

learning from an expert:

①

- I have an MDP
- with UNKNOWN reward
- I see an expert acting.
- Q: learn policy

Versions

① Abel + Ng

assume $R(s) = \sum_i w_i \phi_i(s) = w \cdot \Phi(s)$

↑ feature functions

now any given policy π has a value

$$E_{s_0 \sim D} [V^\pi(s_0)] = E \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \pi \right]$$
$$= w \cdot E \left[\sum_{t=0}^{\infty} \gamma^t \Phi(s_t) \mid \pi \right]$$

So

$$\mu(\pi) = E \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t) \mid \pi \right]$$

↖ determines value of policy
feature expectations.

When we have observations of an expert acting, we have state trajectories (2)

we know

$$\hat{\mu}_E = \frac{1}{M} \sum_{i=1}^M \sum_{t=0}^{\infty} \gamma^t \phi(s_t^{(i)})$$

estimates of feature expects for optimal policy.

strategy: search for policy π^* st.

$$\mu(\pi^*) \approx \hat{\mu}_E$$

Notice, this means

$$\max_{w; \|w\| \leq 1} w_0 (\hat{\mu}_E - \mu(\pi^*)) \text{ small}$$

so search for π^*

$$\text{st } \max_{w; \|w\| \leq 1} [w_0 (\hat{\mu}_E - \mu(\pi^*))] \text{ is small}$$

Search process :

- start w / π^0

- obtain

$$E^{(i)} = \max_{w, \|w\| \leq 1} \left[\min_{j \in 0 \dots i-1} \left\{ w \cdot \left(\mu_E^j - \mu(\pi^i) \right) \right\} \right]$$

$w^{(i)}$ = the w that achieves this

w max difference between expert and best policy so far

- if $E^{(i)} < \epsilon$

- now compute optimal policy for MDP using $R = (w^{(i)})^T \phi$

[known state transition probs — policy/value iteration]

[unknown — RL]

- estimate $\mu^{(i)}$

④

inner loop solves

$$\max_{t, w} t$$

$$w^T \mu_E \geq w^T \mu^{(j)} + t$$

$$\|w\|_2 \leq 1$$

quadratic program