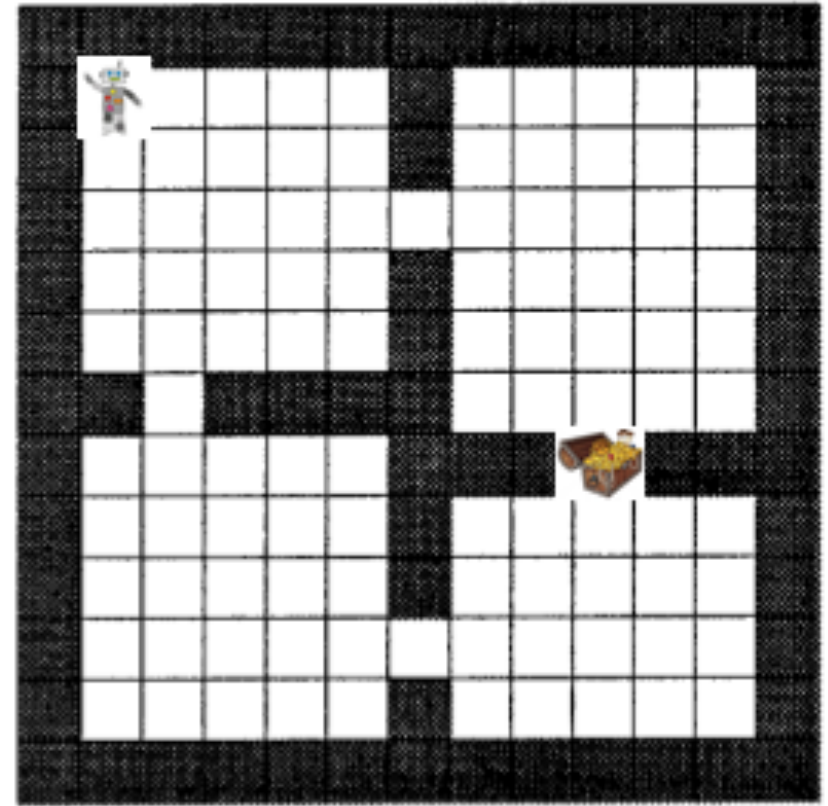
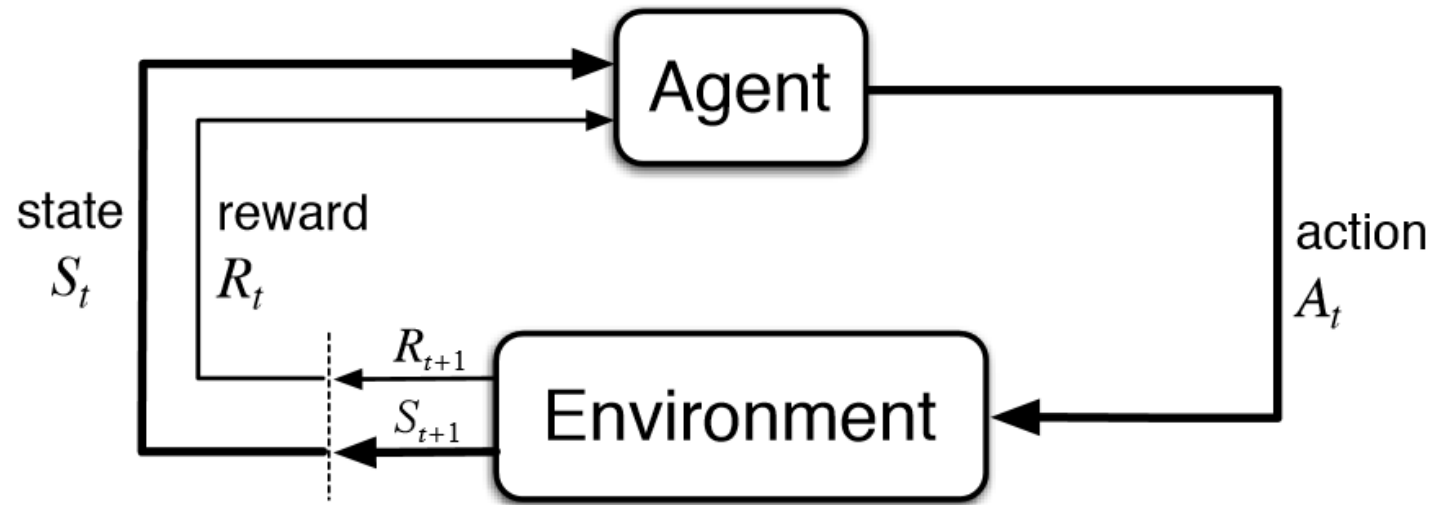


