^m v_i \varepsilon$$

On the other hand, we have:

$$\begin{aligned} P(S) &= \sum_{i=1}^m v_i \varepsilon (1 - \varepsilon)^{m-i} \geq \sum_{i=1}^m v_i \varepsilon (1 - \varepsilon)^m \\ &\geq \sum_{i=1}^m v_i \varepsilon (1 - \varepsilon)^M \end{aligned}$$

From Bernoulli's inequality, we have:

$$(1 - \varepsilon)^M \geq 1 - M\varepsilon = 1 - \frac{M}{M^2 H + 1} > 1 - \frac{1}{MH}$$

Therefore:

$$\begin{aligned} P(S) &> \sum_{i=1}^m v_i \varepsilon (1 - \frac{1}{MH}) = \sum_{i=1}^m v_i \varepsilon - \frac{\sum_{i=1}^m v_i \varepsilon}{MH} \\ &\geq \sum_{i=1}^m v_i \varepsilon - \varepsilon = \varepsilon (\sum_{i=1}^m v_i - 1) \end{aligned}$$

Now let S be an *optimal* solution to the MaxScore problem. Since S satisfies the time constraint, we have $\sum_{i=1}^m w_i \leq T = W$, so S satisfies the weight constraint in the Knapsack problem and is thus a solution therein. Assume in contradiction that there exists a legal solution S' to Knapsack s.t. $\sum_{s_i \in S'} v_i > \sum_{s_i \in S} v_i$. Since the values of the items in knapsack are integers, we know that $\sum_{s_i \in S'} v_i \geq (\sum_{s_i \in S} v_i) + 1$. Thus, as $|S'| \leq \frac{1}{\varepsilon}$:

$$P(S') > \varepsilon (\sum_{s_i \in S'} v_i - 1) \geq \varepsilon (\sum_{s_i \in S} v_i) \geq P(S)$$

As S' satisfies the timing budget in the MaxScore problem, it is a solution better than S , a contradiction. So S is also an optimal solution to the Knapsack problem. \square

References

- Auer, P.; Cesa-Bianchi, N.; and Fischer, P. 2002. Finite-time analysis of the Multiarmed bandit problem. *Mach. Learn.* 47:235–256.
- Garey, M. R., and Johnson, D. S. 1979. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman.
- Gomes, C. P., and Selman, B. 2001. Algorithm portfolios. *Artif. Intell.* 126(1-2):43–62.
- Hoos, H.; Lindauer, M. T.; and Schaub, T. 2014. claspfolio 2: Advances in algorithm selection for answer set programming. *TPLP* 14(4-5):569–585.
- Huberman, B. A.; Lukose, R. M.; and Hogg, T. 1997. An economics approach to hard computational problems. *Science* 275(5296):51–54.
- Kadioglu, S.; Malitsky, Y.; Sellmann, M.; and Tierney, K. 2010. ISAC –Instance-Specific Algorithm Configuration. In *Proceedings of the 2010 Conference on ECAI 2010: 19th European Conference on Artificial Intelligence*, 751–756. Amsterdam, The Netherlands, The Netherlands: IOS Press.
- Kadioglu, S.; Malitsky, Y.; Sabharwal, A.; Samulowitz, H.; and Sellmann, M. 2011. Algorithm selection and scheduling. In *CP (LNCS6876)*, 454–469.
- Kautz, H. A.; Horvitz, E.; Ruan, Y.; Gomes, C. P.; and Selman, B. 2002. Dynamic restart policies. In Dechter, R.; Kearns, M. J.; and Sutton, R. S., eds., *Proceedings of the Eighteenth National Conference on Artificial Intelligence and Fourteenth Conference on Innovative Applications of Artificial Intelligence, July 28 - August 1, 2002, Edmonton, Alberta, Canada.*, 674–681. AAAI Press / The MIT Press.
- Kotthoff, L. 2016. Algorithm selection for combinatorial search problems: A survey. In Bessiere, C.; Raedt, L. D.; Kotthoff, L.; Nijssen, S.; O'Sullivan, B.; and Pedreschi, D.,

eds., *Data Mining and Constraint Programming - Foundations of a Cross-Disciplinary Approach*, volume 10101 of *Lecture Notes in Computer Science*. Springer. 149–190.

Lindauer, M. T.; Hoos, H. H.; Hutter, F.; and Schaub, T. 2015. Autofolio: An automatically configured algorithm selector. *J. Artif. Intell. Res.* 53:745–778.

Rice, J. R. 1976. The algorithm selection problem. *Advances in Computers* 15:65 – 118.

Shperberg, S. S., and Shimony, S. E. 2017. Some properties of batch value of information in the selection problem. *J. Artif. Intell. Res. (JAIR)* 58:777–796.

Stephenson, M., and Renz, J. 2017. Creating a hyper-agent for solving angry birds levels. In Magerko, B., and Rowe, J. P., eds., *Proceedings of the Thirteenth AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment (AIIDE-17), October 5-9, 2017, Snowbird, Little Cottonwood Canyon, Utah, USA.*, 234–240. AAAI Press.

Tziortziotis, N.; Papagiannis, G.; and Blekas, K. 2016. A bayesian ensemble regression framework on the angry birds game. *IEEE Trans. Comput. Intellig. and AI in Games* 8(2):104–115.

Xu, L.; Hutter, F.; Hoos, H. H.; and Leyton-Brown, K. 2008. SATzilla: Portfolio-based algorithm selection for SAT. *J. Artif. Intell. Res.* 32:565–606.

Zilberstein, S., and Russell, S. J. 1996. Optimal composition of real-time systems. *Artif. Intell.* 82(1-2):181–213.