

Reinforcement Learning, MOEA/D and the Sample Complexity

Chen Jingjing, 55766475

July 4, 2020

Abstract

This article includes the introduction and methodologies of Reinforcement Learning(RL) and MOEA/D, as well as the analysis of sample complexity.

Keywords: reinforcement learning, MOEA/D, sample complexity

1 Multi-armed Bandit Problems

1.1 Techniques of Exploration versus Exploitation

I have written the most of them in the iPad.

1.1.1 ϵ -greedy

1.1.2 Optimistic Initial value

The Greedy actions are those that look best at present, but some of the other actions may actually be better. ϵ -greedy action selection forces the non-greedy actions to be tried, but indiscriminately, with no preference for those that are nearly greedy or particularly uncertain. It would be better to choose among the non-greedy actions according to their potential for actually being optimal.

But it's only a trick, which can not be used in general problems.

1.1.3 Upper-Confidence-Bound Action Selection

1.1.4 Gradient Bandit algorithms

Methods above are all in common which considered calculating the estimate action values then selecting actions according to these estimated value.

1.2 Contextual Bandits

Contextual Bandits aim to find the policy: a mapping from situations(states) to the actions that are best in those situations.

Contextual Bandits is a bridge between MAB and general RL for the only one situation in MAB and the states transitions in general RL.

1.3 Articles related to MAB

1.3.1 PAPER 1: Bayesian optimization for modular black-box systems with switching costs [10]

Main idea

The cost function of Bayesian Optimization in previous research depends on current inputs, while in this article, cost depends on the changes between iterations. The optimization problem is cast into a selection problem where optimal variables i (arms) are selected from a set \mathcal{K} to minimize the loss function $l : \mathcal{K} \rightarrow R$. At each iteration, the loss is $l(i^t)$ of arm i^t . GP-UCB as an acquisition function is not a new idea used in Bayesian Optimization. Under certain conditions, the iterative application of this acquisition function will converge to the true global minimum of f .

Slowly moving Bandit(SMB) algorithm took the switching cost into account, $c(i^t, i^{t-1})$, a cost when switching between arms from $t - 1$ to t . SMB aims to minimize the linear combination of the cost and the switching cost.

Comments

Authors extended SMB to the setting of black-box optimization with modular structure. At each iteration, SMB choose an arm based on the probability distribution p_t . It's s conditional sampling to encourage slow switching, to make the arm selection close to the previous choice.

2 The Research about Sample Complexity in Reinforcement Learning

At first, we want to extend the bound in sample complexity in Azar's work to the non-markovian environment. So I searched articles and listened to the conference related to this area. Meanwhile, the theoretical guarantee that huiling proposed is about sample complexity and UCB algorithm. So sample complexity in RL is a line for me to consider.

2.1 Paper 1: The Optimal Sample Complexity of PAC Learning [1]

waiting to update

Main idea

This is a basic article related to the sample complexity in supervised learning.

Comments

2.2 Paper 2: When to Trust Your Model: Model-Based Policy Optimization [5]

Abstract

Designing effective model-based reinforcement learning algorithms is difficult because the ease of data generation must be weighed against the bias of model-generated data. In this paper, we study the role of model usage in policy optimization both theoretically and empirically. We first formulate and analyze a model-based reinforcement learning algorithm with a guarantee of monotonic improvement at each step. In practice, this analysis is overly pessimistic and suggests that real off-policy data is always preferable to model-generated on-policy data, but we show that an empirical estimate of model generalization can be incorporated into such analysis to justify model usage. Motivated by this analysis, we then demonstrate that a simple procedure of using short model-generated rollouts branched from real data has the benefits of more complicated model-based algorithms without the usual pitfalls. In particular, this approach surpasses the sample efficiency of prior model-based methods, matches the asymptotic performance of the best model-free algorithms, and scales to horizons that cause other model-based methods to fail entirely.

Comments

Zhang weinan and Yu yang's work all pay attention to the model-based RL and the algorithms they proposed motivated by this article, then they all analyzed the sample complexity of the new algorithm in their article.

2.3 Paper 3: The Sample-Complexity of General Reinforcement Learning

Main idea Waiting to update. **Comments**

3 Multi-objective Evolutionary Algorithm based on Decomposition

3.1 MOEA/D [2]

There are two main lines in the MOEA/D related works [4].

