

MRAS Bandits: Control Engineering meets Slot Machines

Gayatri Ryali and Varun Sundar

Indian Institute of Technology Madras

1 Introduction

Model Reference Adaptive Search (MRAS) is a well-known algorithm in control engineering that allows for handling of Markov Decision Processes (MDPs) with either uncountable actions spaces or uncountable state spaces. To be accurate however, the MRAS algorithm is designed in the deterministic optimization context of a function $J(x)$ over variable x on space X .

The crux of the MRAS algorithm is to parametrise the sampling distribution in finding a set of candidate solutions. Each step, the parameter is updated so as to approach closer to a reference distribution. By doing so, the hope is that if the parameter updating rule is chosen appropriately, the future sampling process will be more and more concentrated on regions containing high-quality solutions. Convergence guarantees rely on the cornerstone that the reference distribution converges with its level set as the singleton optimal value. Hence, by showing that our parametric distribution closely approximates the reference distribution, we obtain the optimal value x^* . Reference [1] generalises this to situations of noisy observations of $J(x)$. We hope to utilise this model based optimisation approach in formulating a novel bandit algorithm and studying its convergence properties.

2 Objective

To obtain finite time error bounds on the MRAS algorithm (Model reference adaptive search) and utilise the ideas in creating a bandit algorithm.

3 Proposed Steps

We make all references to theorems and lemmas with respect to reference [1].

1. Prove finite time bounds theorem 4.5 of MRAS-2.
2. Explore (1) using simulations to validate the regularity conditions assumed and understand convergence related issues when these regularity conditions donot hold.
3. Formulate a bandit (stochastic bandit) setting using MRAS algorithm.
4. Set up a simulation setting using this algorithm.
5. Provide bounds for regret minimisation, and study the effect of hyperparamters (including discount λ , arms distribution $f(., \theta)$, increasing function $H(x)$ and simulation allocation rule M_k).

6. Explore possibility of MRAS being extended to structured bandit problems.

4 Proposed Algorithm

Algorithm 1 MRAS Bandits

Require: simulation budget n

- 1: parameterized distribution family $f(\cdot, \theta)$
 - 2: initial θ_0 strictly increasing function $H : R \rightarrow R_+$
 - 3: simulation allocation rule $M_k, k = 0, 1, \dots, p$, p partitions.
 - 4:
 - 5: **Denote** : $\hat{\mu}_{j,n}$ as j th sample mean, upto time n .
 - 6: **procedure** M(R)AS Bandit. Loop until: $\sum_{k=0}^{k=p} n_k M_k \geq N$. Here n_k is the no of distinct arms to draw in stage k and M_k the number of times you pull each arm at stage k .
 - 7: **Sample Arm** $a_j(k) \sim (1 - \lambda)f(\mathbf{x}, \theta_k) + \lambda f(\mathbf{x}, \theta_0)$ for $j = 1, \dots, n_k$.
 - 8: **Update sample means** of arm $a_j(k)$, pulling each arm M_k times.
 - 9: **Update parameter** $\theta_{k+1} = \operatorname{argmax}_{\theta \in \Theta} \sum_{j \in \Lambda_k} \frac{H(\hat{\mu}_{j,n})^k}{f(a_j(k), \theta_k)} \ln f(a_j(k), \theta)$.
-

We shall analyse a particular case of this algorithm where $H(x) = e^x$, $f(\cdot; \theta)$ is the categorical distribution, and all $M_k = 1$. This simplifies to:

Algorithm 2 Simplified MRAS Bandits

Require: simulation budget n

- 1: Number of arms K
 - 2: initial $\theta_0 = 1/K$
 - 3:
 - 4: **Denote** : $\hat{\mu}_{j,n}$ as j th sample mean, upto time n .
 - 5: **procedure** M(R)AS Bandit. Loop until: N .
 - 6: **Sample Arm** $a_j(k) \sim (1 - \lambda)f(\mathbf{x}, \theta_k) + \lambda f(\mathbf{x}, \theta_0)$ for $j = 1, \dots, n_k$, that is draw n_k samples.
 - 7: **Update sample means** of arm $a_j(k)$, pulling each arm once.
 - 8: **Update parameter** $\theta_{k+1} = \operatorname{argmax}_{\theta \in \Theta} \sum_{j \in \Lambda_k} \frac{e^{(C\hat{\mu}_{j,t})}}{f(a_j(k), \theta_k)} \ln f(a_j(k), \theta)$.
 - 9: $t = t + n_k$
-

However, note that we need to arrive at tractable solutions to the maximisation step (or with close approximations). For instance the case of the *Dirchlet* distribution entails derivatives of $\log(\Gamma(x))$ which can only be approximately evaluated.

References

- [1] Hyeong Soo Chang, Jiaqiao Hu, Michael C. Fu, and Steven I. Marcus. *Simulation-Based Algorithms for Markov Decision Processes*. Springer Publishing Company, Incorporated, 2nd edition, 2013. ISBN 144715021X, 9781447150213.