

Score vs. Winrate in Score-Based Games: which Reward for Reinforcement Learning?

Luca Pasqualini
University of Siena, Italy
pasqualini@diism.unisi.it

Gianluca Amato
University of Chieti-Pescara, Italy
gianluca.amato@unich.it

Marco Fantozzi
Italy
marco.fantozzi@gmail.com

Rosa Gini
ARS Toscana, Italy
rosa.gini@ars.toscana.it

Alessandro Marchetti¹
University of Chieti-Pescara, Italy
alessandro.marchetti@unicampus.it

Carlo Metta
ISTI-CNR Pisa, Italy
carlo.metta@isti.cnr.it

Francesco Morandin
University of Parma, Italy
francesco.morandin@unipr.it

Maurizio Parton
University of Chieti-Pescara, Italy
maurizio.parton@unich.it

Abstract—In the last years, the DeepMind algorithm AlphaZero has become the state of the art to efficiently tackle perfect information two-player zero-sum games with a win/lose outcome. However, when the win/lose outcome is decided by a final score difference, AlphaZero may play score-suboptimal moves because all winning final positions are equivalent from the win/lose outcome perspective. This can be an issue, for instance when used for teaching, or when trying to understand whether there is a better move. Moreover, there is the theoretical quest for the perfect game. A naive approach would be training an AlphaZero-like agent to predict score differences instead of win/lose outcomes. Since the game of Go is deterministic, this should as well produce an outcome-optimal play. However, it is a folklore belief that “this does not work”.

In this paper, we first provide empirical evidence for this belief. We then give a theoretical interpretation of this suboptimality in general perfect information two-player zero-sum game where the complexity of a game like Go is replaced by the randomness of the environment. We show that an outcome-optimal policy has a different preference for uncertainty when it is winning or losing. In particular, when in a losing state, an outcome-optimal agent chooses actions leading to a higher score variance. We then posit that when approximation is involved, a deterministic game behaves like a nondeterministic game, where the score variance is modeled by how uncertain the position is. We validate this hypothesis in AlphaZero-like software with a human expert.

Index Terms—reinforcement learning, AlphaZero-like algorithms, score-based games

I. INTRODUCTION

The game of Go has been a landmark challenge for AI research since its very beginning. It is very suited to AI, with the importance of patterns and the need for deep exploration, and very tough to actually solve, with its whole-board features and subtle interdependencies of local situations. Nowadays, AI has reached a superhuman level in the game of Go with the well-known DeepMind algorithm AlphaGo Zero [1] (AGZ), a zero-knowledge evolution of AlphaGo [2] (AG). A more general version of the algorithm, AlphaZero [3] (AZ), might

even be able to tackle efficiently the whole class of perfect information two-player zero-sum games.

In perfect information two-player zero-sum games where the win/lose outcome is given by a final score difference and no other rewards during the game, maximizing this score difference is still an open and important question, see the detailed discussion in [4, Introduction]. In fact, AZ-like algorithms play suboptimal moves in the endgame, see for instance [5, moves 210 and 214, page 252]. The open-source clean room implementation of AGZ known as Leela Zero [6] (LZ) is also known to play suboptimal moves, see Section 4.4 in [7].

This phenomenon is rooted in the win/lose reward in the reinforcement learning (RL) pipeline. Giving a reward of 1 (win) or 0 (lose) at the end of the game means that agents maximize just the winrate (that is, the expected win/lose outcome) and could play sub-optimally with respect to the expected score difference. It is a folklore belief that replicating the AZ pipeline using the score difference as a primary target, instead of the Boolean outcome, is unsuccessful. A qualitative argument is that the score difference is unlikely to be a successful reward because without knowledge of the win/lose outcome, the agent will give the same importance to each score point. But when the score difference is close to zero, a single point may change the outcome of the game, thus inducing instability in the agent strength. Note that this only happens because the perfect play is out of reach: if the training was to reach a perfect play, win/lose and score difference rewards would produce agents that are equivalent from the win/lose point of view, because the game of Go is deterministic.

As a matter of fact, there are at least two different RL approaches that proved somewhat successful in maximizing the score: KataGo [8] and SAI [7]. KataGo does include score estimation, but only as a secondary target: the value to be maximized is a linear combination of winrate and expectation of a nonlinear function of the score difference, not the score difference itself. In SAI, the winrate is modeled as a two parameters family of sigmoids $\sigma_{\alpha,\beta}$: while α can be seen as the final score difference, α and β are learned indirectly by training $\sigma_{\alpha,\beta}$ against the classical Boolean reward.

