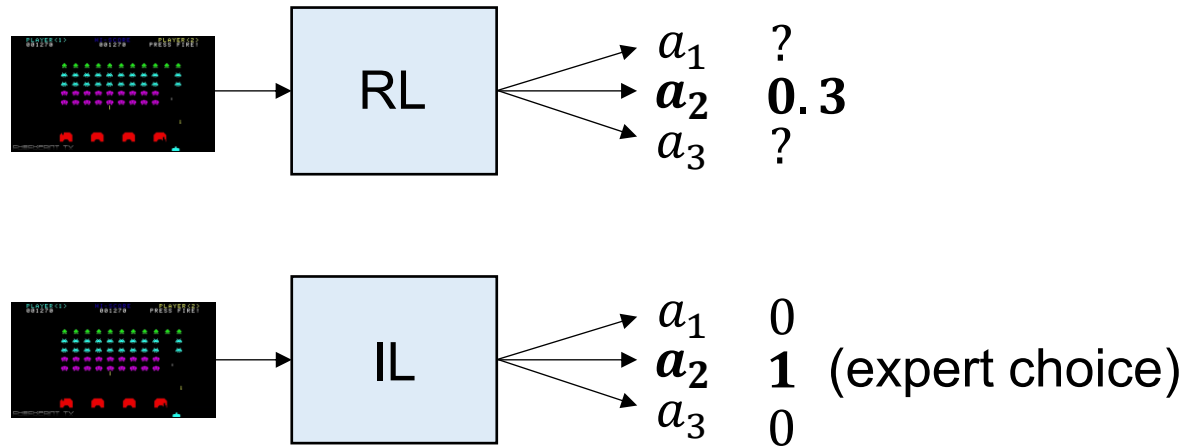


Imitation Learning

Chen-Yu Wei

Imitation Learning \in Supervised Learning



π^* : expert's policy

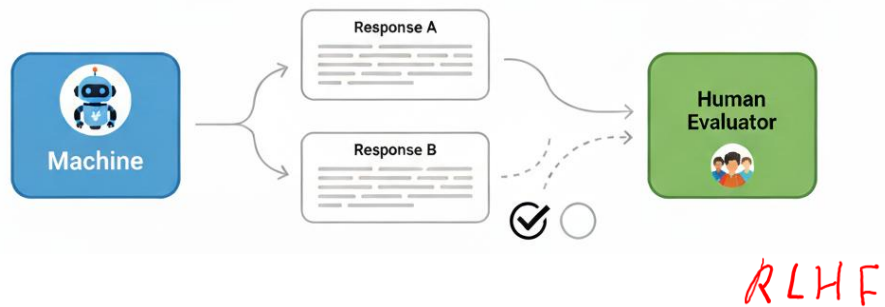
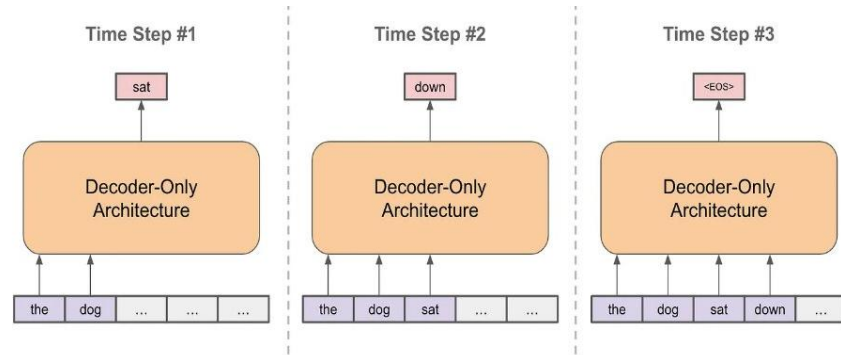
Offline IL: learn from static data generated by the expert $\{(s_1, a_1^*, s_2, a_2^*, \dots, s_H, a_H^*)\}$

Online IL: may interact with MDP and query the expert $\{(s_1, a_1, a_1^*, s_2, a_2, a_2^*, \dots, s_H, a_H, a_H^*)\}$