Playing with temporal abstraction for reinforcement learning

Author: Théophile Gervet

Advisor: Doina Precup

What is reinforcement learning?



$$\mathcal{P}_{ss'}^a = P(S_{t+1} = s' \mid S_t = s, A_t = a)$$

$$\pi(a \mid s) = P(A_t = a \mid S_t = s)$$

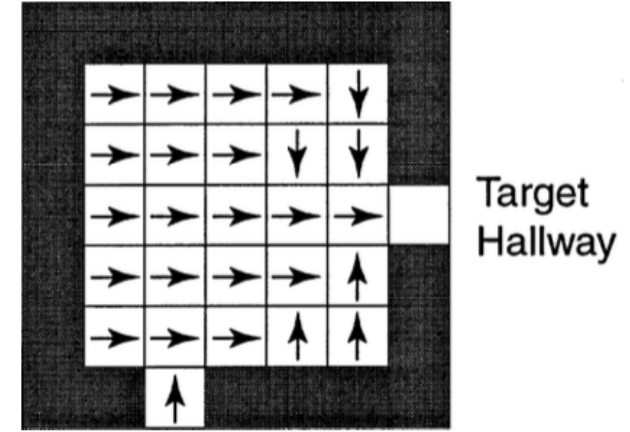
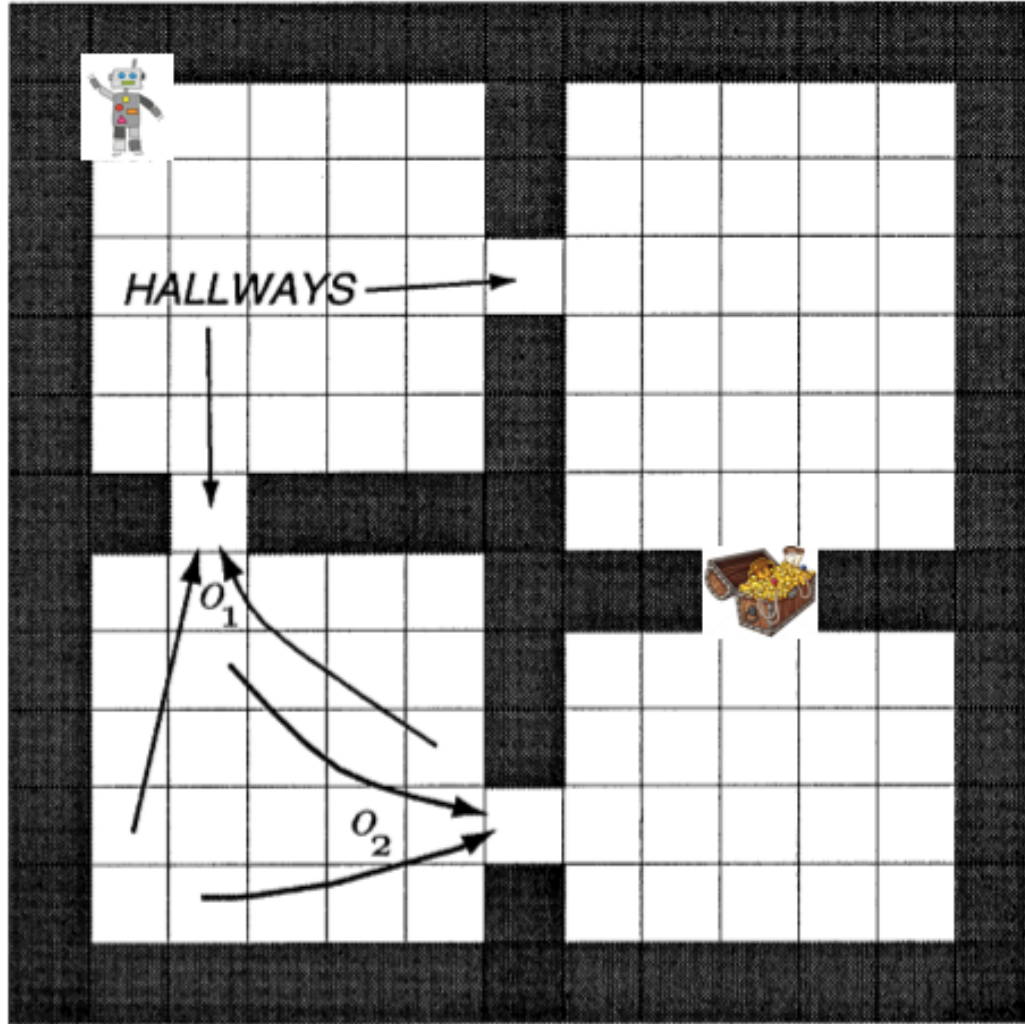
$$\mathcal{R}_s^a = E[R_{t+1} \mid S_t = s, A_t = a]$$