Funded by INdAM groups GNSAGA, GNCS and GNAMPA. Computational resources provided by CLAI lab of Chieti-Pescara.

¹Alessandro Marchetti is a PhD student enrolled in the National PhD in Artificial Intelligence, XXXVII cycle, course on Health and life sciences, organized by Campus Bio Medico University of Rome.

Still, humans do use score estimations instead of win/lose outcome estimations while playing score-based games. The research question we address in this paper is to what extent a win/lose-based optimal play can be achieved by a score-based optimal play. In simpler words: if I play to maximize the score in a win/lose game, do I lose more often?

This very question has different answers in different cases that are better-detailed in II. In a deterministic game sufficiently simple to allow for exact computation of the optimal play, maximizing the expected score difference yields also optimal win/lose-based play. However, if the game is nondeterministic, maximizing the expected score difference does not always give optimality from the point of view of the winrate (even at optimal play). We elaborate on this unexpected behavior and show that score variance is the key to understanding this suboptimality, see IV. In a game that is deterministic but so complex that optimal play can only be roughly approximated, we support the mentioned folklore belief by training LZ-Score (LZS), an instance of LZ on the 9×9 board, using score difference as a target. We show that the training is successful, see III-A, but converges prematurely to a player weaker than a corresponding AGZ-like player, see III-B. Finally, in V, we detail the creation of LZS.

II. CONCEPTUAL FRAMEWORK

We consider two fully competitive agents playing in a perfect information finite sequential game, whose win/lose outcome (hereafter, simply *outcome*) is decided solely by a final score difference. The game can be deterministic, with no chance involved like the game of Go, or nondeterministic with chance events, like backgammon. The agents can be score-based or outcome-based, according to whether they minimize/maximize the expected score difference or the expected win/lose outcome (hereafter, simply *winrate*). For each agent, we call *minimax optimal* any (it is not necessarily unique) theoretical optimal policy for that agent’s target. A backward induction argument shows that, in deterministic games, any minimax-optimal policy for the score-based agent would be also minimax optimal for the outcome-based agent, that is, to win in a win/lose game one can actually maximize the score.

However, in deterministic but very complex games where optimal policy is out of reach, maximizing the score may still lead to a suboptimal play from the winrate perspective. We empirically validate this claim by training LZ-Score (LZS), an AZ-like score-based agent, and by evaluating its performance with a human-in-the-loop and a quantitative approach, see III.

One interpretation of this incoherence between score-based and outcome-based play is that in complex deterministic games, the value approximation is far from being optimal, and therefore the agent plays in partial ignorance. Partial ignorance can be modeled with the fact that, to the agent’s eyes, the game is nondeterministic, and in this case, a minimax-optimal policy for the score-based agents is no longer guaranteed to be minimax optimal for the outcome-based agents. In order to understand this phenomenon, we build a family of nondeterministic games where the minimax-optimal policy for

the score-based agents is not optimal for the outcome-based agents, and elaborate that score variance is the key statistics to understand this behavior, see IV.

III. EVALUATION

A. Qualitative evaluation against a human player

Fifteen games were played between the best LZS network and Carlo Metta, a strong amateur player¹. Ten games were played with 400 visits, that is, the same setting as games in the quantitative evaluation described in III-B, while five games were played with 20,000 visits, to test a stronger player.

A thorough analysis of such games shows that training has been successful in producing a consistent player, which, however, exhibits some unusual characteristics when compared to other artificial agents. The match ended with a score of 14-1 in favor of the human player: although LZS found itself in a position of clear advantage several times, it was only able to win one game, one of those with 20,000 visits. LZS showed some peculiar and not always desirable features. LZS certainly has a direct and aggressive style. It does not seem to admit sacrificing a few stones for better final results, e.g. move 17 in game 6, nor to foresee sacrifice on the opponent’s side, see e.g. move 16 in game 9. This is clearly in contrast with the flexibility shown by other artificial agents.

Another striking situation occurred several times: when in balanced positions, LZS attempted to further increase the score difference, rather than settling for a narrow victory, in such an aggressive and self-delusional way that it resulted in an inevitable defeat. See for instance move 21 in game 8. It may be argued that this phenomenon was a direct effect of the LZS training scheme.

