
A MC-AIXI-CTW Implementation

Group Project

Johannes Kirschner, Kerry Olesen, Jesse Wu

October 31, 2013

1 INTRODUCTION

The AIXI model [Hut05] is an attempt to solve the general AI problem. The AIXI agent interacts with an environment in cycles. Denote by \mathcal{A} , \mathcal{O} and \mathcal{R} an action, observation and reward space respectively. In each cycle, AIXI takes an action $a \in \mathcal{A}$ and receives an observation $o \in \mathcal{O}$ and a reward $r \in \mathcal{R}$. The goal of the agent is to maximize the total future reward. AIXI does not require any previous knowledge of an environment, actions are chosen based on past perceptions, which are used to build a model of the environment. Let \mathcal{M} be the model class of all chronological semi-computable semi-measures and $K(\rho)$ the Kolmogorov Complexity of ρ . Then AIXI chooses in cycle k an action

$$a_k = \arg \max_{a_k} \sum_{o_k r_k} \dots \max_{a_m} \sum_{o_m r_m} (r_k + \dots + r_m) \sum_{\rho \in \mathcal{M}} 2^{-K(\rho)} \rho(o_1 r_1 \dots o_m r_m | a_1 \dots a_m, \rho)$$

Unfortunately the AIXI model is incomputable. For all practical applications, the agent must be approximated. One approach in approximating AIXI is the MC-AIXI-CTW [VNHS09] model. Here the expectimax search is solved by an Monte-Carlo approach. in particular, a variation of the UTC [KS06] algorithm, called ρ UCT, is used to approximate a finite horizon expectimax given an environment model ρ . Within this, the UCB [Aue02] algorithm is used to balance exploration and exploitation. The class of environment models used in the implementation is a mixture of d -th order Markov Decision Process. Notably, Context Tree Weighting allows efficient linear time computation of this rather general class of models [WST95].

In comparison to AIXI, at cycle k MC-AIXI-CTW selects an action

$$a_k = \arg \max_{a_k} \sum_{o_k r_k} \dots \max_{a_m} \sum_{o_m r_m} (r_k + \dots + r_m) \sum_{M \in \mathcal{C}_D} 2^{-\Gamma_D(M)} \rho(o_1 r_1 \dots o_m r_m | a_1 \dots a_m, M)$$

Here \mathcal{C}_D is the class of all prediction suffix trees of maximum depth D , and $\Gamma_D(M)$ is the description length of a context tree M .

In the following we present our implementation of the MC-AIXI-CTW model. In Section 2 we explain how to use the program and specify different options. Section 3 describes the results of the model on several experimental environments.

2 USER MANUAL

Our approximation of aixi is written in C++ and requires g++ for compilation.

2.1 SETUP

Compile:

```
cd aixi
make
```

Run:

```
./aixi file.conf [--option1=value1 --option2=value2 ...]
```

2.2 CONFIGURATION OPTIONS

Options can be either specified in the configuration file or passed directly as `--option=value` to the program. Several configuration files are available, each specifies a particular environment and a set of default options.

AVAILABLE OPTIONS

`--ENVIRONMENT=ENV` Specifies the environment. This option is mandatory. Available environments are

- biased_rock_paper_scissor
- coinflip
- kuhn_poker
- pacman
- tiger

Further optional arguments are listed below:

Option	Effect
--CT-DEPTH=M	Maximum depth of the context tree used for prediction.
--AGENT-HORIZON=N	Number of percept/action pairs to look forward.
--MC-TIMELIMIT=N	Number N of MC simulations per cycle.
--WRITE-CT=FILE	Write CTW to file before agent termination.
--LOAD-CT=FILE	Specifies a (trained) CT to load at initialisation.
--LOG=FILE	Specifies the name of the log file.
--TERMINATE-AGE=N	The number N of agent/environment interaction cycles.
--EXPLORATION=P	Probability $0 \leq P \leq 1$ of taking a random action.
--EXPLORE-DECAY=D	Geometric decay of the exploration rate $0 \leq D \leq 1$
--INTERMEDIATE-CT=[1 0]	Set whether a CT is written at time $t = 2^k$. Default is 1.

3 EXPERIMENTAL RESULTS

3.1 EXPERIMENTAL SETUP

The performance of our agent was evaluated on five sample environments. For each environment the agent was allowed 100000 cycles to learn a model. During the learning process an exploratory constant was used. The performance of the model after various amounts of experience was then evaluated by running the agent without exploration for 5000 cycles, and the average reward per cycle reported.

