

# Behavior Cloning

Setting: Discounted infinite horizon MDP.

$$\mathcal{M} = \{ \mathcal{S}, \mathcal{A}, P, R, \gamma \}$$

unknown  
unobserved

Expert knows optimal policy  $\pi^*$  and dataset

$$\mathcal{D} = \{ (s_i^*, a_i^*) \}_{i=1}^M \sim d_{\pi^*}^{\gamma_0}$$

## Behaviour Cloning

Reduction to supervised learning.

Define a policy class,  $\pi \in \Pi$

e.g. parametric policies  $\Pi = \{ \pi_{\theta} : \theta \in \mathbb{R}^d \}$

Then we estimate a policy with empirical risk minimization

$$\hat{\pi} = \underset{\pi \in \Pi}{\operatorname{argmin}} \sum_{i=1}^M \ell(\pi, s_i^*, a_i^*)$$

Many choices of loss function

- discrete action space: classification
- continuous action space: regression

ex: Negative log likelihood

$$\ell(\pi, s, a) = -\log(\pi(a|s))$$

ex: square loss

$$\ell(\pi, s, a) = \| \pi(s) - a \|_2^2$$

### 3) BC Analysis

Assumption: Supervised ML is successful

$$\mathbb{E}_{s \sim d_{y_0}^{\pi^*}} [\mathbb{1}\{\hat{\pi}(s) \neq \pi^*(s)\}] \leq \varepsilon$$

depend on  
# datapoints,  
complexity of  $\pi$ ,  
whether  $\pi^* \in \Pi$

Train ( $d^{\pi^*}$ ) and Test ( $d^{\hat{\pi}}$ )  
distribution mismatch.

$$(1-\gamma) \left[ 0 + \sum_{t=1}^{\infty} \gamma^t \right] = (1-\gamma) \frac{\gamma}{1-\gamma} = \gamma$$

Recall:  $d_{y_0}^{\pi}(s) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t P_t^{\pi}(s; y_0)$

Recall: Performance Difference Lemma

$$\mathbb{E}_{s \sim y_0} [V^{\pi}(s) - V^{\pi'}(s)] = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{y_0}^{\pi}} \left[ \underbrace{\mathbb{E}_{a \sim \pi(s)} [A^{\pi'}(s, a)]}_{\text{advantage of } \pi \text{ over } \pi'} \right]$$

difference in value  
(cumulative reward)

advantage of  $\pi$  over  $\pi'$

Theorem (BC Performance): Assume  $r(s, a) \in [0, 1]$   
and supervised ML succeeds w/  $\varepsilon$ . Then BC returns  $\hat{\pi}$

$$\mathbb{E}_{s \sim y_0} [V^{\pi^*}(s) - V^{\hat{\pi}}(s)] \leq \frac{2\varepsilon}{(1-\gamma)^2} \quad \left( \sum_{t=0}^{\infty} \gamma^t \leq \frac{1}{1-\gamma} \right)$$

Proof

$$\mathbb{E}_{s \sim y_0} [V^{\pi^*}(s) - V^{\hat{\pi}}(s)] = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{y_0}^{\pi^*}} [A^{\hat{\pi}}(s, \pi^*(s))] \quad \leftarrow \text{deterministic}$$

add 0

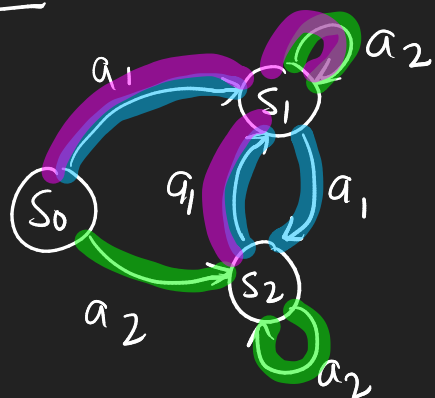
$$= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{y_0}^{\pi^*}} [A^{\hat{\pi}}(s, \pi^*(s)) - A^{\hat{\pi}}(s, \hat{\pi}(s))]$$

$$A \leq \frac{1}{1-\gamma}$$

$$\leq \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\pi^*}} \left[ \frac{2}{1-\gamma} \mathbb{1} \{ \pi^*(s) \neq \hat{\pi}(s) \} \right]$$

$$\leq \frac{2}{(1-\gamma)^2} \cdot \varepsilon$$

### Example: Distribution Shift



$s_0$  is initial state

$$r(s,a) = \begin{cases} 1 & s=s_1 \\ 0 & \text{otherwise} \end{cases}$$

Optimal policy

$$d_{s_0}^{\pi^*}(s) = \begin{cases} 1-\gamma & s=s_0 \\ \gamma & s=s_1 \\ 0 & s=s_2 \end{cases}$$

consider policy  $\hat{\pi}(s_0) = \begin{cases} a_1 & \text{w.p. } 1-\frac{\varepsilon}{1-\gamma} \\ a_2 & \text{w.p. } \frac{\varepsilon}{1-\gamma} \end{cases}$   $\hat{\pi}(s_1) = a_2$   $\hat{\pi}(s_2) = a_2$

Has low SL error

$$\mathbb{E}_{s \sim d_{s_0}^{\hat{\pi}}} \left[ \mathbb{E}_{a \sim \hat{\pi}(s)} [a \neq \pi^*(s)] \right] = \varepsilon$$

Has quadratic performance error

$$V^{\pi^*}(s_0) = \sum_{t=1}^{\infty} \gamma^t = \frac{\gamma}{1-\gamma} \quad V^{\hat{\pi}}(s_0) = \left(1 - \frac{\varepsilon}{1-\gamma}\right) \frac{\gamma}{1-\gamma} + \left(\frac{\varepsilon}{1-\gamma}\right) \cdot 0$$

$$V^{\pi^*}(s_0) - V^{\hat{\pi}}(s_0) = \frac{\varepsilon \gamma}{(1-\gamma)^2}$$

Caused by bad policy in  $s_2 \rightarrow d_{s_0}^{\pi^*}(s_2) = 0$