B. Quantitative evaluation against SAI

To estimate the Elo rating for LZS networks promoted in training, we compared one every four of them against a calibration panel of 32 SAI networks, whose Elo ranged from 683 to 3501, with 400 visits. Elo ratings were then estimated by maximum likelihood, as in Rémi Coulom’s Bayes Elo [9].

In Fig. 1 we compare LZS with run 1 of SAI, trained for a comparable number of self-plays, with the same network architecture and weights. see [7]. This comparison shows that during training LZS had consistently lower values of Elo. The last increase in the size of the network, at 900,000 games, did not yield any relevant improvement in the following 400,000 games, thus confirming that LZS was converging prematurely to a weaker player. After 200,000 games the two curves are approximately parallel, with an Elo difference of around 1,500 points. This is a remarkable difference in strength: the interpretation of the Elo formula means that, at the same level of training, LZS would win against SAI with probability $10^{-(1500/400)}$, i.e. less than once every 5000 games.

¹Player profile on EGD, the European Go Database: https://www.europeangodatabase.eu/EGD/Player_Card.php?&key=14713996.

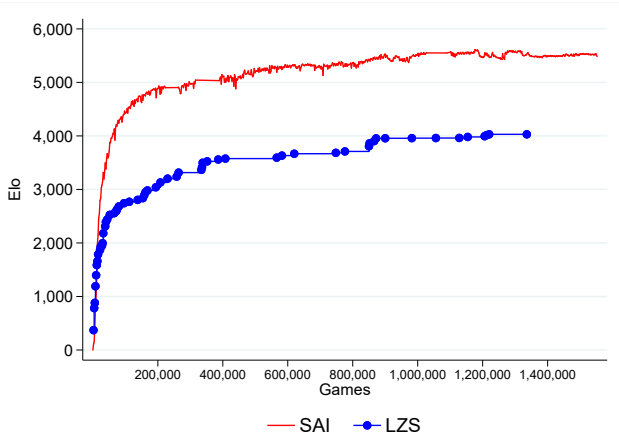


Fig. 1. Elo ratings of all the promoted networks from LZS and from run 1 of SAI 9×9, as a function of the cumulative number of self-play games.

IV. SUBOPTIMALITY OF SCORE-OPTIMAL POLICY IN A WIN/LOSE OUTCOME SETTING

In this section we give a theoretical interpretation of the suboptimality of maximizing the score, considering two much simpler settings in which the complexity is replaced by the randomness. In both settings, we show that suboptimality is a common behavior, and it is associated with the fact that the outcome-optimal policy has a different preference for uncertainty when it is winning or losing. When in a losing state, an outcome-optimal agent chooses actions leading to a higher score variance. According to strong human players, this is in fact what human players do too.

In the rest of this section, we consider a single agent instead of two fully-competitive agents. By the analysis in [10, Section 4.3], this is not a restrictive hypothesis.

A. Multi-armed bandit

We start by proving that in the multi-armed bandit case, score-based optimality does not imply outcome-based optimality.

Proposition 1. *Consider a 2-armed bandit where each action $i \in \{1, 2\}$ yields a random score reward $r_{\text{score}}(i) \in \mathbb{R}$, with Gaussian distribution $\mathcal{N}(\mu_i, \sigma_i^2)$ and an outcome reward $r_{\text{outcome}}(i) \in \{0, 1\}$ defined as $r_{\text{outcome}}(i) := \mathbb{1}_{r_{\text{score}}(i) > 0}$. Then there are choices of the parameters for which:*

$$\mathbb{E}[r_{\text{score}}(1)] > \mathbb{E}[r_{\text{score}}(2)], \mathbb{E}[r_{\text{outcome}}(1)] < \mathbb{E}[r_{\text{outcome}}(2)] \quad (1)$$

Indeed this happens if and only if $\mu_1 > \mu_2 > 0$ and $\sigma_1 > \sigma_2 \frac{\mu_1}{\mu_2}$, or $\mu_2 < \mu_1 < 0$ and $\sigma_1 < \sigma_2 \frac{\mu_1}{\mu_2}$.

