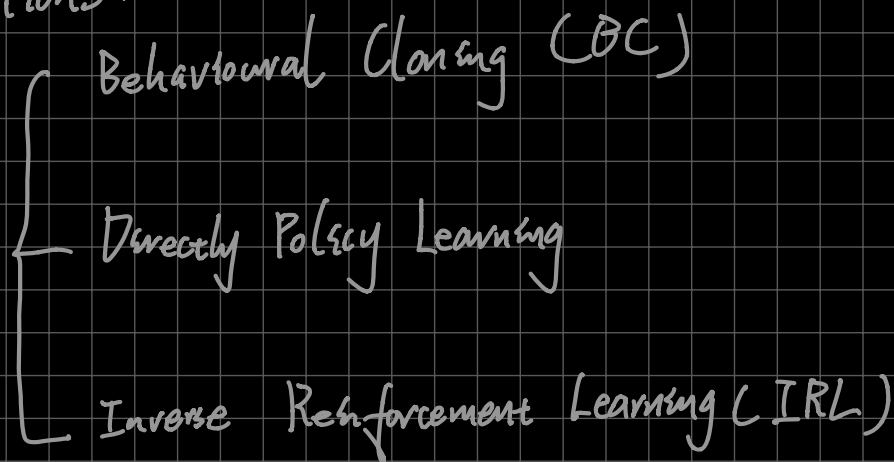


# Imitation Learning (IL)

What:

- Part of Machine Learning, applied in Reinforcement Learning
- Used when manually designing a reward function that satisfies the desired behaviour is extremely complicated.  
*usually complex, graphical/sequential actions*
- Learn the optimal policy directly from expert decisions (which is known as trajectory:  $\tau = (s_0, a_0, s_1, a_1, \dots)$ )

Classifications:



Differences between IL algorithms:

- ① Learning algorithm
- ② Loss function

Advants · Drawbacks:

## ① Behavioural Cloning:

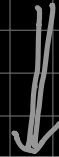
Advantages:

Simple and efficient in short term planning.

Drawbacks:

MDP doesn't satisfy the IID assumption  $\Rightarrow$   
(独立同分布)

Accumulated Reward can lead to unknown states,  
which causes unknown behaviours



expert has never visited, so model is never trained on.

## ② Directly Policy Learning

Advantages:

Efficient in long term planning / complex planning

Drawbacks:

Needs iterative experts.

### ③ Inverse Reinforcement Learning (IRL)

#### Advantages:

- Efficient in long term planning / complex planning
- No need iterative experts.

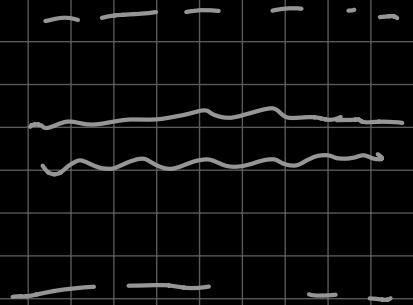
#### Drawbacks:

Difficult to train.

Real World Application by Pieter Abbeel:

[ High variability / long tail  
Always changing  
Need to know what don't know

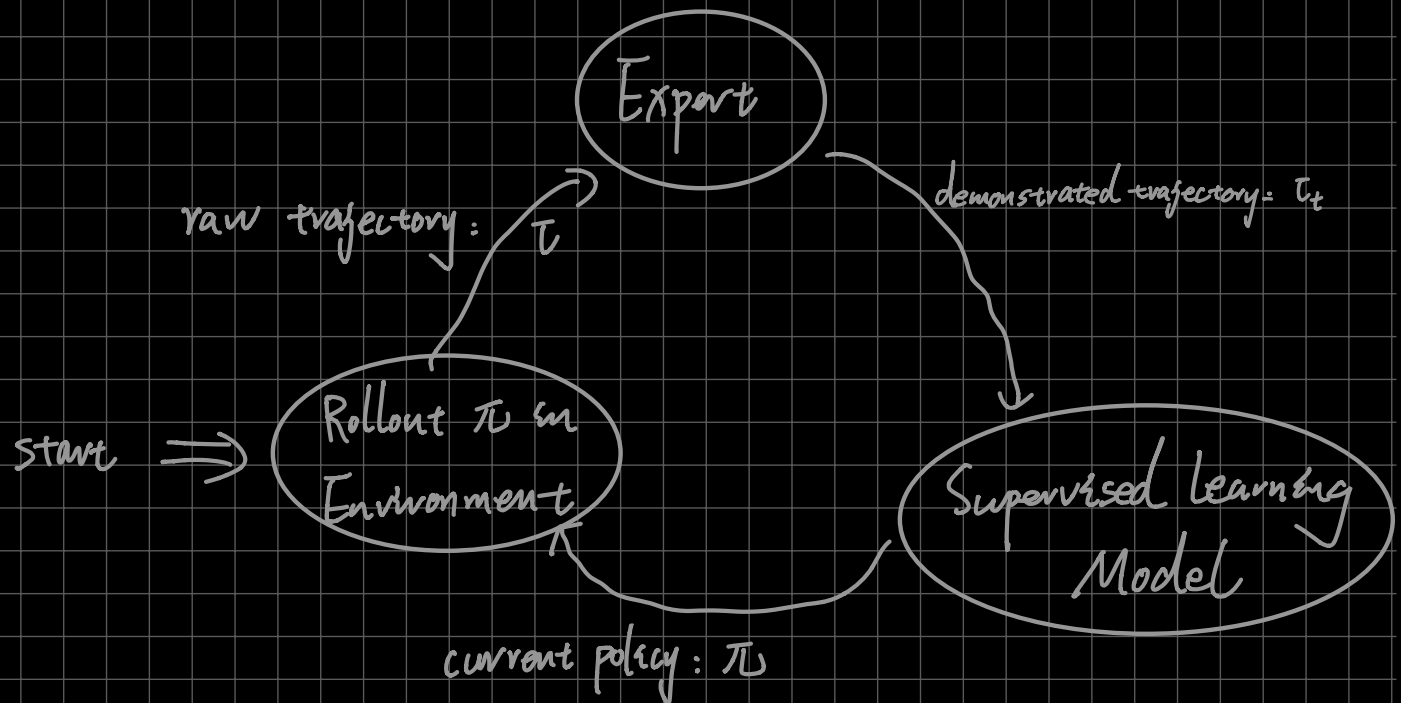
⇒



Research:  $0 \rightarrow 1$  - 70% of benchmark then move on  
Real World Application:  $0.99 \rightarrow 0.999$ , 90% of benchmark is not enough.

# Direct Policy Learning:

Process:



Algorithm:

Initial predictor:  $\pi_0$

For  $m = 1$ :

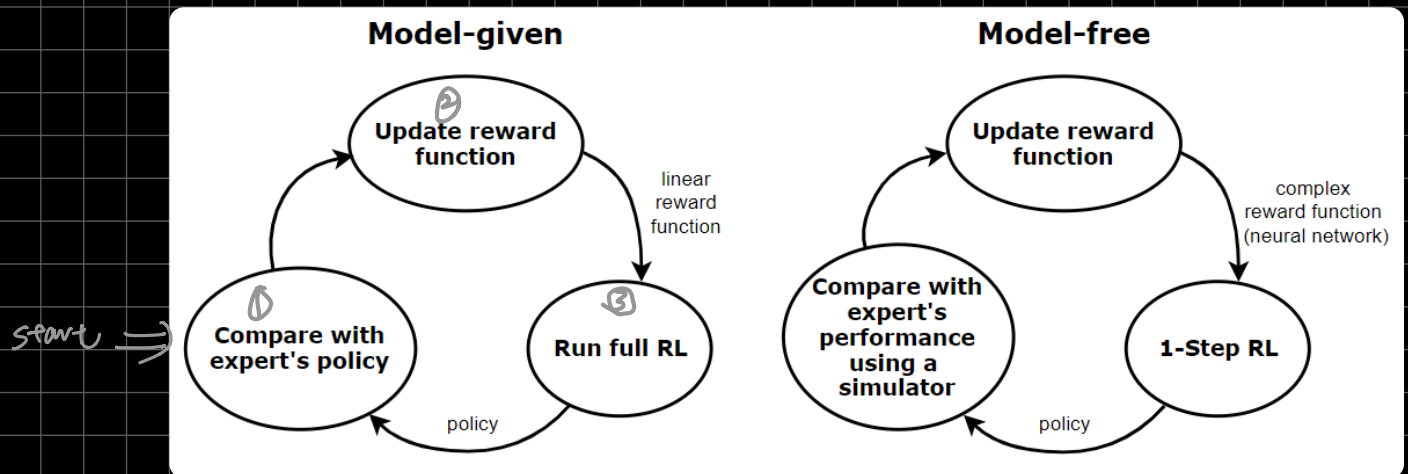
- Collect trajectories  $\tau$  by rolling out  $\pi_{m-1}$
- Estimate state distribution  $P_m$  using  $s \in \tau$
- Collect interactive feedback  $\{\pi^*(s) \mid s \in \tau\}$
- Data Aggregation (e.g. Dagger)
  - Train  $\pi_m$  on  $P_1 \cup \dots \cup P_m$
- Policy Aggregation (e.g. SEARN & SMILE)
  - Train  $\pi'_m$  on  $P_m$
  - $\pi_m = \beta \pi'_m + (1 - \beta) \pi_{m-1}$

