

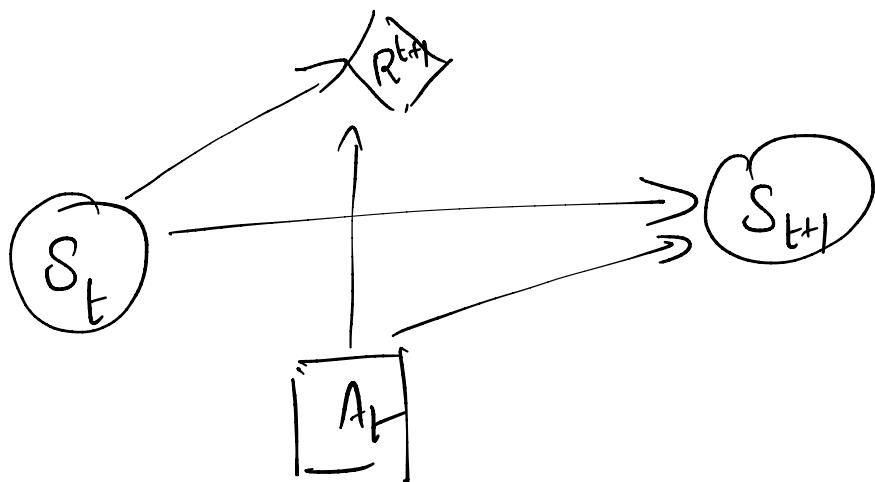
Lecture-12

- Partial observability
 - Offline RL
-

- Partial observability :-

MDP framework :-

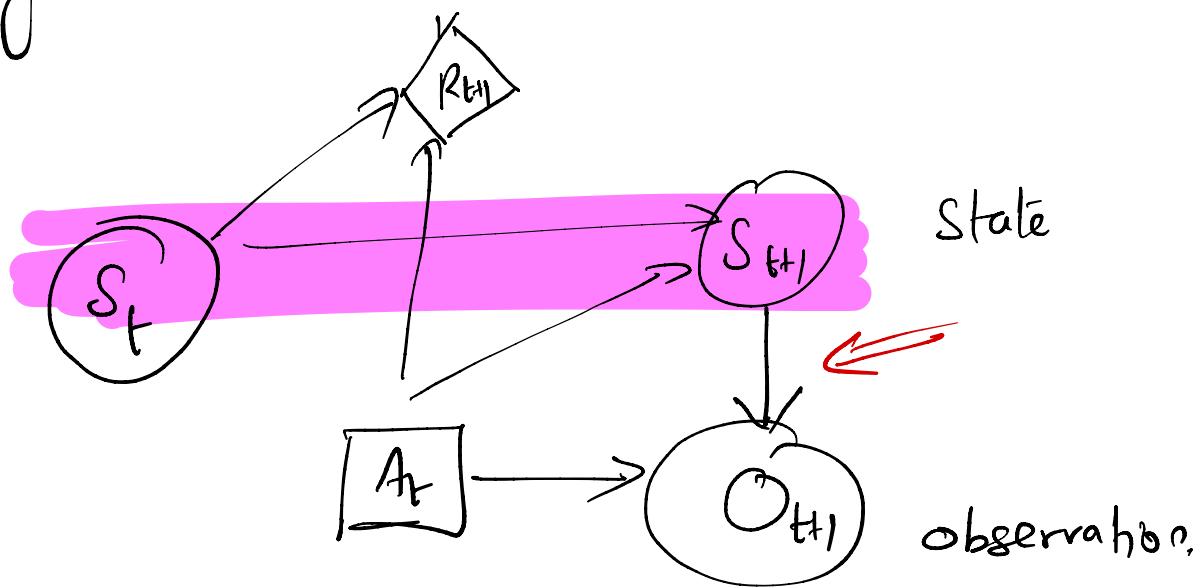
$$\Pr(S_{t+1} | s_0, s_1, \dots, s_t) = \Pr(S_{t+1} | s_t)$$



$$P_{ij}^a = \Pr(s_{t+1} = j \mid s_t = s_i, a_t = a)$$

$$\bar{T}(s, a, s') = \Pr(s_{t+1} = s' \mid s_t = s, a_t = a)$$

Partially Observable MDP:-



\mathcal{O} = be a set of observations -

$$Z : S \times A \mapsto \Delta(\mathcal{O})$$

$$Z(s', a, o') = \Pr(O_{t+1} = o' \mid S_{t+1} = s', A_t = a)$$

$$\underline{\text{PoMDP}}: \langle S, A, T, R, O, Z \rangle$$

POMDP

history: $\langle S_0, O_0, A_0 \rangle \langle S_1, O_1, A_1 \rangle \dots \langle S_t, O_t, A_t \rangle$

\mathcal{H} - set of all complete histories.