Proof. Compute $\mathbb{E}[r_{\text{outcome}}(i)] = P(r_{\text{score}}(i) > 0) = \Phi(\frac{\mu_i}{\sigma_i})$, where Φ denotes the cumulative distribution function. By the monotonicity of Φ and assuming (1) we get $\mu_1 > \mu_2$ and $\frac{\mu_1}{\sigma_1} < \frac{\mu_2}{\sigma_2}$. This implies that μ_1 and μ_2 must have the same sign, and from that, we deduce the thesis. \square

We can interpret Proposition 1 as a selective preference for lower or higher variance: if the two rewards have different

optimal actions, then either the environment is “winning” (positive μ_i ’s and over 50% probabilities of outcome 1), thus $\sigma_1 > \sigma_2$, and the *outcome* reward prefers the lower variance action, or the environment is “losing”, thus $\sigma_1 < \sigma_2$, and it prefers the higher variance action.

B. Suboptimality in nondeterministic MDPs

More generally, let $\mathcal{M} := (\mathcal{S}, \mathcal{A}, r, \mathcal{P})$ be a finite, episodic, and tabular Markov Decision Process (MDP), where \mathcal{S} is the state space, \mathcal{A} is the action space, r is the reward function and $\mathcal{P} : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ is the transition model. To model a game where the win/lose outcome is given by a final score, we make the additional assumption that \mathcal{S} contains a subset of “leaves” \mathcal{S}_L , that is, states s from which every action a takes to the terminal state with probability 1 and reward $r(s)$ depending on s but not on a . The reward is received only at the end of the game: $r(s) = 0$ for all $s \in \mathcal{S} - \mathcal{S}_L$. We denote by r_{score} this *score* reward, and by $\mathcal{M}_{\text{score}} := (\mathcal{S}, \mathcal{A}, r_{\text{score}}, \mathcal{P})$ the game where the reward is given by the score.

Correspondingly, $\mathcal{M}_{\text{outcome}} := (\mathcal{S}, \mathcal{A}, r_{\text{outcome}}, \mathcal{P})$ denotes a game with the same state space, action space, and transitions as $\mathcal{M}_{\text{score}}$, but with a win/lose outcome. Here r_{outcome} is related to r_{score} by $r_{\text{outcome}}(s) := \mathbb{1}_{r_{\text{score}}(s) > 0}$.

Any policy $\pi : \mathcal{S} \rightarrow \text{PD}(\mathcal{A})$ induces value functions on $\mathcal{M}_{\text{score}}$ and $\mathcal{M}_{\text{outcome}}$. We denote these value functions by $v_{\pi}^{\text{score}}, v_{\pi}^{\text{outcome}}, q_{\pi}^{\text{score}}$ and q_{π}^{outcome} . In particular, if π_{score} and π_{outcome} are optimal policies for $\mathcal{M}_{\text{score}}$ and $\mathcal{M}_{\text{outcome}}$ respectively, then $v_{\pi_{\text{score}}}^{\text{outcome}}(s)$ represents the probability of winning, starting from s and following a score-based optimal play, in the game $\mathcal{M}_{\text{outcome}}$. With this notation, our research question becomes whether strict inequality may sometimes hold in

$$v_{\pi_{\text{score}}}^{\text{outcome}}(s) \leq v_{\pi_{\text{outcome}}}^{\text{outcome}}(s) = \sup_{\pi} v_{\pi}^{\text{outcome}}(s). \quad (2)$$

Claim. *In a nondeterministic MDP, π_{score} can be suboptimal in $\mathcal{M}_{\text{outcome}}$.*

To prove this claim, we exhibit a class of MDPs that are exactly computable and for which the strict inequality can be numerically verified. We consider tree-like MDPs with two parameters b (state branch) and d (depth), that we call *action-shared tree* MDPs (an example in Fig. 2). More precisely, these trees are defined as follows: nonterminal states are partitioned by levels $\ell = 0, \dots, d$, with b^{ℓ} states at level ℓ . Each state s has a unique set C_s of b “children” states at level $\ell + 1$ to which one can transition, independently of the action. From the b^d leaves at level d , every action a takes to the terminal state with probability 1 and a reward independent of a .

In our numerical experiments, we fixed the state branch b and the depth d , and then we randomly generated the remaining parameters (all independently): the rewards $r_{\text{score}}(s)$ for all leaves s , as integers taken uniformly in $[-b^d, b^d]$; the transition probabilities $(p(s'|s, a))_{s' \in C_s}$ for all nonterminal (s, a) , from a b -dimensional Dirichlet distribution with all parameters equal to 1, and $p(\cdot|s, a) = 0$ outside C_s .

