

Contents

| | |
|---|----------|
| 1 Preliminaries | 1 |
| 1.1 Markov Decision Processes | 1 |
| Bibliography | 4 |

Chapter 1

Preliminaries

This is an introduction...

1.1 Markov Decision Processes

A [Markov Decision Process \(MDP\)](#) is a way to model the interaction between an agent, which is the learner and the decision maker, and an environment. When the agent performs an action a , the environment responds presenting a new state s to the agent and rewarding the agent with a certain scalar signal called *immediate reward*. Importantly, the environment dynamics of a [MDP](#) are *stationary*: they do not depend upon time. Moreover, the current state s_{h+1} at time step h must depend only on the previous state s_h and action a_h . This property, called Markovian Property, entails that the agent is somehow pushed to forget the states and actions of the past. However, there is another mechanism of [MDPs](#) which encourages the agent to take into account the long term consequences of its choices. In fact, the objective of the agent is to maximize over time the *cumulative reward* or *return*, which is the cumulated sum of the immediate rewards obtained after each action undertaken by the agent. Evidently, sometimes it is more profitable to sacrifice immediate reward in order to reach a higher cumulative reward in the long term. This mechanism pushes the agent to take into account the future in its current decisions.

Hence, [MDPs](#) are a powerful tool capable of modelling many challenges that we encounter in life, science and engineering. Take, for example, a hungry newborn infant in his mother's arms. Through various attempts, the infant moves around its arms, head and mouth looking for food. In this ways, he

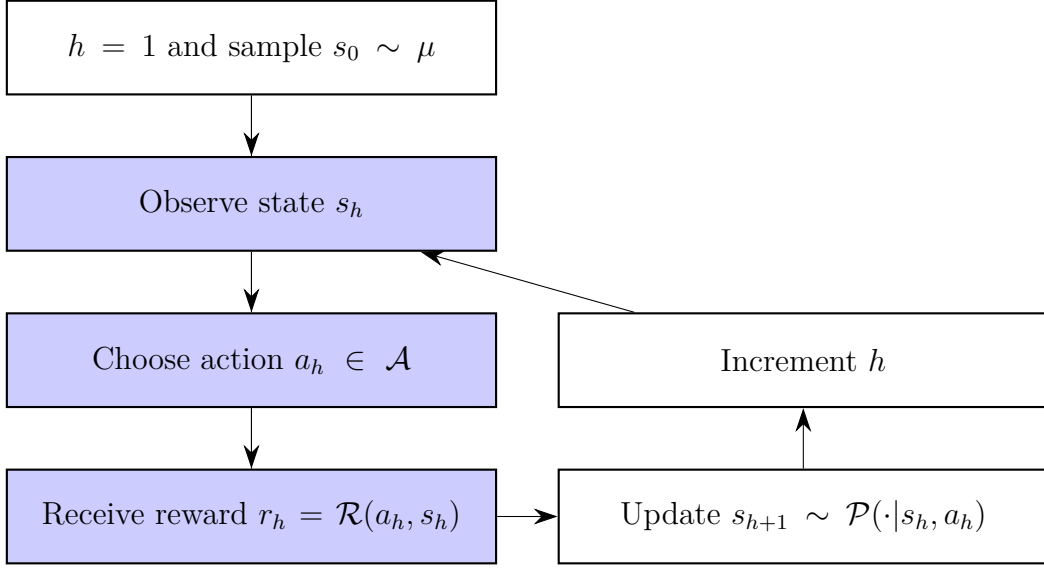


Figure 1.1: Interaction protocol for Markov decision processes

changes its state until the moment he eventually receives a positive reward, the mother's milk. Little by little, the baby will learn the precise sequence of actions and states that rewards him with milk.

Formally, a discrete-time continuous **MDP** [64, 73] is a tuple $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, \mu \rangle$, where:

- $\mathcal{S} \in \mathbb{R}^{d_{\mathcal{S}}}$ is the $d_{\mathcal{S}}$ -dimensional continuous *state space*, i.e., the set of all possible observable states of the environment;
- $\mathcal{A} \in \mathbb{R}^{d_{\mathcal{A}}}$ is the $d_{\mathcal{A}}$ -dimensional continuous *action space*, i.e., the set of all the possible actions that the agent can perform. Sometimes, not all the actions are performable in all states. In such cases, we can define the set $\mathcal{A}(s)$ for all $s \in \mathcal{S}$, s.a $\mathcal{A} = \bigcup_{s \in \mathcal{S}} \mathcal{A}(s)$;
- $\mathcal{P} : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ is a function called *transition model* such that, for every $s \in \mathcal{S}$ and $a \in \mathcal{A}$, it assigns a probability measure $\mathcal{P}(\cdot|s, a)$ over \mathcal{S} . Its corresponding probability density function is $P(\cdot|s, a)$;
- $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow [-R_{\max}, R_{\max}]$ is a bounded *reward* function such that, every time the agent chooses action $a \in \mathcal{A}$ in state $s \in \mathcal{S}$, it receives a reward $r = \mathcal{R}(a, s)$. R_{\max} is the maximum absolute reward that the agent can receive;
- $\gamma \in (0, 1]$ is a discount factor to apply to future rewards;