The parameters used for learning each environment are shown in Figure 3.1. The experiments themselves were performed using a 3.47GHz CPU with 4GB of RAM.

Domain	CTW depth	m	ϵ	γ	ρ UCT Simulations
Biased Rock-Paper-Scissor	32	4	0.999	0.99999	500
Kuhn-Poker	42	2	0.99	0.9999	500
Partial Observable Pacman	96	4	0.9999	0.99999	500
Tiger	96	5	0.99	0.9999	500

Figure 3.1: MC-AIXI-CTW model learning configuration

m , ϵ , γ and ρ UCT Simulations correspond to the options "agent-horizon", "exploration", "explore-decay" and "mc-timelimit" as described previously.

3.2 RESULTS

All environments except for Pacman have a known optimal policy and reward. Figure 3.2 shows the average reward obtained by the agent in four environments, using a model with various cycles of experience. Generally the agent's performance improves with training, and approaches optimality. However our results are not quite as successful as those given in the original paper [VNHS09]. There are two main explanations for this.

Only 500 simulations were used during each cycle of model evaluation, while Veness et al. indicated that up to 25000 simulations are required for near optimal performance on some environments. Fewer simulations provided poorer estimates for the ρ UCT sampling process, and may occasionally lead to the selection of the incorrect action.

The most likely contributor to the difference in performance is due the implementation of model prediction. Veness et al. implement a factored CTW model, which allows for far greater predictive capabilities by using a chain of action-conditional prediction suffix trees. Each tree essentially deals with a single bit of an environment’s percept space. In contrast, our implementation uses only a single tree.

While far simpler, a single tree allows only a single model for predicting percepts, and notably cannot distinguish between observation and reward bits. Consequently larger depth trees must be used in order to extract information from percepts, and the implementation generates an overall environment model which is more complex than necessary. Rather than predicting observations and rewards individually, no differentiation between the two is possible using a single CTW tree. Clearly then, such an agent cannot be expected to perform as well as an agent which models percepts separately.

Performance on Pacman is significantly worse than on the other environments. While an optimal policy and reward for this domain are currently unknown, average rewards of approximately 2 are entirely possible. From a visual inspection of the agent playing pacman, the agent appears to have learnt how to move without bumping into walls, but still generally fails to actively find and eat food. Such poor behaviour can be explained by the lack of a factored CTW implementation, but also by a lack of training. At 100000 cycles the agent’s performance is in fact comparable to the implementation of Veness et al. With one million cycles to refine an environment model, the agent may approach positive average rewards.

Overall, we find that even with a rather limited environment model, our agent still manages to learn and improve average reward on a variety of test domains.



Figure 3.2: Average Reward per Cycle vs Experience

3.3 FURTHER EXPERIMENTS

0.2 COIN VS 0.8 COIN To see how well the agent handles a sudden change of the environment, we conducted the following experiment. The agent is trained in a simple 0.8 biased coinflip environment and the generated context tree is saved. This data is then loaded in a new 0.2 biased coinflip environment. The agent was evaluated during 20000 cycles with the following settings:

Domain	CTW depth	m	ϵ	ρ UCT Simulations	bias p
Coinflip A	16	2	0.999	100	0.8
Coinflip B	16	2	0.999	100	0.2

These values were chosen to give reasonable speed and accuracy. In detail, the following experiments were conducted:

1. environment A with exploration=0.2, *save context tree*
2. environment A with exploration=0
3. environment A with exploration=0, *context tree A loaded*
4. environment B, exploration = 0.2 (expected: same result as in 1.)
5. environment B, exploration = 0 (expected: same result as in 2.)
6. environment B with exploration = 0, *context tree A loaded*
7. environment B with exploration = 0.2, *context tree A loaded*

Experiments 1-5 should be seen as reference values for the settings in 6-7 where the agent acts in environment B with the model of environment A. The most interesting results are summarized in Figure 3.3. Compared to the normal learning case (4) the agent does not perform well anymore in the new environment if exploration is suppressed (6) and it only slowly learns the new strategy. If exploration is permitted the new environment can be learned faster (7). To reach a 0.75 average reward it took the agent 1600 cycles in setting (4), 10700 cycles in setting (6) and 5000 cycles in setting (7). Full results are presented in Figure 3.4.