For each MDP so generated, we computed by value iteration the optimal policies $\pi_{\text{score}}, \pi_{\text{outcome}}$ and the optimal

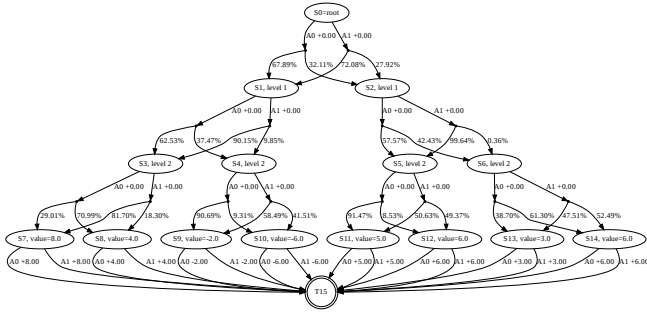


Fig. 2. A randomly generated action-shared tree MDP with state branch $b = 2$, depth $d = 3$ and 2 actions for each state. Numbers on edges denote rewards (nonzero rewards only on leaves) and transitions.

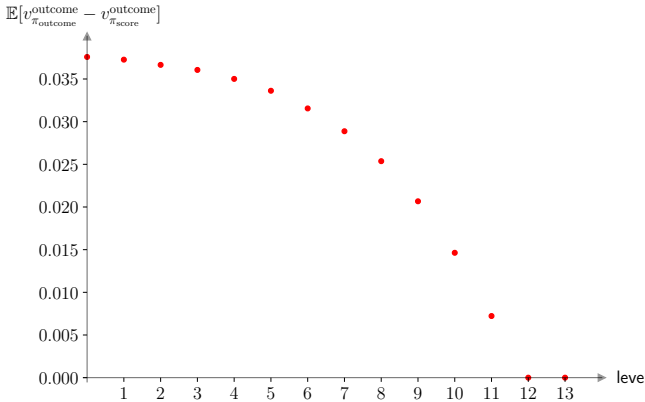


Fig. 3. The average loss on 2,000 runs in the outcome value function $v_{\pi_{outcome}}^{\text{outcome}}$ passing from the outcome-optimal policy to the score-optimal policy, as a function of the states level.

value function $v_{\pi_{outcome}}^{\text{outcome}}$. Then, we computed $v_{\pi_{score}}^{\text{outcome}}$ by policy evaluation. We finally obtained the average of the difference $v_{\pi_{outcome}}^{\text{outcome}}(s) - v_{\pi_{score}}^{\text{outcome}}(s)$ as a function of the level of the state, across the states of each MDP and across all the generated MDPs. The results are depicted in Fig. 3 and show that there are MDPs and states for which (2) holds with the strict inequality, thus proving the claim.

C. Fondness for the variance depends on the odds

As seen for the multi-armed bandit setting, when the score-optimal policy π_{score} and the win/lose outcome-optimal policy $\pi_{outcome}$ differ, it is possible to interpret the differences in terms of variances of the scores of the available actions. A similar interpretation is possible in the more general MDP setting, as stated in the following claim.

Claim. When $\pi_{outcome}$ is losing, it prefers actions leading to a higher score variance in \mathcal{M}_{score} , and vice versa when it is winning.

Though this claim is too vague to be proved rigorously, we can support it with the following experiment, performed on the same simulated action-shared tree MDPs of the previous section.

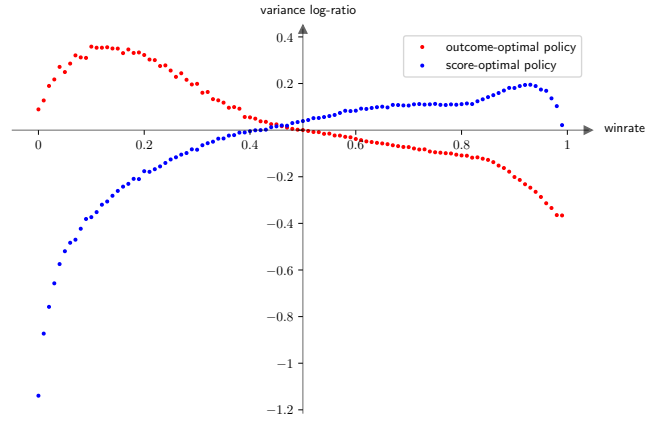


Fig. 4. Log-ratio for the score variances of the chosen and discarded action, when the two policies do not agree (median over 100 bins). The two curves are from the point of view of either policy, with variances computed along their own trajectories. The x -axis is the best winrate, i.e. the outcome value from the point of view of the outcome-optimal policy.