- $\mu \in \Delta(\mathcal{S})$ is the initial state distribution, such that the initial state is drawn as $s_0 \sim \mu$.

In the more general case, the state space \mathcal{S} and the action space \mathcal{A} , which are the sensor and actuator possibilities of the agent, respectively, can be both continuous and discrete, finite or infinite. In what follows, we will focus on the continuous case because it is the most relevant to our work.

The time is typically modelled as a discrete sequence of decision steps represented by the natural numbers: $\mathcal{H} = \{0, 1, \dots, H\}$, where $H \in \mathbb{N}$ is the *horizon* of a given task, which can be either infinite $H = \infty$ or finite. The agent-environment interaction ends when either the horizon is reached or a *terminal state* is reached. Tasks that always end in a finite amount of steps (an *episode*) are called *episodic*. A *trajectory* is the sequence of states and actions $\tau = \{s_0, a_0, s_1, a_1, \dots, s_{H-1}, a_{H-1}\}$ up to the last time step of the episode.

The core of a [MDP](#) agent is the *policy* π , a mapping from perceived states of the environment to possible actions. A policy can also take as input the past history of set and actions, but in our work we only consider *memoryless policies*. Such mappings can be deterministic $\pi : \mathcal{S} \rightarrow \mathcal{A}$ or stochastic $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$, such that the current action is drawn as $a_h \sim \pi(\cdot | s_h)$. Note that the memoryless property of policies does not imply that the agent cannot take into account what it has seen before when making an action choice. In fact, the agent has multiple ways to leverage its experience and learn through time, as we will see later on.

Bibliography

- [1] ABBASI-YADKORI, Y., PÁL, D., AND SZEPESVÁRI, C. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems* (2011), pp. 2312–2320.
- [2] AGRAWAL, R. Sample mean based index policies by $o(\log n)$ regret for the multi-armed bandit problem. *Advances in Applied Probability* 27, 4 (1995), 1054–1078.
- [3] ATAN, O., TEKIN, C., AND SCHAAR, M. Global multi-armed bandits with hölder continuity. In *Artificial Intelligence and Statistics* (2015), pp. 28–36.
- [4] AUER, P. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research* 3, Nov (2002), 397–422.
- [5] AUER, P., CESA-BIANCHI, N., AND FISCHER, P. Finite-time analysis of the multiarmed bandit problem. *Machine learning* 47, 2-3 (2002), 235–256.
- [6] AUER, P., ORTNER, R., AND SZEPESVÁRI, C. Improved rates for the stochastic continuum-armed bandit problem. In *International Conference on Computational Learning Theory* (2007), Springer, pp. 454–468.
- [7] BACH, F. R., AND BLEI, D. M., Eds. *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015* (2015), vol. 37 of *JMLR Workshop and Conference Proceedings*, JMLR.org.
- [8] BELLEMARE, M., SRINIVASAN, S., OSTROVSKI, G., SCHAUL, T., SEXTON, D., AND MUNOS, R. Unifying count-based exploration and intrinsic motivation. In *Advances in Neural Information Processing Systems* (2016), pp. 1471–1479.

