

MDPs with discounted cost as
MDPs with finite random duration

“Death and Discounting”
by Adam Shwartz

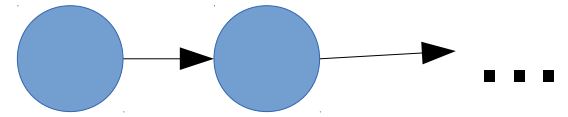
COMP 767 - Reinforcement Learning
January 20th

~Nicolas Angelard-Gontier

- They show that:
 - 1) MDPs with discounted cost \sim MDPs with finite random duration. Discount factor \sim life-span
 - 2) “an objective function which is a linear combination of several discounted costs does NOT, in general, model processes with several time scales”
- Focus on 1)

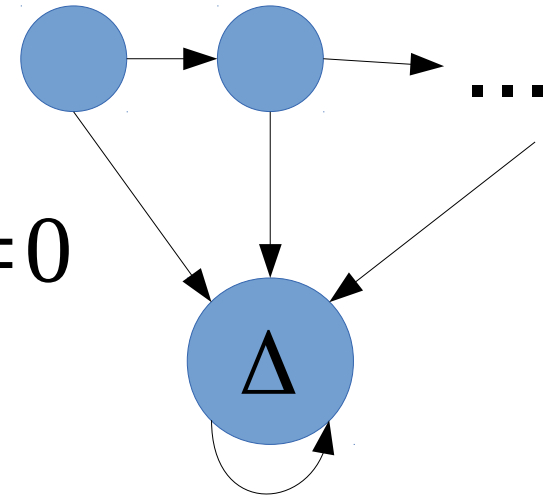
- M1: Classical discounted MDP with

$$V(x; \pi) = E_x^\pi \sum_{n=0}^{\infty} \beta^n c(x_n, a_n)$$



- M2: Now add an absorbing state Δ ('cemetery') with:

- $p(\Delta|x, a) = 1 - \beta$
- $p(y|x, a) = \beta p(y|x, a)$
- $p(\Delta|\Delta, a) = 1$ and $c_\Delta(\Delta, a) = 0$



$$V_\Delta(x; \pi_2) = E_x^{\tilde{\pi}_2} \sum_{n=0}^{\infty} c_\Delta(x_n, a_n)$$

- Let $\{\xi_n\}$ be a sequence of $\{0,1\}$ i.i.d. rnd.var. with mean β so that $x_n = \Delta \Leftrightarrow \xi_n = 0$.
 T = first time step t when $\xi_t = 0$

- Let $\{\zeta_n\}$ be a sequence of $\{0,1\}$ i.i.d. rnd.var. with mean β so that $x_n = \Delta \Leftrightarrow \zeta_n = 0$.
 T = first time step t when $\zeta_t = 0$
- Lemma: Can replace discount factor of M1 by geometric time-horizon:

$$\begin{aligned}
 V(x; \pi) &= E_x^\pi \sum_{n=0}^{\infty} \beta^n c(x_n, a_n) \\
 &= E_x^\pi \sum_{n=0}^{\infty} \left(\prod_{t=0}^{n-1} \zeta_t \right) c(x_n, a_n) \\
 &= E_x^\pi \sum_{n=0}^{T-1} c(x_n, a_n)
 \end{aligned}$$

- Policies in M2 rely on richer info:

$$\tilde{h}_n = x_0 \zeta_0 a_0 \dots x_{n-1} \zeta_{n-1} a_{n-1} x_n$$

- Given a policy π_2 in M2, we can define π in M1:

$$- \pi(.|x_0 a_0 \dots x_n) = \pi_2(.|x_0 1 a_0 \dots 1 x_n)$$

- Conversely, given π , we can define π_2 :

$$- \text{if } \zeta_t = 1 \forall t < n \text{ then } \pi_2(.|\tilde{h}_n) = \pi(.|h_n)$$

$$- \text{else } \pi_2(.|\tilde{h}_n) = a_\Delta$$

- Thm: M1 and M2 are equivalent in the following sense:

$$\begin{aligned} - E_x^\pi \beta^n c(x_n, a_n) &= E_x^\pi \left(\prod_{t=0}^{n-1} \zeta_t \right) c(x_n, a_n) \\ &= \tilde{E}_x^{\pi_2} \left(\prod_{t=0}^{n-1} \zeta_t \right) c_\Delta(x_n, a_n) = \tilde{E}_x^{\pi_2} c_\Delta(x_n, a_n) \end{aligned}$$

$$\text{so: } V(x, \pi) = V_\Delta(x, \pi_2)$$