- 1) For assessing the preferences of the win/lose outcome-optimal policy, let $\pi^+ := \pi_{outcome}$ and $\pi^- := \pi_{score}$.
- 2) For every state-action pair (s, a) , we compute the variance $\text{Var}_{\pi^+}(r_{score}|s, a)$, that is, the variance of the final score over trajectories following π^+ from (s, a) . This quantity can be computed efficiently as $\mathbb{E}_{\pi^+}[r_{score}^2|s, a] - (q_{\pi^+}^{\text{score}}(s, a))^2$ by policy evaluation and value iteration.
- 3) For every state s , let $y(s)$ be the difference of variance (in a log scale) between the actions preferred by $\pi^+(s)$ and $\pi^-(s)$, namely, $y(s) := \log(\text{Var}_{\pi^+}(r_{score}|s, \pi^+(s))) - \log(\text{Var}_{\pi^-}(r_{score}|s, \pi^-(s)))$.
- 4) Consider only those states where $y(s) \neq 0$, that is, those states for which the action chosen by $\pi^+(s)$ and $\pi^-(s)$ actually differ. Draw a binned graph of the median of $y(s)$ against $v_{\pi_{outcome}}^{\text{outcome}}(s)$. The median is computed across the states in the same x -axis bin in each MDP and across all generated MDPs.
- 5) For assessing the preferences of the score-optimal policy restart from item 1, but with the role of π^+ and π^- reversed: $\pi^+ := \pi_{score}$ and $\pi^- := \pi_{outcome}$.

The results are depicted in Fig. 4, with $\pi_{outcome}$ in red. The figure shows that when $\pi_{outcome}$ chooses differently from π_{score} , if it is losing, then the chosen action will typically lead to a more uncertain game, with higher variance in the score. If it is winning, then the opposite happens. Moreover π_{score} , even computing variances following its own trajectories, agrees with the evaluation and typically chooses the opposite criterion.

V. LZ-SCORE

Leela Zero (LZ) [6] is an open-source Go program with no human-provided knowledge, known as one of the most faithful re-implementations of AlphaGo Zero (AGZ) [1]. For all intents and purposes, it is considered an open-source AGZ.

LZ was initially released on 25 October 2017. The neural networks powering the agent of LZ were trained by a

distributed effort, which was coordinated at the LZ website: members of the community provided computing resources by running the client, which generates self-play games and submits them to the server.

The self-play games were used to train newer networks, which were then matched against the current best and possibly promoted according to a process known as gating.

The reinforcement learning pipeline of LZ ended on 15 February 2021. LZ is available at <https://zero.sjeng.org/> and <https://github.com/leela-zero/leela-zero>.

We developed LZ-Score (LZS) as a fork of LZ, suitably modified to change the reward from the outcome of the game to the score difference. While Go is usually played on a 19×19 board, we trained LZ-Score on the 9×9 board, which is still largely used for the game of Go and is of much simpler complexity.

SAI is a fork of LZ described in [4, 11, 7, 12]. While LZ was trained on the 19×19 board only, SAI was trained on the 9×9 board, hence we choose SAI 9×9 as a benchmark to assess the strength of LZ-Score. For the purposes of this work, SAI 9×9 can be considered an AGZ-like software.

A. Setup of LZ-Score

The overall architecture of LZ was replicated in LZS, the main difference being in the quantity used as target and reward. To use the integer score difference in the place of the Boolean game outcome, both during the training (as a target) and in the inference (as a reward), the following issues were addressed.

