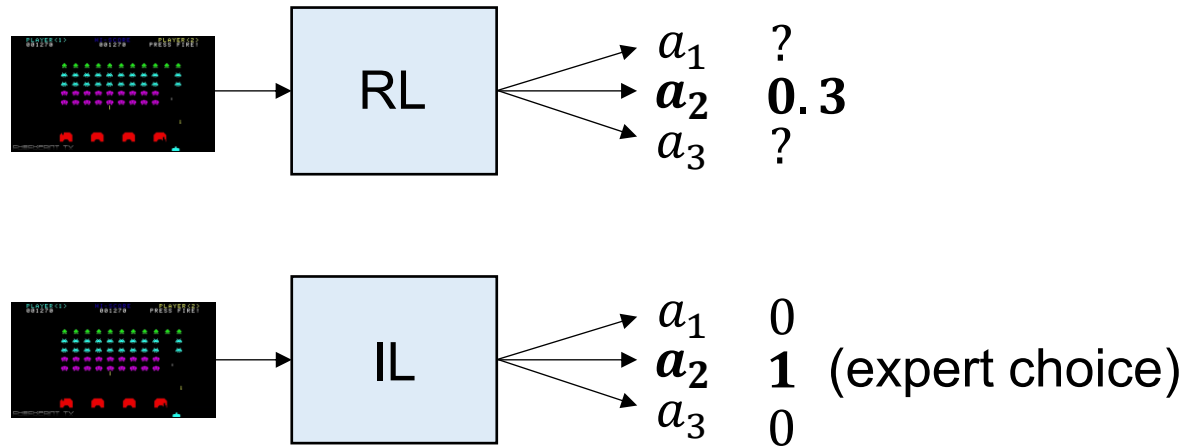


# Imitation Learning

Chen-Yu Wei

# Imitation Learning $\in$ Supervised Learning



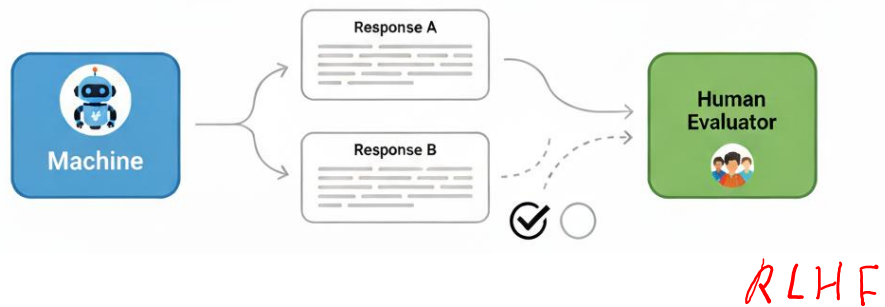
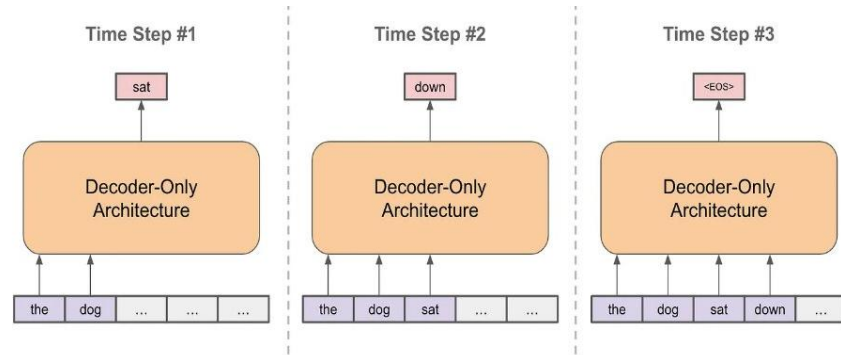
Offline IL: learn from static data generated by the expert  $\{(s_1, a_1^*, s_2, a_2^*, \dots, s_H, a_H^*)\}$

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

- Direct Imitation: directly learn policy to imitate the expert
  - Behavior cloning
  - DAgger
  - Direct preference optimization (preference feedback)
- Inverse RL: learn reward function from expert, and perform RL on it
  - Adversarial IRL ([paper](#))
  - MaxEnt IRL ([paper](#))
  - RLHF (preference feedback)