# Inversed Reinforcement Learning (IRL)

What:

Estimate the parameterized reward function represented by a neural network from expert's decisions to cause the expert's decisions

Classification:



Algorithm:

(i) Brief Explanation

- Collect expert demonstrations:  $D = \{\tau_1, \tau_2, \dots, \tau_m\}$
- In a loop:
  1. ☐ Compare  $\pi$  with  $\pi^*$  (expert's policy)
  2. ☐ Learn reward function:  $r_\theta(s_t, a_t)$
  3. ☐ Given the reward function  $r_\theta$ , learn  $\pi$  policy using RL
- STOP if  $\pi$  is satisfactory

## (ii) Detailed Algorithm:

Input:  $\underbrace{E}_{env}, \underbrace{X}_{state\ space}, \underbrace{A}_{action\ space}, \underbrace{D = \{T_1, T_2, \dots, T_m\}}_{expert\ demonstration}$

Output:  $\underbrace{R(x) = W^{*T} X}_{\substack{\text{assume linear function,} \\ \text{easy to turn to non-linear} \\ \text{by adding kernel function.}}}, \underbrace{\pi}_{\substack{\text{given reward function, optimal} \\ \text{policy in RL.}}}$

Process:

$$1. \quad \bar{X}^* = \frac{\overbrace{(S_1^1 + S_2^1 + \dots + S_{n_1}^1)}^{T_1's\ states} + (S_1^2 + S_2^2 + \dots + S_{n_2}^2) + \dots + (S_1^m + S_2^m + \dots + S_{n_m}^m)}{n_1 + n_2 + \dots + n_m}$$

Expert Demonstration can be viewed as a sample of optimal policy. Thus, we can use ML to calculate the estimate expected  $\bar{X}$ . Then for optimal reward function  $R(x) = W^{*T} X$ , we have:

$W^{*T} (\underbrace{\bar{X}^*}_{\text{expected } x \text{ from expert}} - \underbrace{\bar{X}^\pi}_{\text{expected } x \text{ from } \pi}) \geq 0$ . Thus, we can update the reward function  $R(x)$  by:

$$\begin{cases} W^* = \arg \max_W \min_{\pi} W^T (\bar{X}^* - \bar{X}^\pi) \\ s.t. \quad \|W\| \leq 1 \end{cases} \quad (2)$$

2.  $\pi \leftarrow$  random policy

3. for  $t = 1, 2, \dots, T$  do

4. Rollout  $\pi$  to get trajectories:  $T_1^\pi, T_2^\pi, \dots$

5.

$$\bar{X}_t^\pi \leftarrow \textcircled{1}$$

6.

$$W^* \leftarrow \textcircled{2}$$

7.

Given reward function  $R(x) = W^* x$ , get the optimal solution  $\pi$ .

8. end for