- **Target for the neural network.** The value head of LZ’s neural network ends with a linear layer with a tanh activation and is trained against the outcome of the game in $\{-1, +1\}$, with an l^2 loss term; in LZS with board size $N = 9$, we removed the activation and changed the target to the score difference of the game, scaled by a factor $1/N^2$ in order to keep it inside $[-1, 1]$. The loss term was amplified by a factor of 20 to make its average size similar to the one of LZ. The score difference was computed with Tromp-Taylor rules at the game completely finished (no resignation allowed).
- **Reward for the MCTS.** For LZ, the exploration/exploitation tree search is driven by a UCT formula:

$$a^* = \operatorname{argmax}_a (Q(a) + U(a)), \quad Q(a) = \frac{1}{|V_a|} \sum_{i \in V_a} S(i)$$

where $S(i)$ is the scaled score difference estimated by the value head for node i , $Q(a)$ is the mean action value over the set V_a of the visited MCTS tree nodes that follow a , and $U(a)$ is the same as in [1], with the c_{puct} constant raised from 1 to 1.5 to account for the larger variability of the scaled scores of siblings with respect to the winrate probabilities.

- **Game ending.** LZ employs a set of heuristics during self-plays to avoid uselessly long games. For example, most games are allowed to end by resignation, when

the winrate drops below 5%, and there are some hard-coded conditions for passing. These heuristics do not work anymore in our new setting: to maximize the score, the agent should keep playing to the very end, so that the computationally amenable Tromp-Taylor score becomes equivalent to the proper area score of the game of Go.

Scaling down the board from size $N = 19$ to size $\rho N = 9$ with $\rho < 1$ yields several benefits.

- Average number of legal moves at each position scales down by ρ^2 . The number of good, meaningful moves might scale similarly.
- Average length of games scales down by ρ^2 .
- The suggested number of visits in the MCTS tree might scale down by as much as ρ^4 , as can be inferred from the previous two.
- The suggested number of layers in the residual convolutional neural network stack scales by ρ .
- The fully connected layers at the end of the neural network stack are smaller.

To summarize the performance benefits, the total speed improvement for self-play games can be estimated to be in the order of ρ^9 , that is about 10^{-3} . In fact, our resources allowed us to implement a pipeline with a larger number of visits. We thus chose the same number of visits of 9×9 SAI, in order to grant a fair comparison.

B. Reinforcement learning pipeline

Our training was composed of multiple phases, inspired by the original LZ training process, and by SAI well-documented pipeline. In particular, we followed almost the same training procedure as in the first run of 9×9 SAI [7], with the only exception that we reintroduced *gating* (i.e. the best network is replaced only if a stronger one emerges). Gating was removed in the passage from AlphaGo Zero to AlphaZero, as it was seen as not being strictly necessary: since we are experimenting with score-based reinforcement learning, we considered it more conservative to keep it in our pipeline. In fact, while we train the network to maximize its score, through gating we are assessing its capabilities in winning the game. Each training cycle is called a generation, and it is composed of the following steps.

- 1) **Self-play.** LZS engine uses the current best network to play a total of 2,000 *self-play games* (against itself and sharing the tree search), with the following parameters: fixed number of visits v depending on the generation (see the table below), random choice of the first 15 moves (with probability proportional to the visits number), Dirichlet noise at the root. The agent is allowed to pass only if there is no other move with an expected score larger than the Tromp-Taylor score.

generations	v	w	generations	v	w
0–15	100	4	16–31	150	8
32–47	250	12	48–63	400	16
64–79	600	20	80–∞	850	24

- 2) **Training.** Starting from the current best network, new training is started over the self-play games of a window

of the most recent w generations (see the table above). The number of steps adapts to the number of games inside the window, with a typical number of 30,000 steps. The batch size is 512.

- 3) **Gating.** The newly trained network is matched against the current best, for a total of 400 games, with the same number of visits for self-play games and no randomness in the move choice. The agent is allowed to pass when UCT selects this move. If at least 55% of the 400 matches are won by the candidate network, it is promoted to be the new best network.

All games were played with *komi* 7.5, that is, an additional 7.5 points were added to white player’s Tromp-Taylor score, to compensate for the fact that black plays the first move.

C. Outcome of the learning pipeline

The learning pipeline was guided by the intermediate results of the gating matches. In fact, following LZ, at each generation we estimated the Elo rating of the new network from the 400 games in the gating match, anchoring arbitrarily the first random network to Elo 0. This rating is clearly overrated (for example every promoted network is constrained to have at least 34.8 points, –corresponding to a 55% winrate– more than the previous one), and so we refer to it as *uncalibrated Elo*. This estimate was nonetheless adequate to assess when the learning pipeline stalled: when no new network overcame the gating for too many generations (empirically chosen), we scaled up the network size or decreased the learning rate of the training’s optimizer.