Bellman equations

$$\begin{aligned} V^\pi(s) &= E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots \mid S_t = s] \\ &= E[R_{t+1} + \gamma V^\pi(S_{t+1})] \\ &= \sum_a \pi(a|s) [\mathcal{R}_s^a + \gamma \sum_{s'} \mathcal{P}_{ss'}^a V^\pi(s')] \end{aligned}$$

$$\begin{aligned} V^*(s) &= \max_{\pi} V_\pi(s) \\ &= \max_a E[R_{t+1} + \gamma V^*(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \max_a [\mathcal{R}_s^a + \gamma \sum_{s'} \mathcal{P}_{ss'}^a V^*(s')] \end{aligned}$$

What is temporal abstraction?



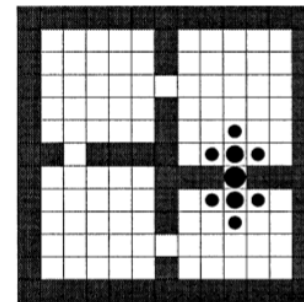
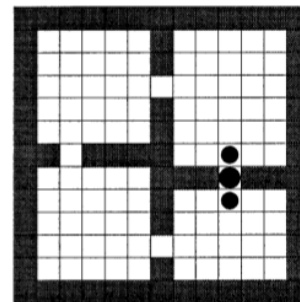
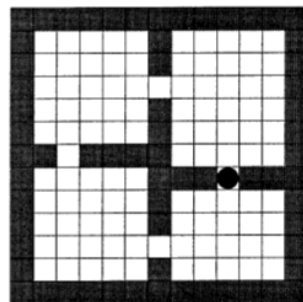
$$o = (\pi, \beta, I)$$

$$\begin{cases} \pi : S \times A \rightarrow [0, 1] \\ \beta : S \rightarrow [0, 1] \\ I \subseteq S \end{cases}$$

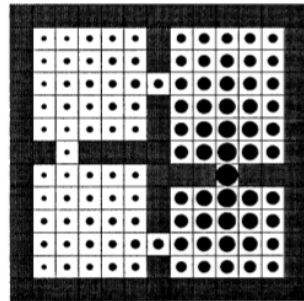
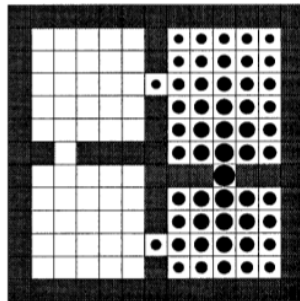
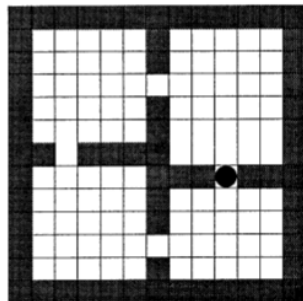
$$\mu(o \mid s) = P(O_t = o \mid S_t = s)$$

