

Recoverability and Safe Policies

Let T be the recall time

s_0 is the desired home state

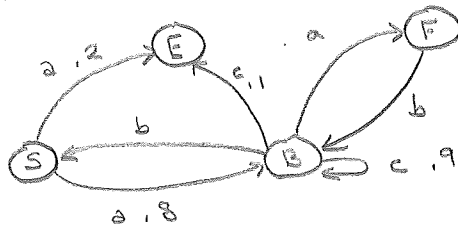
def: π_0 is δ -safe wrt s_0 & stopping time T iff

$$\exists \pi_r \text{ s.t. } E_{\pi} E_{s_0, \pi_0}^P [E_{s_T, \pi_r}^P [B_0]] \geq \delta$$

in general NP hard

→ Alternative ideas of safety?

1) Safe states or actions?



$\delta = 0.8$ safe π : a a b a b ...

unsafe a c c c c

2) 1 return π for each possible potential sample from belief
2 MDPs, prob of each is 0.5

see examples in paper

if total expected reward under any MDP drawn from the belief is bounded $\in [0, 1]$

then $(P_{sas'} - E_B[P_{sas'}]) E_{s', \pi, P}[V] \geq P_{sas'} - E_B[P_{sas'}]$ if $P_{sas'} - E_B[P_{sas'}] < 0$

and if cap at 0 then always an underestimate of term

thm 2: let B be a belief st. for any policy π and any starting state s , the total expected reward in any MDP drawn from the belief is between 0 and 1; ie. $0 \leq E_{s, \pi}^P[V] \leq 1$, B almost surely. Then the following bound holds for any policy π & starting state s

$$E_B E_{s, \pi}^P \sum_{t=0}^{\infty} R_{s_t, A_t} \geq E_{s, \pi}^P \sum_{t=0}^{\infty} (E_B[R_{s_t, A_t}] + \sigma_{s_t, A_t}^B) \text{ where}$$

$$\sigma_{s, a}^B \equiv \sum_{s'} E_B[\min(0, P_{sas'} - E_B[P_{sas'}])]$$

lemma 3: given belief B & policy π , there exists a policy dependent correction $\sigma^{B, \pi}$ such that the MDP w/ transition measure $p \equiv E_B P$ and rewards $r + \sigma^{B, \pi}$ where $r \equiv E_B R$ has the same expected total return as the belief for any initial distribution. Formally

$$\forall p \quad E_B E_{p, \pi}^P \sum_{t=0}^{\infty} R_{s_t, A_t} = E_{p, \pi}^P (r_{s, a} + \sigma_{s, a}^{B, \pi})$$

$$\sigma_{s, a}^{B, \pi} \equiv \sum_{s'} E_B [(P_{s, a, s'} - E_B [P_{s, a, s'}]) E_{s', \pi, p} [V]]$$

proof: Markov property under belief B reads

$$(1A) \quad E_B E_{s, \pi}^P [V] = \sum_a \pi_{s, a} \underbrace{E_B [R_{s, a}]}_{\equiv r_{s, a}} + \sum_a \pi_{s, a} \sum_{s'} E_B [P_{s, a, s'} E_{s', \pi}^P [V]]$$

Markov property assuming expected transition frequencies & expected rewards with corrections is

$$(1B) \quad E_{s, \pi}^P [\bar{V}] = \sum_a \pi_{s, a} (r_{s, a} + \sigma_{s, a}^{B, \pi}) + \sum_a \pi_{s, a} \sum_{s'} p_{s, a, s'} E_{s', \pi}^P [\bar{V}]$$

$$\text{let } \Delta_s \equiv E_B E_{s, \pi}^P [V] - E_{s, \pi}^P [\bar{V}]$$

subtract, 1A - 1B

$$\Delta_s = - \sum_a \pi_{s, a} \sigma_{s, a}^{B, \pi} + \sum_a \pi_{s, a} \sum_{s'} (E_B [P_{s, a, s'} E_{s', \pi}^P [V]] - p_{s, a, s'} E_{s', \pi}^P [\bar{V}])$$

where used defn of $r_{s, a}$

$$= \sum_a \pi_{s, a} \sum_{s'} E_B [\underbrace{(E_B [P_{s, a, s'}] - p_{s, a, s'})}_{\equiv P_{s, a, s'}^{(1)}} E_{s', \pi, p}^{(2)} [V]] \quad \text{sub in defn of } \sigma$$

$$+ \sum_a \pi_{s, a} \sum_{s'} (E_B [P_{s, a, s'} E_{s', \pi, p}^{(u)} [V]] - p_{s, a, s'} E_{s', \pi}^P [\bar{V}]) \quad (2)$$

$$= \sum_a \pi_{s, a} \sum_{s'} p_{s, a, s'} E_B E_{s', \pi, p} [V] \quad p_{s, a, s'} \text{ is indep of } B$$

$$- \sum_a \pi_{s, a} \sum_{s'} p_{s, a, s'} E_{s', \pi}^P [\bar{V}] \quad \text{and (2) \& (5) cancel}$$

$$= \sum_a \pi_{s, a} \sum_{s'} p_{s, a, s'} \Delta_{s'} \quad \text{defn of } \Delta_{s'}$$

Δ satisfies the same eqn as a value function in a MDP w/ transition measure p and zero rewards. Since the value func in such an MDP is uniquely defined & identically 0, $\Delta_s = 0$. \square

but in general don't know V & transformation requires V
if have a bound on V , can use to get a lower bound

$$\max_{\pi_0, \pi_r} E_{s_0, \pi_0}^{\pi_r} \sum_t (r_{s_t, a_t} + \gamma V_{s_t, a_t}^B) \quad (5A)$$

$$\text{subject to } E_{s_0, \pi_0}^{\pi_r} \sum_{t=0}^{\infty} (1-\gamma) V_{s_t} + \gamma \sigma_{s_t, a_t}^B \geq \delta \quad (5B)$$

$$\text{and } V_s = E_{s, \pi_r}^{p, (1-1_{s=s_0})} \sum_{t=0}^{\infty} (1_{s_t=s_0} + (1-1_{s_t=s_0}) \sigma_{s_t, a_t}^B) \quad (5C)$$

value func of MDP with transition measure $p(1-1_{s=s_0})$ and reward $1_{s=s_0} + (1-1_{s=s_0}) \sigma_{s, a}^B$ under policy π_r

The return policy π_r only appears in Equation 5C.

Therefore the above optimization can be done in 2 steps

Step 1: solve for optimal return policy π_r^* given expected transition measure p & also compute optimal value function V_s^*

$$V_s^* = \max_{\pi_r} E_{s, \pi_r}^{p, (1-1_{s=s_0})} \sum_{t=0}^{\infty} (1_{s_t=s_0} + (1-1_{s_t=s_0}) \sigma_{s_t, a_t}^B)$$

Step 2: compute optimal exploration policy π_0^* under strict safety constraint by solving a constrained MDP

$$\max_{\pi_0} E_{s_0, \pi_0}^{\pi_r} \sum_t (r_{s_t, a_t} + \gamma V_{s_t}^B)$$

$$\text{s.t. } E_{s_0, \pi_0}^{\pi_r} \sum_t ((1-\gamma) V_{s_t}^* + \gamma \sigma_{s_t, a_t}^B) \geq \delta$$

they solved as a linear program

grid world experiment

exploration: adapted R-mex

safety costs

$$\sigma_{s, a}^B = -2 E_p [P_{s, a}] (1 - E_p [P_{s, a}]) = -2 \text{Var}_p [P_{s, a}]$$