- [9] BENGIO, S., WALLACH, H. M., LAROCHELLE, H., GRAUMAN, K., CESA-BIANCHI, N., AND GARNETT, R., Eds. *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 3-8 December 2018, Montréal, Canada* (2018).
- [10] BOUCHERON, S., LUGOSI, G., AND MASSART, P. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- [11] BRAFMAN, R. I., AND TENNENHOLTZ, M. R-max-a general polynomial time algorithm for near-optimal reinforcement learning. *Journal of Machine Learning Research* 3, Oct (2002), 213–231.
- [12] BROCKMAN, G., CHEUNG, V., PETTERSSON, L., SCHNEIDER, J., SCHULMAN, J., TANG, J., AND ZAREMBA, W. Openai gym, 2016.
- [13] BUBECK, S., CESA-BIANCHI, N., ET AL. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning* 5, 1 (2012), 1–122.
- [14] BUBECK, S., CESA-BIANCHI, N., AND LUGOSI, G. Bandits with heavy tail. *IEEE Transactions on Information Theory* 59, 11 (2013), 7711–7717.
- [15] BUBECK, S., AND MUNOS, R. Open loop optimistic planning. In *COLT 2010 - The 23rd Conference on Learning Theory, Haifa, Israel, June 27-29, 2010* (2010), pp. 477–489.
- [16] BUBECK, S., STOLTZ, G., SZEPESVÁRI, C., AND MUNOS, R. Online optimization in x-armed bandits. In *Advances in Neural Information Processing Systems* (2009), pp. 201–208.
- [17] CESA-BIANCHI, N., AND LUGOSI, G. Combinatorial bandits. *Journal of Computer and System Sciences* 78, 5 (2012), 1404–1422.
- [18] CHAPELLE, O., AND LI, L. An empirical evaluation of thompson sampling. In *Advances in neural information processing systems* (2011), pp. 2249–2257.
- [19] CHEN, W., WANG, Y., YUAN, Y., AND WANG, Q. Combinatorial multi-armed bandit and its extension to probabilistically triggered arms. *The Journal of Machine Learning Research* 17, 1 (2016), 1746–1778.

- [20] CHOWDHURY, S. R., AND GOPALAN, A. Online learning in kernelized markov decision processes. *CoRR abs/1805.08052* (2018).
- [21] COCHRAN, W. G. *Sampling techniques*. John Wiley & Sons, 2007.
- [22] CORTES, C., MANSOUR, Y., AND MOHRI, M. Learning bounds for importance weighting. In *Advances in Neural Information Processing Systems 23*, J. D. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S. Zemel, and A. Culotta, Eds. Curran Associates, Inc., 2010, pp. 442–450.
- [23] DANN, C., AND BRUNSKILL, E. Sample complexity of episodic fixed-horizon reinforcement learning. In *Advances in Neural Information Processing Systems* (2015), pp. 2818–2826.
- [24] DANN, C., LATTIMORE, T., AND BRUNSKILL, E. Unifying pac and regret: Uniform pac bounds for episodic reinforcement learning. In *Advances in Neural Information Processing Systems* (2017), pp. 5713–5723.
- [25] DEISENROTH, M. P., NEUMANN, G., PETERS, J., ET AL. A survey on policy search for robotics. *Foundations and Trends® in Robotics 2*, 1–2 (2013), 1–142.
- [26] DORATO, P., ABDALLAH, C. T., CERONE, V., AND JACOBSON, D. H. *Linear-quadratic control: an introduction*. Prentice Hall Englewood Cliffs, NJ, 1995.
- [27] DY, J. G., AND KRAUSE, A., Eds. *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018* (2018), vol. 80 of *JMLR Workshop and Conference Proceedings*, JMLR.org.
- [28] ESPEHOLT, L., SOYER, H., MUNOS, R., SIMONYAN, K., MNIH, V., WARD, T., DORON, Y., FIROIU, V., HARLEY, T., DUNNING, I., LEGG, S., AND KAVUKCUOGLU, K. IMPALA: scalable distributed deep-rl with importance weighted actor-learner architectures. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018* (2018), pp. 1406–1415.
- [29] GIL, M., ALAJAJI, F., AND LINDER, T. Rényi divergence measures for commonly used univariate continuous distributions. *Information Sciences 249* (2013), 124–131.

- [30] HAARNOJA, T., TANG, H., ABBEEL, P., AND LEVINE, S. Reinforcement learning with deep energy-based policies. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017* (2017), pp. 1352–1361.
- [31] HAARNOJA, T., ZHOU, A., ABBEEL, P., AND LEVINE, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018* (2018), pp. 1856–1865.
- [32] HERSHEY, J. R., AND OLSEN, P. A. Approximating the kullback leibler divergence between gaussian mixture models. In *Acoustics, Speech and Signal Processing, 2007. ICASSP 2007. IEEE International Conference on* (2007), vol. 4, IEEE, pp. IV–317.
- [33] HOUTHOOFT, R., CHEN, X., DUAN, Y., SCHULMAN, J., DE TURCK, F., AND ABBEEL, P. Vime: Variational information maximizing exploration. In *Advances in Neural Information Processing Systems* (2016), pp. 1109–1117.
- [34] IONIDES, E. L. Truncated importance sampling. *Journal of Computational and Graphical Statistics* 17, 2 (2008), 295–311.
- [35] JAKSCH, T., ORTNER, R., AND AUER, P. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research* 11, Apr (2010), 1563–1600.
- [36] JIN, C., ALLEN-ZHU, Z., BUBECK, S., AND JORDAN, M. I. Is q-learning provably efficient? In *Advances in Neural Information Processing Systems* (2018), pp. 4868–4878.
- [37] KALAI, A. T., AND MOHRI, M., Eds. *COLT 2010 - The 23rd Conference on Learning Theory, Haifa, Israel, June 27-29, 2010* (2010), Omnipress.
- [38] KALLUS, N. Instrument-armed bandits. In *Algorithmic Learning Theory* (2018), pp. 529–546.
- [39] KAUFMANN, E., KORDA, N., AND MUNOS, R. Thompson sampling: An asymptotically optimal finite-time analysis. In *International Conference on Algorithmic Learning Theory* (2012), Springer, pp. 199–213.

