



MTA
SZTAKI

Online Learning in Non-Stationary Markov Decision Processes



Gergely Neu

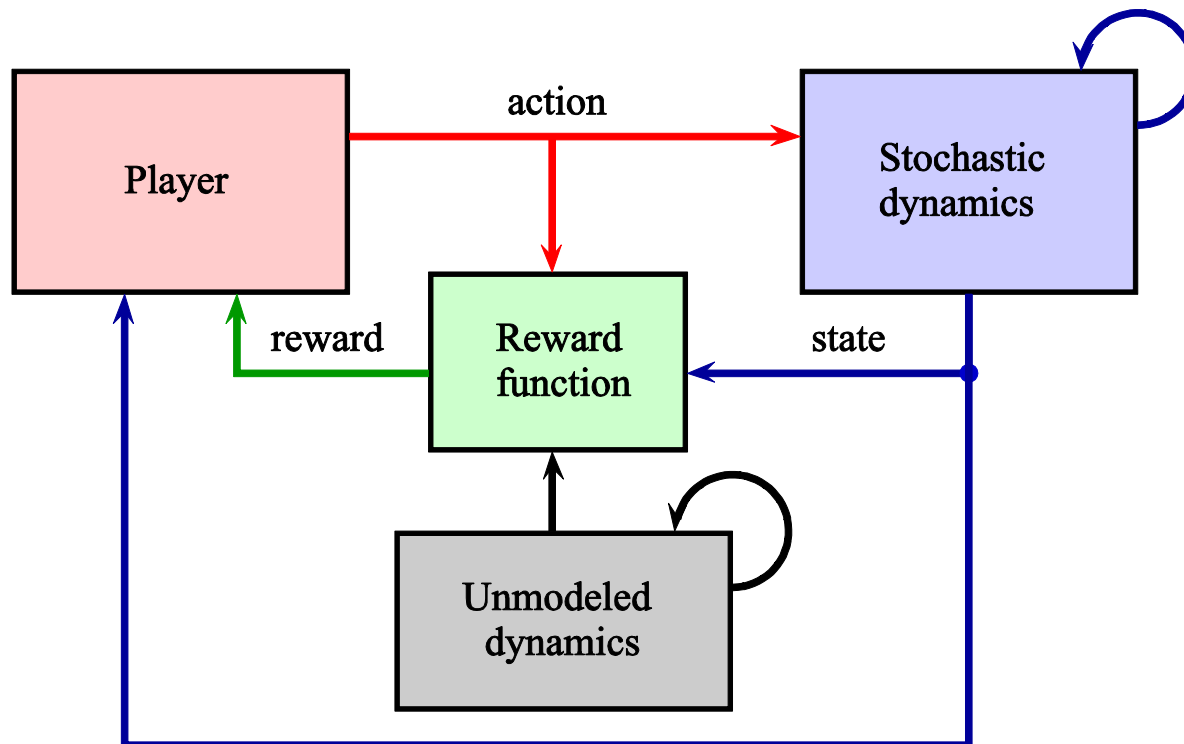
gergely.neu@gmail.com
MTA SZTAKI, Hungary

Csaba Szepesvári

szepesva@ualberta.ca
University of Alberta,
Canada

András György

gyorgy@ualberta.ca
University of Alberta,
Canada



The learning problem

Goal: minimize total expected regret:

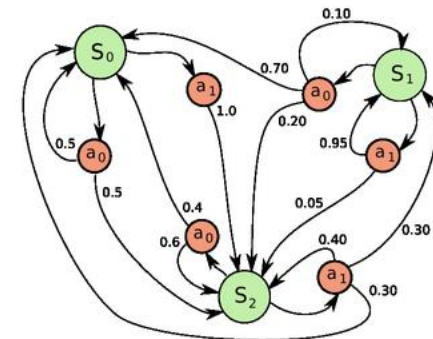
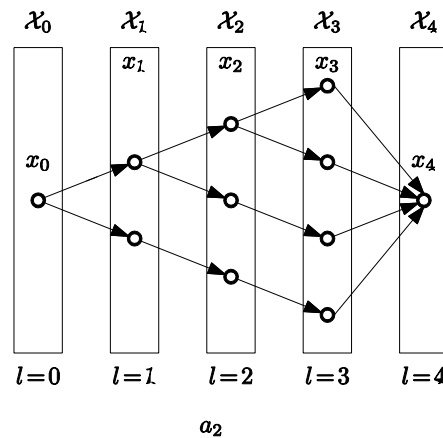
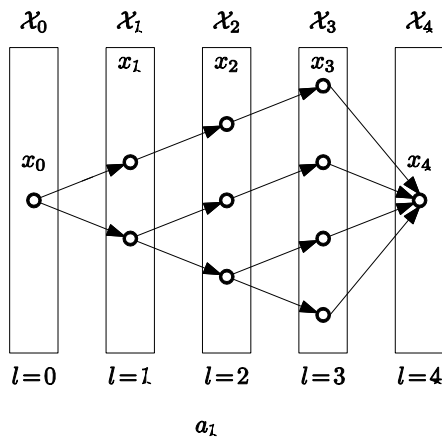
$$\hat{L}_T = \max_{\pi} E[R_T^{\pi} - \hat{R}_T] \rightarrow \min$$

Total reward of policy π :

$$R_T^{\pi} = \sum_{t=1}^T r_t(x'_t, \pi(x'_t))$$

Total reward of agent:

$$\hat{R}_T = \sum_{t=1}^T r_t(x_t, a_t)$$



Results

	Whole r_t observed	Only $r_t(x_t, a_t)$ observed
P known	<ul style="list-style-type: none">• Even-Dar et al. (2006,2009)<ul style="list-style-type: none">• Unichain MDP• $\hat{L}_T = O(\tau^2 \sqrt{T \log A })$	<ul style="list-style-type: none">• Neu et al. (2010a, 2013a)<ul style="list-style-type: none">• Loop-free SSP• $\hat{L}_T = O(L^2 \sqrt{T A /\alpha})$• Neu et al. (2010b, 2013b)<ul style="list-style-type: none">• Unichain MDP• $\hat{L}_T = O(\tau^{3/2} \sqrt{T A /\alpha'})$
P unknown	<ul style="list-style-type: none">• Jaksch et al. (2010)<ul style="list-style-type: none">• Stochastic rewards• Connected MDP• $\hat{L}_T = \tilde{O}(D X \sqrt{T A })$• Neu et al. (2012)<ul style="list-style-type: none">• Loop-free SSP• $\hat{L}_T = \tilde{O}(L X A \sqrt{T})$	