

2024.09.09 (1) Behavioral Cloning
(2) Inverse RL.

• Imitation Learning

→ Given: Expert or expert's demonstration

→ Goal: Train a policy to mimic demonstrations

function maps
(observation) state → action

"Learn" human intents, preference, underlying reward functions²⁹



→ We need (1) demonstration (2) Environment / Simulator (3) Policy Class
or demonstrator (which function family?)
(4) Loss Function (5) Learning Algorithm

• Problem Setup: Markov Decision Process without Reward.

- State space S [can be partially observable], $s \in S$

- Action space A , $a \in A$

- An expert policy π^* : $S \rightarrow p(a|s)$ or a] top trace of π^* is Trajectory
State distribution over actions. Demonstration
 $T = (S_0, a_0, S_1, a_1, S_2, a_2, \dots)$
↳ "Roll Out":

- Transition model $P_{sa}(S_{t+1}|S_t, a_t)$: characteristic/nature of simulator or environment

Goal:

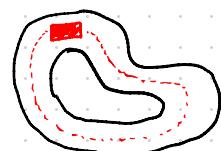
↳ Inverse RL: Can we recover reward from π^* ? (assume experts are optimal)

↳ Behavior Cloning: Can we train π_θ using supervised learning?

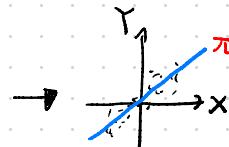
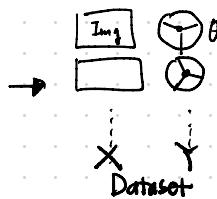
Ex) Learning to Drive by Imitation.

• Input: RGB Image, Output: Steering Angle in $[-1, 1]$

• Behavior Cloning: Offline & Supervised Imitation Learning.



Expert Trajectories.



→ Learned Policy $\pi_\theta(s) \rightarrow p(a)$ or a
(state → action)

Supervised Learning,
"Fixed Policy Class"
(i.e. SUM, NN, ...)

• Formulation of Behavior Cloning Algorithm

- Algorithm's Input : A restricted policy class Π

- Learning Objective : Train $\Pi_{\theta^*}(s)$, where $\theta^* = \arg \min_{\theta} E(s, a^*) \sim P^*$

$$\text{OR } \hat{\Pi} \doteq \arg \min_{\Pi \in \mathcal{P}} \sum_{(s, a^*) \sim P^*} L(a^*, \pi_\theta(s))$$

$$\text{OR } \hat{\Pi} \doteq \arg \min_{\Pi \in \mathcal{P}} \sum_{T \in \mathcal{T}} \sum_{(s, a^*) \in T} \frac{L(a^*, \pi_\theta(s))}{\pi_\theta(s)}$$

$$\text{General IL: minimize } E_{S \sim P(S|D)} L(a^*, \pi_\theta(s))$$

Assume perfect imitation so far,

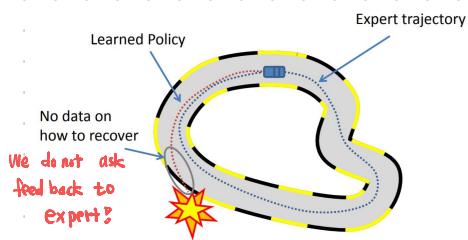
Learn to continue this perfect imitation,

Minimize 1-step deviation error along the expert traj.

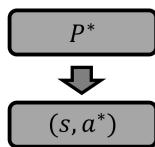
(mismatch)

* Distribution Shift in BC.

→ Data distribution mismatch



IID Assumption
(supervised learning)



Reality

$$P_0$$

$$S \xrightarrow{\quad} \pi_\theta$$

* Supervised learning assumes (S, a) are i.i.d.

(independent & identically distributed)

→ If error for each time is t , clear in time period T
is expected to bounded in GT.

<Training> $S \sim P^*$

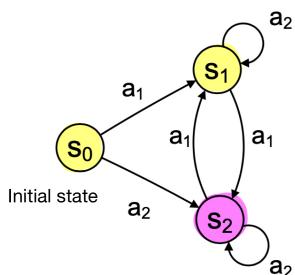
<Inference> $S \sim P(S | \pi_\theta)$: unseen states?

makes mistake & cannot recover.