Starting from a random 2×64 network (with 2 residual convolutional blocks of 64 filters), we moved to 4×128 at 2,000 games; to 8×160 at 154,000 games; to 10×192 at 748,000 games; to 12×256 at 982,000 games. The learning rate started at 0.02 and decreased once to 0.002 at 850,000 games. Fig. 5 shows the pipeline results in terms of the uncalibrated Elo rating.

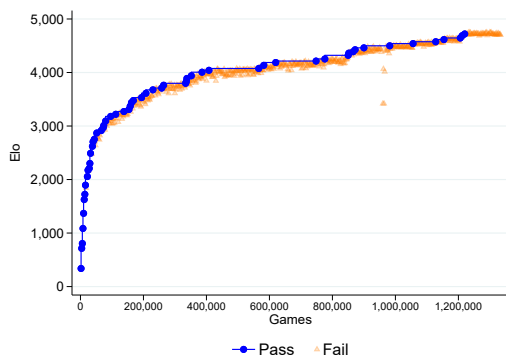


Fig. 5. Uncalibrated Elo ratings of LZS networks during the learning pipeline. Blue circles and pink triangles represent networks that respectively passed and failed gating. In this case, the rating of new networks was estimated only with the gating matches, yielding noisy and positively biased values. See also Fig. 1 for a more robust estimate.

VI. CONCLUSIONS

In this paper, we discuss when and why maximizing the score in a score-based win/lose game can be suboptimal in terms of winrate. Since this cannot happen in a deterministic game, we start with the simplest nondeterministic case, that is, the multi-armed bandit. We prove that in the multi-armed bandit, a suboptimal behavior is related to the score variance and to whether the agent is winning or losing: when losing, the outcome-based agent prefers actions leading to a larger score variance, and vice versa. We then define a class of MDPs modeling score-based win/lose games, and empirically show that this relation between suboptimal play and score variance is still valid. Finally, we train a score-based AlphaGo-like agent, and empirically show that it still behaves suboptimally, despite the fact that the game of Go is deterministic. This suggests that for certain problems, very complex deterministic games can be successfully modeled as nondeterministic, and provides sound and quantified evidence of the limitations of training a DRL score-based agent in a win/lose game, a folklore belief that we had not been able to find in the literature.

REFERENCES

- [1] David Silver et al. “Mastering the game of Go without human knowledge”. In: *Nature* 550.7676 (2017), pp. 354–359.
- [2] David Silver et al. “Mastering the game of Go with deep neural networks and tree search”. In: *Nature* 529.7587 (2016), pp. 484–489.
- [3] David Silver et al. “A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play”. In: *Science* 362.6419 (2018), pp. 1140–1144.
- [4] Francesco Morandini et al. “SAI a Sensible Artificial Intelligence that plays Go”. In: *IJCNN*. 2019, pp. 1–8. DOI: 10.1109/IJCNN.2019.8852266.
- [5] Antti Törmänen. *Invisible: the games of AlphaGo*. Hebsacker Verlag, 2017. ISBN: 978-3937499062.
- [6] Gian-Carlo Pascutto and contributors. *Leela Zero*. <http://zero.sjeng.org>.
- [7] Francesco Morandini et al. “SAI: A Sensible Artificial Intelligence That Plays with Handicap and Targets High Scores in 9×9 Go”. In: *ECAI 2020*. Vol. 325. 2020, pp. 403–410. DOI: 10.3233/FAIA200119.
- [8] David J Wu. “Accelerating Self-Play Learning in Go”. In: *arXiv* (2019). arXiv: 1902.10565.
- [9] Rémy Coulom. *Bayesian Elo Rating*. <http://www.remi-coulom.fr/Bayesian-Elo>.
- [10] Michael L. Littman. “Markov Games As a Framework for Multi-agent Reinforcement Learning”. In: *Proceedings of ICML*. 1994, pp. 157–163. URL: <http://dl.acm.org/citation.cfm?id=3091574.3091594>.
- [11] Francesco Morandini et al. “SAI: a Sensible Artificial Intelligence that plays with handicap and targets high scores in 9×9 Go (extended version)”. In: *arXiv* (2019). arXiv: 1905.10863.
- [12] Gianluca Amato and contributors. *SAI: a fork of Leela Zero with variable komi*. <https://github.com/sai-dev/sai>.