# **Direct Imitation**

# Behavior Cloning: Reduction to Classification

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

$\text{Range}(Q) =$

$$V^{\pi^*}(\rho) - V^{\hat{\pi}}(\rho) = \sum_{h=1}^H \sum_s d_h^{\pi^*}(s) (Q_h^{\hat{\pi}}(s, \pi_h^*(s)) - Q_h^{\hat{\pi}}(s, \hat{\pi}_h(s)))$$

$\max_{s, a, b} |Q(s, a) - Q(s, b)|$

$$\leq \text{Range}(Q^{\hat{\pi}}) \sum_{h=1}^H \sum_s d_h^{\pi^*}(s) \mathbb{I}\{\pi_h^*(s) \neq \hat{\pi}_h(s)\}$$

$$= H \text{Range}(Q^{\hat{\pi}}) \times \frac{1}{H} \sum_{h=1}^H \sum_s d_h^{\pi^*}(s) \mathbb{I}\{\pi_h^*(s) \neq \hat{\pi}_h(s)\}$$

# Behavior Cloning with Logistic Loss

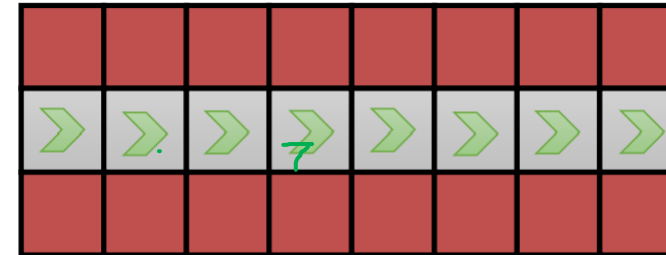
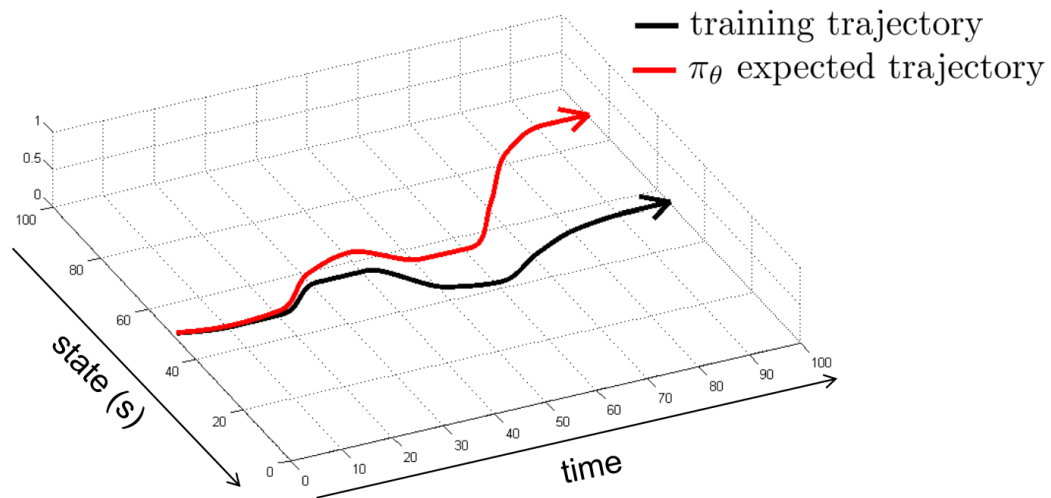
## Behavior Cloning (Offline IL)

Run expert policy  $\pi^*$  and obtain  $(s_1, a_1), \dots, (s_N, a_N)$

Obtain policy  $\pi_\theta$  by finding

$$\theta = \operatorname{argmin}_{\theta} \frac{1}{N} \sum_{i=1}^N \log \frac{1}{\pi_{\theta}(a_i | s_i)}$$