∴ BC : Simple, efficient but No long-term planning.

(USE) → 1-step deviations not too bad, Learning Reactive behavior,
expert traj. cover S.

(DON'T) → 1-step deviation lead to disaster.
has long-term objective.



: Training

If $\hat{\pi}(s_0) = a_1$ w.prob $\frac{6}{1+r}$
 a_2 w.prob $\frac{r}{1+r}$

$$\hat{\pi}(s_1) = a_2, \hat{\pi}(s_2) = a_2$$

If mistake is made at s_0 , we stay in s_2 , oo.

: Test,

distribution of
states, visited by π^*

$$P^* = P(s | \pi^*)$$

$$L(a^*, \pi_\theta(s))$$

Loss function (can vary)

i.e.) Negative Log-likelihood (NLL)

$$L(a^*, \pi_\theta(s)) = -\log \pi(a^* | s)$$

i.e.) Square loss

$$L(a^*, \pi_\theta(s)) = \| \pi_\theta(s) - a^* \|_2^2$$

OR KL-Divergence...

Depends on the form of Policy.

How to overcome distribution shift in BC.

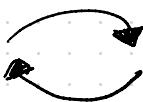
→ Interactive Expert.

→ Query π^* at any s , Build loss as $L(\pi^*(s), \pi_\theta(s))$

↳ provides feedback on state visited by π_θ .

(Naive Approach) : Not guaranteed to converge.

Fix P
Estimate π_θ
Solve $\arg\min_{\theta} \mathbb{E}_{s \sim P} L(\pi^*(s), \pi_\theta(s))$



Fix π_θ
Estimate P , empirically
by rolling out π_θ

(Sequential Learning Reductions) Harder → easier.

1) set initial predictor π_0

2) for $1:m$,

→ Collect T via rolling out π_{m-1} for multiple times

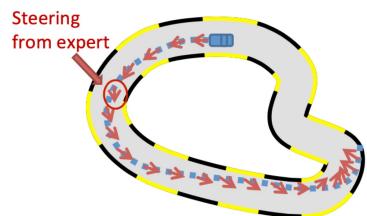
→ Estimate state distribution P_m with SGT .

→ Collect feedback $\{\pi^*(s) | SGT\}$

↳ Interactive expert

→ Data Aggregation (i.e. DAgger)

→ Policy Aggregation (i.e. SEARN & SMILE)



* DAgger.

```

Initialize  $D \leftarrow \emptyset$ .
Initialize  $\hat{\pi}_1$  to any policy in  $\Pi$ .
for  $i = 1$  to  $N$  do
    Let  $\pi_i = \beta_i \pi^* + (1 - \beta_i) \hat{\pi}_i$ . "mixing" policy.
    Sample  $T$ -step trajectories using  $\pi_i$ .
    Get dataset  $D_i = \{(s, \pi^*(s))\}$  of visited states by  $\pi_i$ 
    and actions given by expert.
    Aggregate datasets:  $D \cup D_i$ .
    Train classifier  $\hat{\pi}_{i+1}$  on  $D$ .
end for
Return best  $\hat{\pi}_i$  on validation.
  
```

* SMILE & SEARN.

: at iter n ,

$$\underline{\pi^n} = (1-\alpha) \underline{\pi^{n-1}} + \alpha \underline{\hat{\pi}^n}$$

current
policy

old
policy

new policy
trained by
querying π^*
under P_{n-1} .

$$\underline{\hat{\pi}^n} = \arg\min_{\pi \in \Pi} \mathbb{E}_{s \sim P_{n-1}} [L(a^*, \pi(s))]$$

$$\text{after } N \text{ iter, final policy is } \underline{\pi^N} = \frac{\pi^N - (1-\alpha)^N \pi^*}{1 - (1-\alpha)^N}$$

ETC

Learning by
Watching (Imitation) : Need to define correspondence
doing (demonstration) Avoid correspondence → Tele-operating, kinesthetic Teaching.
Key frame demonstration

i.e. body matching.

Object-based,
end effector-based
Shadowing, Retargeting,