Goal: output a policy $\hat{\pi}$ such that $V^{\pi^*}(\rho) - V^{\hat{\pi}}(\rho)$ is small

Examples

- Language models



- Robotics



Types



$$\pi(a|s) \leftrightarrow \pi^*(a|s)$$

- Direct Imitation: directly learn policy to imitate the expert

- Behavior cloning
- DAgger
- Direct preference optimization (preference feedback)

- Occupancy matching

- DICE
- GAIL

$$d^{\pi}(s,a) \leftrightarrow d^{\pi^*}(s,a)$$

$$(s_1, a_1^*, s_2, a_2^*, \dots)$$

- Inverse RL: learn an MDP (or just reward function) from expert, and perform RL on it

- Adversarial IRL ([paper](#))
- MaxEnt IRL ([paper](#))
- RLHF (preference feedback)



Behavior Cloning: Reduction to Classification

$(s_1, a_1^*, s_2, a_2^*, \dots)$ are from expert

Relate $V^{\pi^*}(\rho) - V^{\hat{\pi}}(\rho)$ to $\mathbb{E}^{\pi^*} \left[\frac{1}{H} \sum_{h=1}^H \mathbb{I}\{\hat{\pi}(s_h) \neq a_h^*\} \right]$

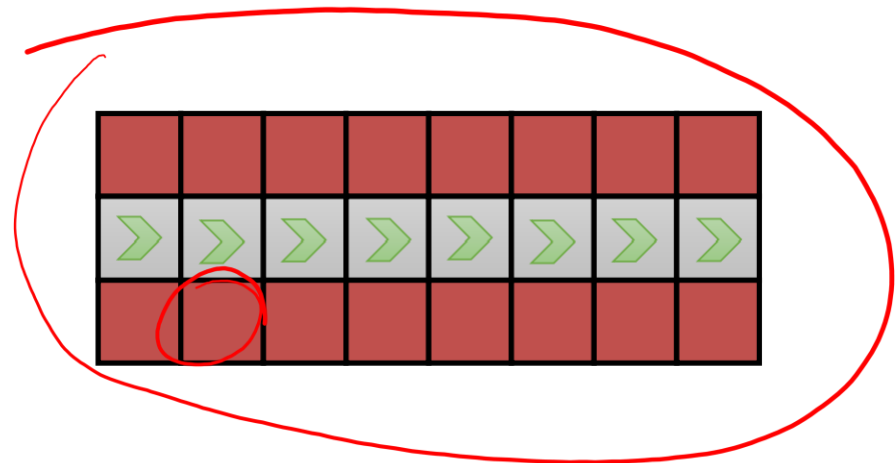
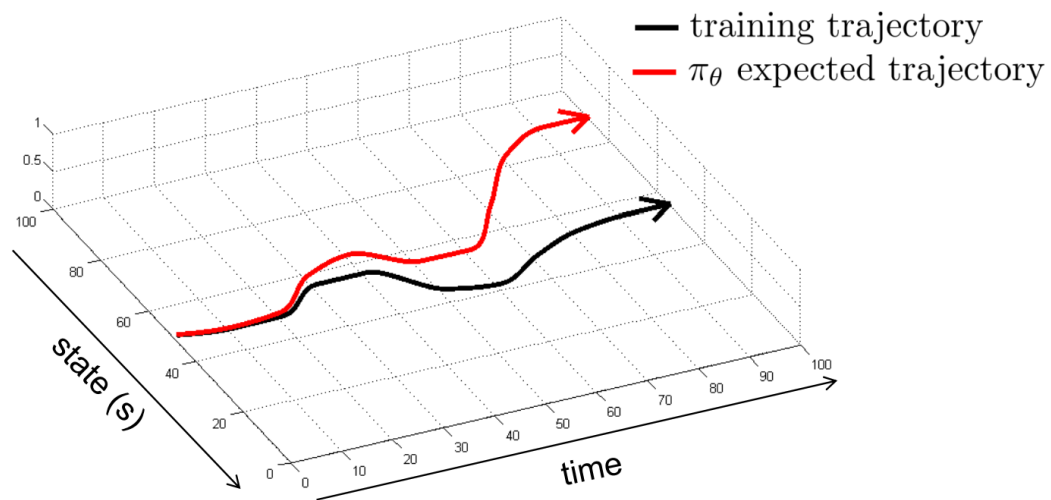
$$V^{\pi^*}(\rho) - V^{\hat{\pi}}(\rho) = \sum_h \sum_{s,a} d^{\pi^*}(s) \left(\pi^*(a|s) - \pi(a|s) \right) \underbrace{\left(Q^{\pi^*}_{(s,a)} \right)}_{\leq H}$$