# Behavior Cloning: Reduction to Classification



Issue: **distribution shift**

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

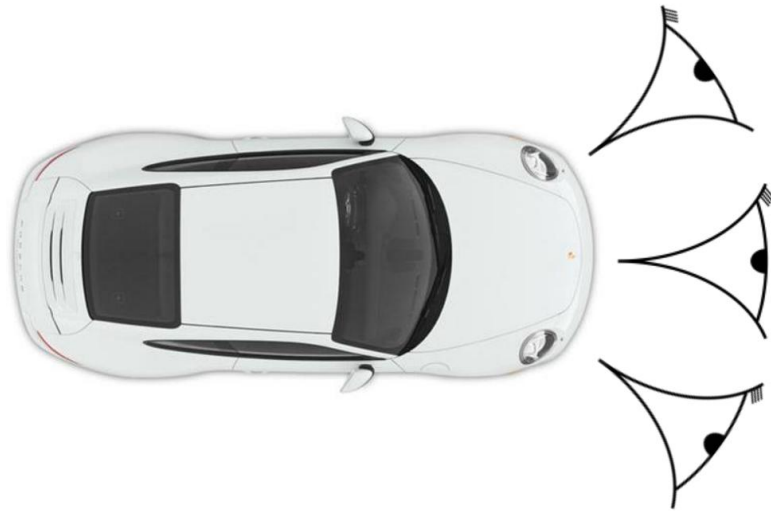
$$\text{total loss} = \sum_h \mathbb{I}(\text{off-track at step } h)$$

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



# Solution

- Data augmentation



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

# Solution: Interact with Expert (Online IL)

## Dataset Aggregation (DAgger)

For  $k = 1, 2, \dots$

Train  $\pi_\theta(a|s)$  with dataset  $\mathcal{B}$

Run  $\pi_\theta$  to collect states  $s_1, s_2, \dots, s_N$

Ask the expert to label actions, giving  $(s_1, a_1), \dots, (s_N, a_N)$

Add  $(s_1, a_1), \dots, (s_N, a_N)$  to  $\mathcal{B}$

# DAgger with Logistic Loss

## DAgger (Online IL)

Run expert policy  $\pi^*$  and obtain  $(s, a)$ -pairs and push them to  $\mathcal{B}$

For  $k = 1, 2, \dots$

For  $m = 1, 2, \dots, M$ :

Sample a batch  $b$  from  $\mathcal{B}$

$$\theta \leftarrow \theta - \alpha \frac{1}{|b|} \sum_{(s,a) \in b} \log \frac{1}{\pi_{\theta}(a|s)}$$

Use  $\pi_{\theta}$  to generate states  $s_1, s_2, \dots, s_N$

Ask expert to provide labels  $(s_1, a_1), (s_2, a_2), \dots, (s_N, a_N)$

Push them to  $\mathcal{B}$

# Analysis

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

$$\begin{aligned} V^{\pi^*}(\rho) - V^{\hat{\pi}}(\rho) &= \sum_{h=1}^H \sum_s d_{\hat{\pi}}^h(s) \left( Q_h^{\pi^*}(s, \pi_h^*(s)) - Q_h^{\pi^*}(s, \hat{\pi}_h(s)) \right) \\ &\leq \text{Range}(Q^{\pi^*}) \sum_{h=1}^H \sum_s d_{\hat{\pi}}^h(s) \mathbb{I}\{\pi_h^*(s) \neq \hat{\pi}_h(s)\} \\ &= \textcolor{red}{H\text{Range}(Q^{\pi^*})} \times \frac{1}{H} \sum_{h=1}^H \sum_s d_{\hat{\pi}}^h(s) \mathbb{I}\{\pi_h^*(s) \neq \hat{\pi}_h(s)\} \end{aligned}$$

# Imitation with Preference Feedback

Direct Preference Optimization (DPO) ([link](#))

Given context  $x$  and two potential actions, the expert chooses the **preferred** one (the preferred one is denoted as  $y_w$  and the other is  $y_l$ ).

$$\mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_{\theta}(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi_{\theta}(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right]$$

$\pi_{\text{ref}}$  is a reference policy that  $\pi_{\theta}$  is initialized as.

The previously discussed frameworks of online / offline IL can be directly applied here.

# Inverse RL

# Inverse Reinforcement Learning (IRL)

Find a reward / value function that explains the behavior of the expert.