-
- [40] KEARNS, M., AND SINGH, S. Near-optimal reinforcement learning in polynomial time. *Machine learning* 49, 2-3 (2002), 209–232.
 - [41] KLEINBERG, R., SLIVKINS, A., AND UPFAL, E. Multi-armed bandits in metric spaces. In *Proceedings of the fortieth annual ACM symposium on Theory of computing* (2008), ACM, pp. 681–690.
 - [42] KLEINBERG, R., SLIVKINS, A., AND UPFAL, E. Bandits and experts in metric spaces. *arXiv preprint arXiv:1312.1277* (2013).
 - [43] KLEINBERG, R. D. Nearly tight bounds for the continuum-armed bandit problem. In *Advances in Neural Information Processing Systems* (2005), pp. 697–704.
 - [44] LAI, T. L., AND ROBBINS, H. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics* 6, 1 (1985), 4–22.
 - [45] LAKSHMANAN, K., ORTNER, R., AND RYABKO, D. Improved regret bounds for undiscounted continuous reinforcement learning. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015* (2015), pp. 524–532.
 - [46] LATTIMORE, T., AND HUTTER, M. Near-optimal pac bounds for discounted mdps. *Theoretical Computer Science* 558 (2014), 125–143.
 - [47] LATTIMORE, T., AND SZEPESVÁRI, C. Bandit algorithms.
 - [48] LATTIMORE, T., AND SZEPESVÁRI, C. *Bandit Algorithms*. Cambridge University Press (preprint), 2019.
 - [49] MCCLUSKEY, L., WILLIAMS, B. C., SILVA, J. R., AND BONET, B., Eds. *Proceedings of the Twenty-Second International Conference on Automated Planning and Scheduling, ICAPS 2012, Atibaia, São Paulo, Brazil, June 25-19, 2012* (2012), AAAI.
 - [50] MEDINA, A. M., AND YANG, S. No-regret algorithms for heavy-tailed linear bandits. In *Proceedings of The 33rd International Conference on Machine Learning* (New York, New York, USA, 20–22 Jun 2016), M. F. Balcan and K. Q. Weinberger, Eds., vol. 48 of *Proceedings of Machine Learning Research*, PMLR, pp. 1642–1650.
 - [51] MERSEREAU, A. J., RUSMEVICHIENTONG, P., AND TSITSIKLIS, J. N. A structured multiarmed bandit problem and the greedy policy. *IEEE Transactions on Automatic Control* 54, 12 (2009), 2787–2802.

- [52] METELLI, A. M., PAPINI, M., FACCIO, F., AND RESTELLI, M. Policy optimization via importance sampling. In *Advances in Neural Information Processing Systems* (2018), pp. 5447–5459.
- [53] MNIH, V., BADIA, A. P., MIRZA, M., GRAVES, A., LILICRAP, T. P., HARLEY, T., SILVER, D., AND KAVUKCUOGLU, K. Asynchronous methods for deep reinforcement learning. In *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016* (2016), pp. 1928–1937.
- [54] OK, J., PROUTIERE, A., AND TRANOS, D. Exploration in structured reinforcement learning. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 3-8 December 2018, Montréal, Canada.* (2018), pp. 8888–8896.
- [55] ORTNER, R., AND RYABKO, D. Online regret bounds for undiscounted continuous reinforcement learning. In *Proceedings of the 25th International Conference on Neural Information Processing Systems-Volume 2* (2012), Curran Associates Inc., pp. 1763–1771.
- [56] OSBAND, I., RUSSO, D., AND ROY, B. V. (more) efficient reinforcement learning via posterior sampling. In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States.* (2013), pp. 3003–3011.
- [57] OSBAND, I., VAN ROY, B., AND WEN, Z. Generalization and exploration via randomized value functions. In *International Conference on Machine Learning* (2016), pp. 2377–2386.
- [58] OWEN, A., AND ZHOU, Y. Safe and Effective Importance Sampling. *Journal of the American Statistical Association* (Mar. 2000), 135–143.
- [59] OWEN, A. B. *Monte Carlo theory, methods and examples*. 2013.
- [60] OWEN, A. B. *Monte Carlo theory, methods and examples*. 2013.
- [61] PANDEY, S., CHAKRABARTI, D., AND AGARWAL, D. Multi-armed bandit problems with dependent arms. In *Proceedings of the 24th international conference on Machine learning* (2007), ACM, pp. 721–728.