$$\leq 2H \sum_h \sum_s d^{\pi^*}(s) \mathbb{I}\{ \pi^*(s) \neq \pi(s) \}$$

$$\leq \underbrace{O(H^2)}_{\text{indicator loss}} \cdot \underbrace{\frac{1}{H} \sum_h d^{\pi^*}(s) \mathbb{I}\{ \pi^*(s) \neq \pi(s) \}}_{\text{indicator loss}}$$

indicator loss

Behavior Cloning: Reduction to Classification



Issue: ~~distribution shift~~

$$\frac{1}{H} \sum_h \mathbb{I}(z_h(s_h) \neq \pi(s_h)) \leq \epsilon$$

$$\text{total loss} = \sum_h \mathbb{I}(\text{off-track at step } h)$$

$$\begin{aligned} &= \epsilon \times \mathbb{I}\{\text{go off-track at 1st step}\} \times H \\ &+ (1-\epsilon) \epsilon \mathbb{I}\{\text{2nd step}\} \times (H-1) \\ &+ (1-\epsilon)^2 \epsilon \times (H-2) \\ &\vdots \end{aligned} \left. \begin{aligned} &\approx \epsilon H \\ &\approx \epsilon H \\ &\approx \epsilon H \end{aligned} \right\} \epsilon H^2$$

Behavior Cloning: Reduction to Classification

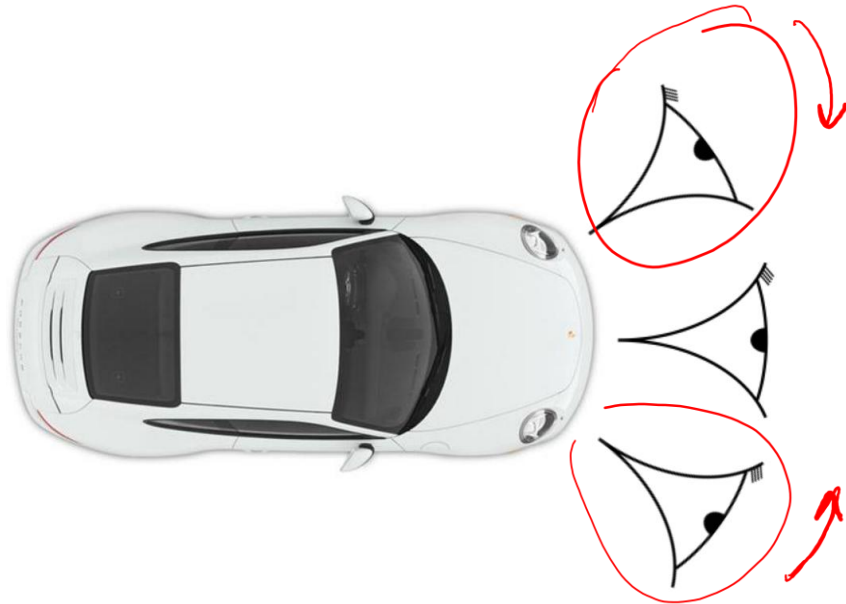
The bound might be pessimistic

- Single mistake may not lead to catastrophic failure



Solution

- Data augmentation



Bojarsky et al. End to End Learning for Self-Driving Cars. 2016


Solution: Interact with Expert (Online IL)

DAgger: Dataset Aggregation

goal: collect training data from $p_{\pi_\theta}(\mathbf{o}_t)$ instead of $p_{\text{data}}(\mathbf{o}_t)$

how? just run $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$

but need labels \mathbf{a}_t !

- 
1. train $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
 2. run $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ to get dataset $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
 3. Ask human to label \mathcal{D}_π with actions \mathbf{a}_t
 4. Aggregate: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$