BIASED RPS VS TIGER A similar experiment was also conducted on the biased rock-paper-scissor and tiger domains. The agent was evaluated during 20000 cycles with the following settings:

Domain	CTW depth	m	ϵ	γ	ρ UCT Simulations
Tiger	32	5	0.999	0.99	500
rps	32	5	0.999	0.99	500

The following test runs where recorded:

1. tiger environment, *save context tree*
2. tiger environment with exploration=0, context tree loaded
3. rock-paper-scissor environment

4. rock-paper-scissor environment with exploration=0
5. rock-paper-scissor environment, *tiger context tree loaded*
6. rock-paper-scissor environment with exploration=0, *tiger context tree loaded*

The results for (3), (5), (6) are presented in Figure 3.3. In this case the context tree structure from a different environment doesn't seem to affect the agent. Moreover, with exploration turned off the agent performs better in spite of the 'wrong' model. To explain this we have to assume that the context tree of the tiger environment is similar to the one generated in rock-paper-scissor, so forgetting old percepts is not necessary. Results of all test runs can be found in Figure 3.4.

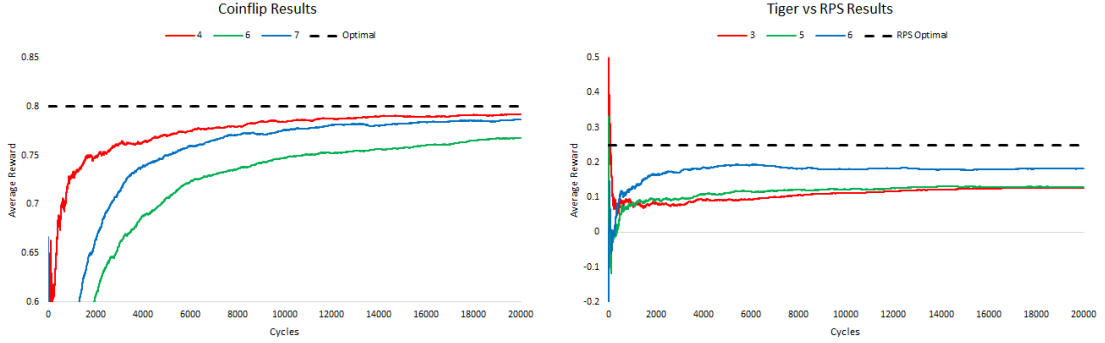


Figure 3.3: Further experiments

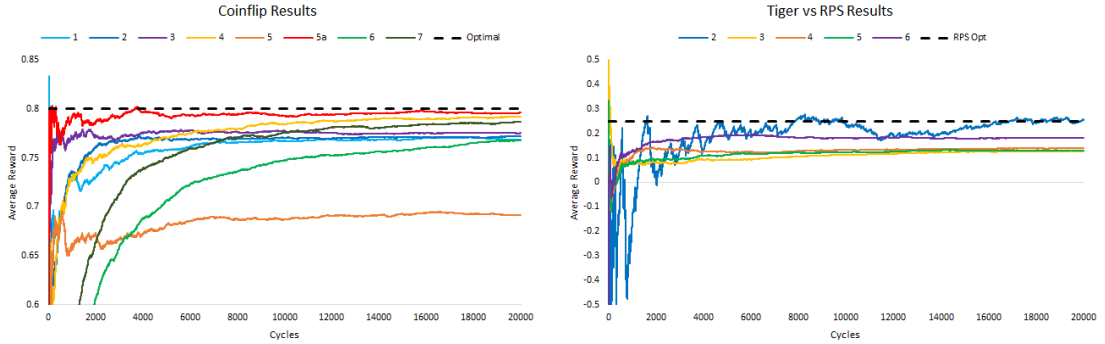


Figure 3.4: Further experiments, all results. Result (5) in the coinflip experiment seems to be an outlier, the agent was run again with the same settings resulting in (5a). The reason for this is probably the use of a single context tree instead of a factored context tree.

REFERENCES

- [Aue02] Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3:397–422, 2002.
- [Hut05] Marcus Hutter. *Universal Artificial Intelligence: Sequential Decisions based on Algorithmic Probability*. Springer, Berlin, 2005.
- [KS06] Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In *In: ECML-06. Number 4212 in LNCS*, pages 282–293. Springer, 2006.
- [VNHS09] Joel Veness, Kee Siong Ng, Marcus Hutter, and David Silver. A monte carlo aixo approximation. *CoRR*, abs/0909.0801, 2009.
- [WST95] Frans M. J. Willems, Yuri M. Shtarkov, and Tjalling J. Tjalkens. The context tree weighting method: Basic properties. *IEEE Transactions on Information Theory*, 41:653–664, 1995.