-
- [62] PETERS, J., AND SCHAAL, S. Reinforcement learning of motor skills with policy gradients. *Neural networks* 21, 4 (2008), 682–697.
- [63] PRECUP, D., AND TEH, Y. W., Eds. *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017* (2017), vol. 70 of *Proceedings of Machine Learning Research*, PMLR.
- [64] PUTERMAN, M. L. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- [65] RÉNYI, A. On measures of entropy and information. Tech. rep., Hungarian Academy of Sciences Budapest Hungary, 1961.
- [66] ROBBINS, H. Some aspects of the sequential design of experiments. In *Herbert Robbins Selected Papers*. Springer, 1985, pp. 169–177.
- [67] SARITAÇ, A. Ö., AND TEKIN, C. Combinatorial multi-armed bandit problem with probabilistically triggered arms: A case with bounded regret. In *Signal and Information Processing (GlobalSIP), 2017 IEEE Global Conference on* (2017), IEEE, pp. 111–115.
- [68] SCHULMAN, J., LEVINE, S., ABBEEL, P., JORDAN, M., AND MORITZ, P. Trust region policy optimization. In *International Conference on Machine Learning* (2015), pp. 1889–1897.
- [69] SEHNKE, F., OSENDORFER, C., RÜCKSTIESS, T., GRAVES, A., PETERS, J., AND SCHMIDHUBER, J. Policy gradients with parameter-based exploration for control. In *International Conference on Artificial Neural Networks* (2008), Springer, pp. 387–396.
- [70] SILVER, D., LEVER, G., HEES, N., DEGRIS, T., WIERSTRA, D., AND RIEDMILLER, M. Deterministic policy gradient algorithms. In *ICML* (2014).
- [71] SRINIVAS, N., KRAUSE, A., KAKADE, S., AND SEEGER, M. Gaussian process optimization in the bandit setting: No regret and experimental design. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)* (Haifa, Israel, June 2010), J. Fürnkranz and T. Joachims, Eds., Omnipress, pp. 1015–1022.

-
- [72] STREHL, A. L., LI, L., AND LITTMAN, M. L. Reinforcement learning in finite mdps: Pac analysis. *Journal of Machine Learning Research* 10, Nov (2009), 2413–2444.
- [73] SUTTON, R. S., AND BARTO, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- [74] SUTTON, R. S., MCALLESTER, D. A., SINGH, S. P., AND MANSOUR, Y. Policy gradient methods for reinforcement learning with function approximation. In *Advances in neural information processing systems* (2000), pp. 1057–1063.
- [75] THOMPSON, W. R. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika* 25, 3/4 (1933), 285–294.
- [76] VAKILI, S., LIU, K., AND ZHAO, Q. Deterministic sequencing of exploration and exploitation for multi-armed bandit problems. *arXiv preprint arXiv:1106.6104* (2011).
- [77] VEACH, E., AND GUIBAS, L. J. Optimally combining sampling techniques for Monte Carlo rendering. In *Proceedings of the 22nd annual conference on Computer graphics and interactive techniques - SIGGRAPH '95* (1995), ACM Press, pp. 419–428.
- [78] WANG, Z., ZHOU, R., AND SHEN, C. Regional multi-armed bandits. In *International Conference on Artificial Intelligence and Statistics* (2018), pp. 510–518.
- [79] WEINSTEIN, A., AND LITTMAN, M. L. Bandit-based planning and learning in continuous-action markov decision processes. In *Proceedings of the Twenty-Second International Conference on Automated Planning and Scheduling, ICAPS 2012, Atibaia, São Paulo, Brazil, June 25-19, 2012* (2012).
- [80] YU, X., SHAO, H., LYU, M. R., AND KING, I. Pure exploration of multi-armed bandits with heavy-tailed payoffs.