

# A MC-AIXI-CTW Implementation Group Project

---

Johannes Kirschner

Kerry Olesen

Jesse Wu

October 20, 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. More specifically, AIXI chooses in cycle  $k$  an action

$$a_k = \arg \max_{a_k} \sum_{o_k r_k} \dots \max_{a_m} \sum_{o_m r_m} (r_k + \dots + r_m) \xi(o_1 r_1 \dots o_m r_m | a_1 \dots a_m)$$

TODO: explain what  $\xi$  is, and why it is good (because it is so large)

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. The UTC [KS06] 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 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 ...]
```

Include trained ctw data?? I think this is a good idea Johannes

### 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. Available environments are

- `biased_rock_paper_scissor`
- `coinflip`
- `kuhn_poker`
- `pacman`
- `tiger`

`--MC-TIMELIMIT=N` The number  $N$  of MC simulations per cycle.

`--WRITE-CT=FILE` Write CTW to file before agent termination.

`--LOAD-CT=FILE` Specifies a (trained) CTW for the agent to load at initialisation.

`--LOG=FILE`

`--TERMINATE-AGE=N` The number  $N$  of agent/environment interaction cycles.

--EXPLORATION=P    Probability  $0 \leq P \leq 1$  that a action is chosen randomly.

--EXPLORE-DECAY=D    Decay  $0 \leq D \leq 1$  of exploration constant.  $P$  is multiplied by  $D$  in each cycle.

### 3 EXPERIMENTAL RESULTS

#### 3.1 EXPERIMENTAL SETUP

TODO: reword the following paragraphs

We evaluate the performance of our agent on five sample environments. For each environment the agent was allowed TODO: 10000 cycles to learn a model. During the learning process an exploratory constant was used. The performance of the model after at various cycle intervals was then evaluated by running the agent without exploration for TODO: 5000 cycles, and the average reward per cycle reported.

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

Domain	CTW depth	$m$	$\epsilon$	$\gamma$	$\rho$ UCT Simulations
Biased Rock-Paper-Scissor	32	3	0.1	0.999	100
Coinflip	32	3	0.1	0.999	100
Kuhn-Poker	32	3	0.1	0.999	100
Pacman	32	3	0.1	0.999	100
Tiger	32	3	0.1	0.999	100

Figure 3.1: MC-AIXI-CTW model learning configuration

#### 3.2 RESULTS

Present Results/Graphs of each environment. Maybe mention optimal result and/or scale results accordingly

All environments except for Pacman have a known optimal policy and reward.

#### 3.3 FURTHER EXPERIMENTS

Include results to do with forgetting past model - changing environments

How does the exploration constant affect learning? (tiger env)

How does aixi handle a change of the environment?

– Compare 0.3 coinflip to 0.7 coinflip

– biased rps vs tiger

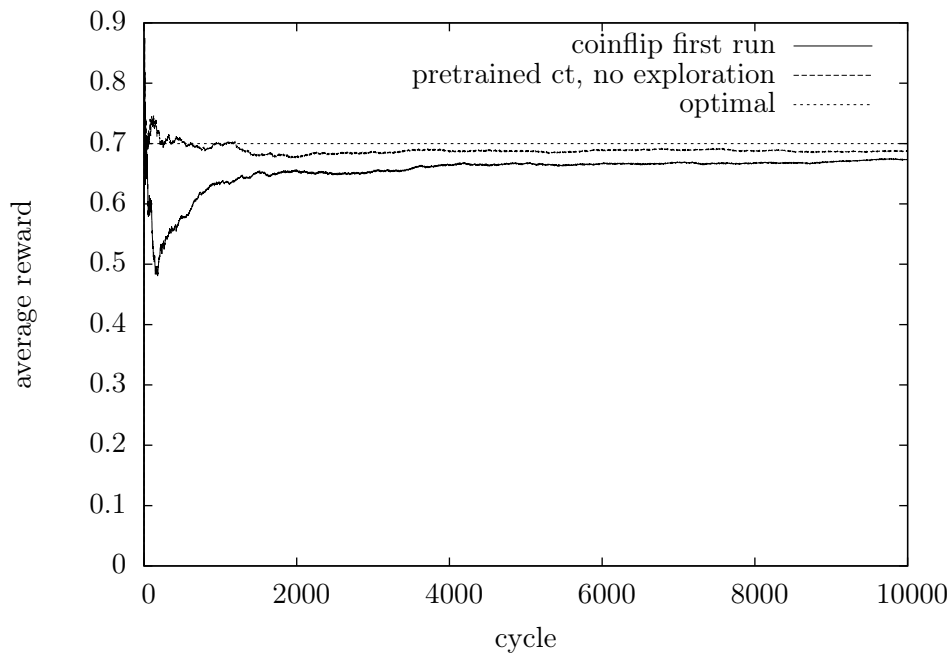


Figure 3.2: Results on the coinflip environment

### 3.4 DISCUSSION

How well did it do.

Include statistics about cycles required for optimal performance, time per cycle as in the VNHS paper [VNHS09].

### REFERENCES

- [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.