We will concentrate on the selection of weighted vectors next week.

3.2 MOEA/D-M2M [11]

MOEA/D-M2M decomposes the multi-objective problems(MOP) into several small MOPs. The decision space is divided by k unit direction vectors into k regions. Compared to MOEA/D, it can promote the diversity of the population.

3.3 MOEA/D-VLP [8]

Motivated by metameric genetic algorithm(MGA) [13], Test problems proposed by Lihui were all composed by the shape function $\alpha_i(x_I)$ plus distance function $\beta_i(x_{II})$. The challenge of these kind of test problems I found is that MOEA/D, MOEA/D-M2M and NSGA-II in approximating the higher dimension.

MOEA/D-VLP combined the idea of MOEA/D and MOEA/D-M2M, but not the same. It decompose the original MOP into small MOPs according to the dimension.

Meanwhile, I read Zhenkun's article [9] about the construction of the shape function $h_k(x_I)$ and distance function $g_k(x_{II})$.

3.4 Multi-objective MAB

3.4.1 PAPER 1: Adaptive operator selection with bandits for a multi-objective evolutionary algorithm based on decomposition [12]

Using MAB to perform AOS.

3.4.2 PAPER 2: waiting to explore more articles in this area

References

- [1] Hanneke S. The optimal sample complexity of PAC learning[J]. The Journal of Machine Learning Research, 2016, 17(1): 1319-1333.
- [2] Zhang Q, Li H. MOEA/D: A multiobjective evolutionary algorithm based on decomposition[J]. IEEE Transactions on evolutionary computation, 2007, 11(6): 712-731.

- [3] Li K, Fialho A, Kwong S, et al. Adaptive operator selection with bandits for a multiobjective evolutionary algorithm based on decomposition[J]. *IEEE Transactions on Evolutionary Computation*, 2013, 18(1): 114-130.
- [4] Trivedi A, Srinivasan D, Sanyal K, et al. A survey of multiobjective evolutionary algorithms based on decomposition[J]. *IEEE Transactions on Evolutionary Computation*, 2016, 21(3): 440-462.
- [5] Janner M, Fu J, Zhang M, et al. When to trust your model: Model-based policy optimization[C]//*Advances in Neural Information Processing Systems*. 2019: 12519-12530.
- [6] Li K, Malik J. Learning to optimize[J]. *arXiv preprint arXiv:1606.01885*, 2016.
- [7] Lattimore T, Hutter M, Sunehag P. The sample-complexity of general reinforcement learning[C]//*Proceedings of the 30th International Conference on Machine Learning*. *Journal of Machine Learning Research*, 2013.
- [8] Li H, Deb K, Zhang Q. Variable-length Pareto optimization via decomposition-based evolutionary multiobjective algorithm[J]. *IEEE Transactions on Evolutionary Computation*, 2019, 23(6): 987-999.
- [9] Z. Wang, Y. Ong and H. Ishibuchi, "On Scalable Multiobjective Test Problems With Hardly Dominated Boundaries," in *IEEE Transactions on Evolutionary Computation*, vol. 23, no. 2, pp. 217-231, April 2019, doi: 10.1109/TEVC.2018.2844286.
- [10] Chi-Heng Lin, Joseph D. Miano, Eva L. Dyer: "Bayesian optimization for modular black-box systems with switching costs", 2020; [<http://arxiv.org/abs/2006.02624> arXiv:2006.02624].
- [11] H. Liu, F. Gu and Q. Zhang, "Decomposition of a Multiobjective Optimization Problem Into a Number of Simple Multiobjective Subproblems," in *IEEE Transactions on Evolutionary Computation*, vol. 18, no. 3, pp. 450-455, June 2014, doi: 10.1109/TEVC.2013.2281533.
- [12] Handoko S D, Nguyen D T, Yuan Z, et al. Reinforcement learning for adaptive operator selection in memetic search applied to quadratic assignment problem[C]//*Proceedings of the Companion Publication of the 2014 Annual Conference on Genetic and Evolutionary Computation*. 2014: 193-194.
- [13] Ryserkerk, M.L., Averill, R.C., Deb, K. et al. Solving metameric variable-length optimization problems using genetic algorithms. *Genet Program Evolvable Mach* 18, 247–277 (2017). <https://doi.org/10.1007/s10710-016-9282-8>.