Planning with options

Primitive
options
 $\mathcal{O}=\mathcal{A}$



Hallway
options
 $\mathcal{O}=\mathcal{H}$



Initial Values

Iteration #1

Iteration #2

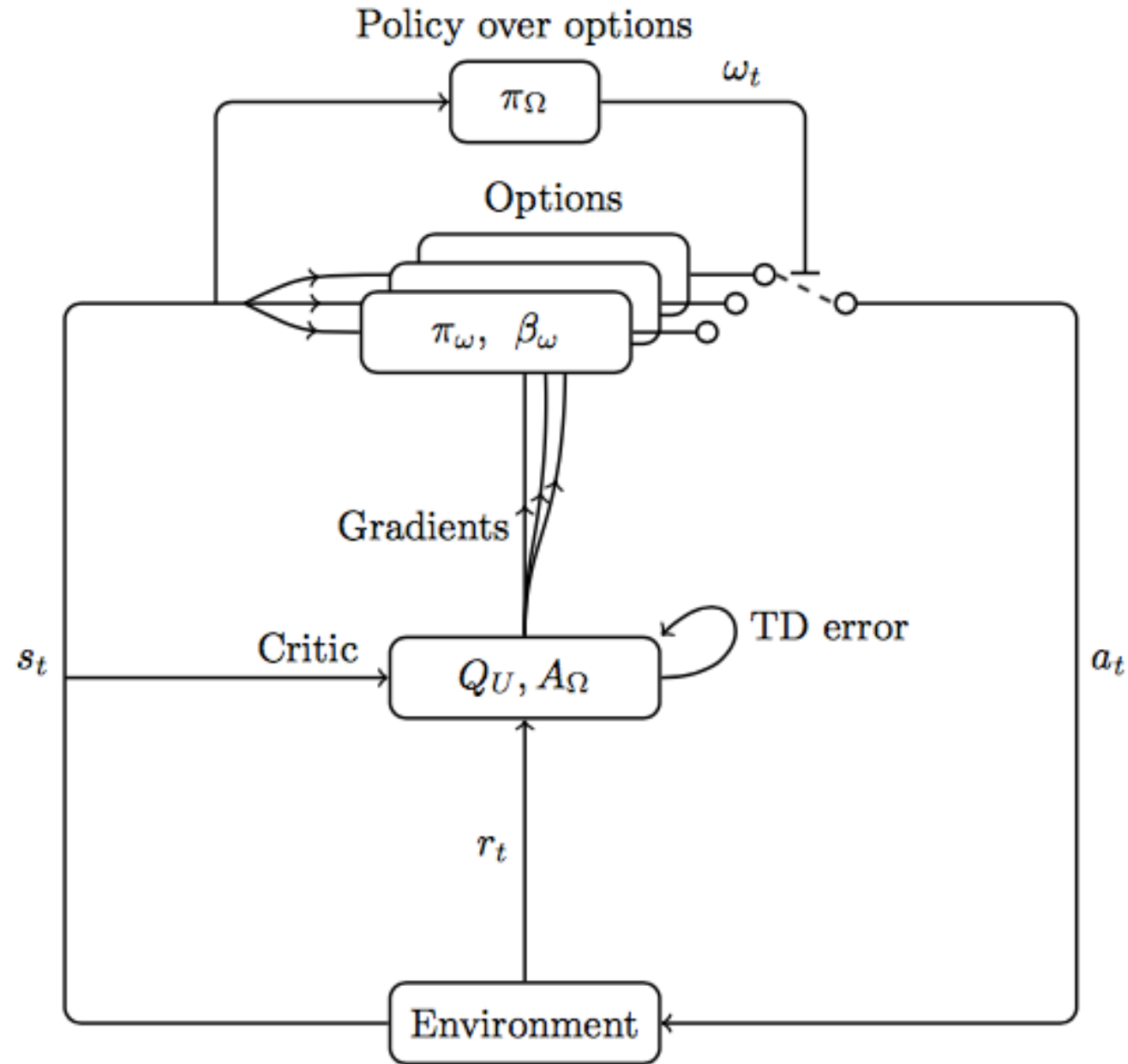
$$\mathcal{R}_s^o = E[R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{k-1} R_{t+k} \mid o \text{ initiated in state } s \text{ at time step } t]$$