- **Behavior assumption:** Assume some connections between expert policy and reward/value. E.g.,  
$$\pi^*(a|s) \propto \exp(Q^*(s, a))$$
$$P^{\pi^*}(\tau) \propto \exp(\sum_h R^*(s_h, a_h))$$
- **Maximum likelihood:** Given expert trajectories, find  $\phi$  that maximizes likelihood
- **Policy training:** Find a policy that approximately maximizes the expected reward / value under  $\phi$ . E.g.,  
$$\pi(a|s) \propto \exp(Q_\phi(s, a))$$
$$\theta = \operatorname{argmax}_{\theta} \mathbb{E}_{s \sim D} [\mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [R_\phi(s, a)] - \beta \operatorname{KL}(\pi_\theta(\cdot|s), \pi_{\text{ref}}(\cdot|s))]$$

Mathematically, this “learn reward from expert + learn policy from reward” paradigm can always be condensed as “learn policy from expert”

# IRL with Preference Feedback

$$p^*(y_1 \succ y_2 \mid x) = \frac{\exp(r^*(x, y_1))}{\exp(r^*(x, y_1)) + \exp(r^*(x, y_2))} \quad \text{Behavior assumption}$$

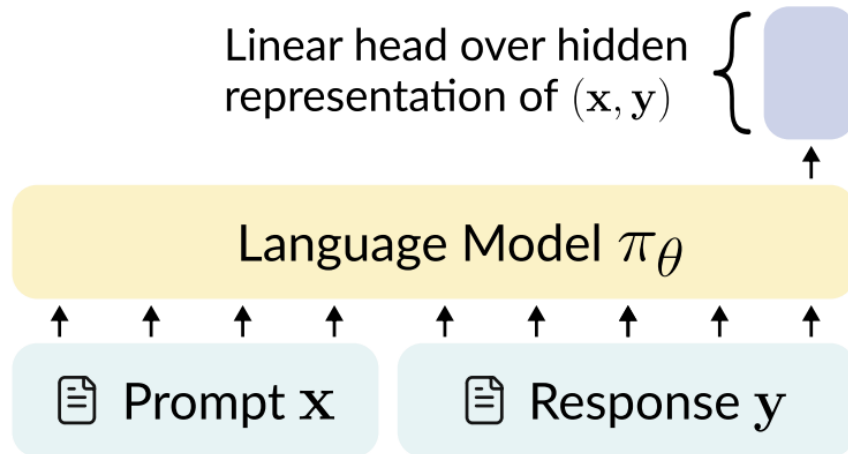
$$\mathcal{L}_R(r_\phi, \mathcal{D}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log \sigma(r_\phi(x, y_w) - r_\phi(x, y_l))] \quad \text{Maximum likelihood}$$

$$\max_{\pi_\theta} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(y|x)} [r_\phi(x, y)] - \beta \mathbb{D}_{\text{KL}}[\pi_\theta(y \mid x) \parallel \pi_{\text{ref}}(y \mid x)] \quad \text{Policy training}$$

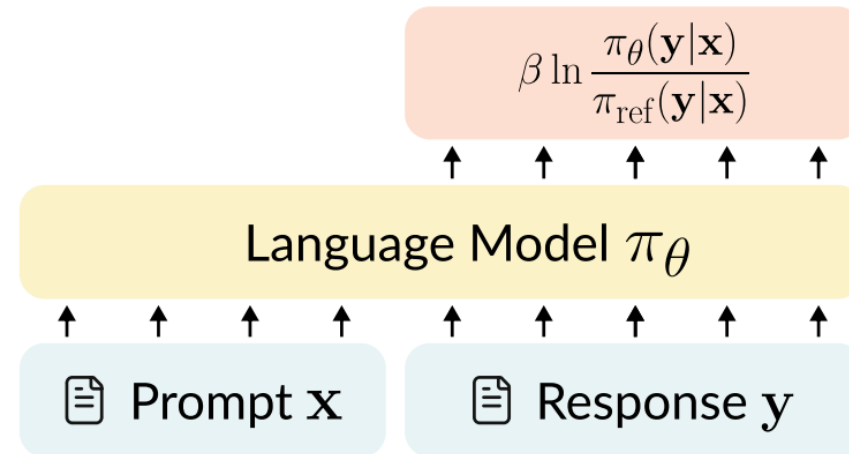


# Benefits of having an explicit reward model

## Explicit Reward Model (EX-RM)



## Implicit Reward Model (IM-RM)



Razin et al., 2025. Why is Your Language Model a Poor Implicit Reward Model?