$$\mathcal{P}_{ss'}^o = \sum_{k=1}^{\infty} \gamma^k P\{o \text{ terminates in } s' \text{ after } k \text{ steps}\}$$

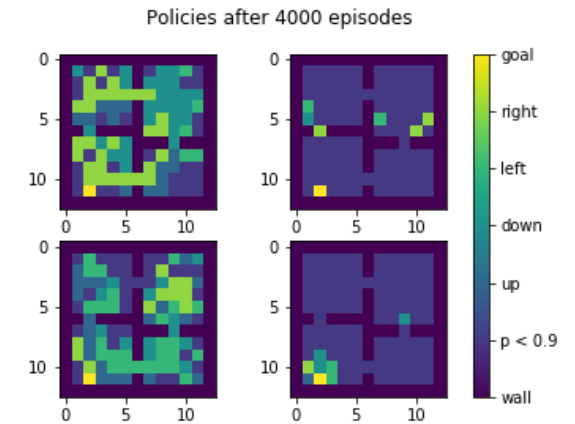
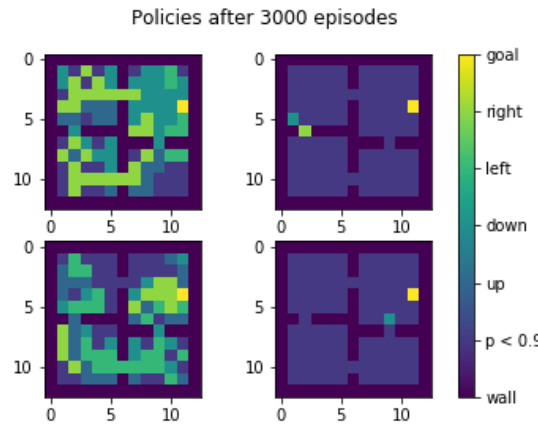
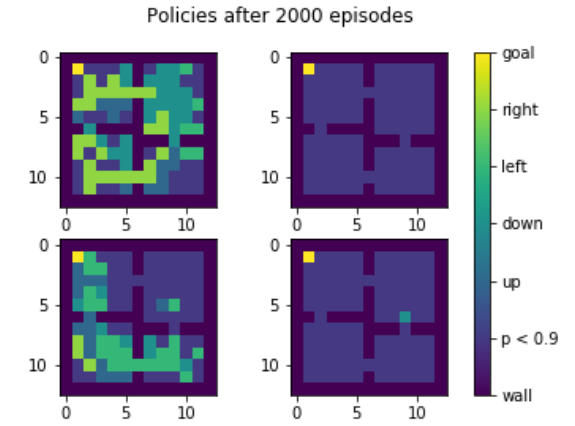
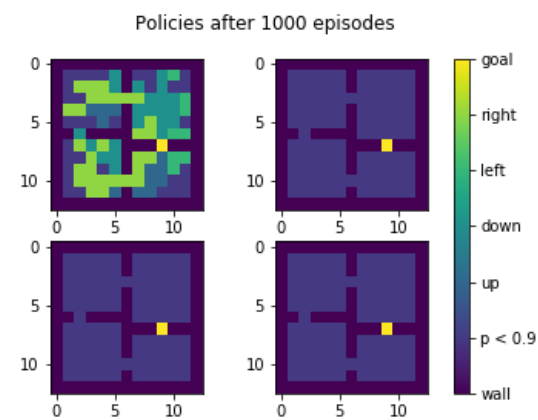
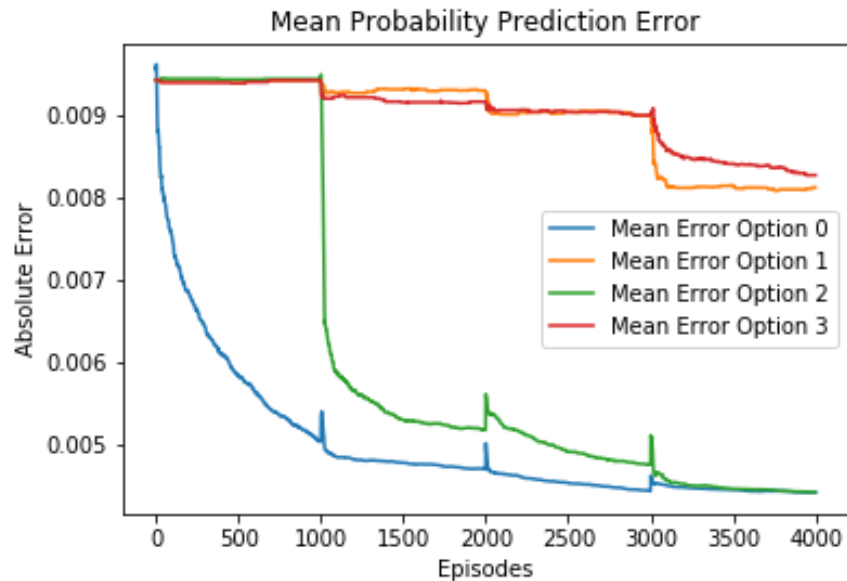
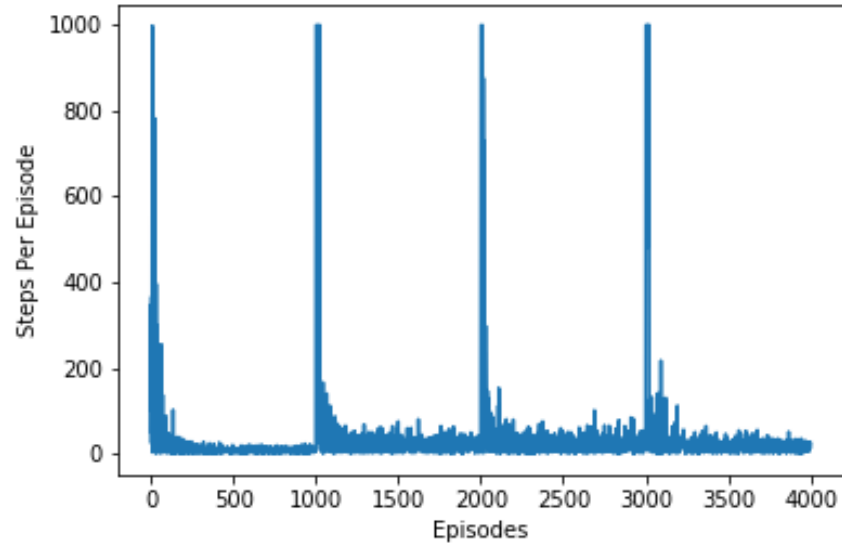
$$\begin{aligned} V^\mu(s) &= E[R_{t+1} + \dots + \gamma^{k-1} R_{t+k} + \gamma^k V^\mu(S_{t+k})] \\ &= \sum_o \mu(o \mid s) [\mathcal{R}_s^o + \sum_{s'} \mathcal{P}_{ss'}^o V^\mu(s')] \end{aligned}$$

Option-critic

- Learn options directly from data
- Policy gradient algorithm
- Similar to Actor-Critic



Transfer learning + model learning experiments



References

- *Between MDPs and Semi-MDPs: A framework for temporal abstraction in reinforcement learning*, Richard S. Sutton, Doina Precup, Satinder Singh – Artificial Intelligence, 1999
- *The Option-Critic Architecture*, Pierre-Luc Bacon, Jean Harb, Doina Precup - JMLR 2016