System history

\mathcal{HS}

$V: \mathcal{HS} \rightarrow \mathbb{R}$

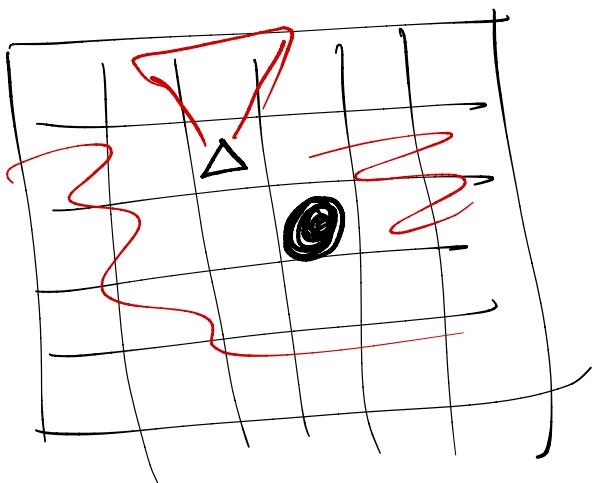
$\langle S_0, A_0 \rangle \langle S_1, A_1 \rangle \dots \langle S_t, A_t \rangle$

Observation history

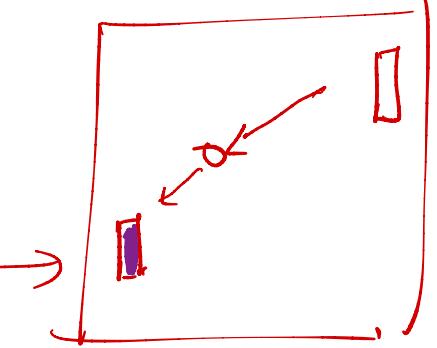
\mathcal{HO}

$\langle A_0, O_1 \rangle \langle A_1, O_2 \rangle \dots \langle A_{t-1}, O_t \rangle$

b_0



$$h_t = \langle a_0, o_1 \rangle \langle a_1, o_2 \rangle \dots \langle a_{t-1}, o_t \rangle$$

$$a_t = \pi(h_t)$$
$$\rightarrow \begin{bmatrix} & o \\ 0 & \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ \text{---} \end{bmatrix}$$


$$z(s', a, o') \quad t-1 \quad t$$
$$z(o' | s', a)$$

POMDP \rightarrow belief state MDP. $\langle s_1, s_2, s_3 \rangle$

Belief state MDP:-

$b_1 (0, 0.5, 0.5)$
 $b_2 (0.9, 0.05, 0.05)$
⋮

$B = \Delta(s)$ — Continuous state space.

A = set of actions is same as the original POMDP.

$R^b : B \times A \rightarrow \mathbb{R}$

$$R^b(b, a) = \sum_{s \in S} b(s) R(s, a)$$

$T^b : B \times A \rightarrow B$

$$T^b(b, a, b') = \Pr(b' | b, a)$$

$$= \sum_{o \in O} \underbrace{\Pr(b' | a, b_o)}_{\text{---}} \underbrace{\Pr(o | a, b)}_{\text{---}}$$

$$= \sum_{o \in O} \Pr(b' | a, b, o) \sum_{s' \in S} \sum_{s \in S} T(s, a, s') b(s) =$$

$$= \sum_{o \in O} \Pr(b' | a, b, o) \sum_{s' \in S} Z(s', a, o) \sum_{s \in S} T(s, a, s') b(s)$$

$$\Pr(b' | a, b, o) = \begin{cases} 1 & \text{if } b_0^a = b' \\ 0 & \text{otherwise.} \end{cases}$$

$$b_0^a(s') = \frac{Z(s', a, o)}{\Pr(o | a, b)}$$

$$Q(b, a) = \sum_{s \in S} b(s) R(s, a) + \gamma \sum_{o \in O} \Pr(o | a, b) V^*(b_0^a)$$

$$\pi^*(b) = \operatorname{argmax}_{a \in A} Q(b, a)$$

Value iteration:

$$V_{t+1}(s) = \max_{a \in A} \left[R(s, a) + \gamma \sum_{s' \in S} T(s, a, s') V_t(s') \right]$$

POMDP value iteration:

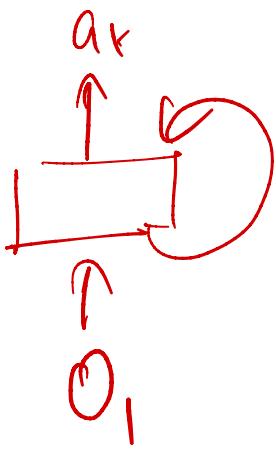
$$V_{t+1}(b) = \max_{a \in A} \left[R^b(b, a) + \gamma \sum_{b' \in B} T^b(b, a, b') V_t(b') \right]$$



$$V_{t+1}(b) = \max_{a \in A} \left[\sum_{s \in S} b(s) R(s, a) + \gamma \sum_{o \in O} P_o(o|a, b) V_t(o) \right]$$

POMDPs using Recurrent Neural Networks:



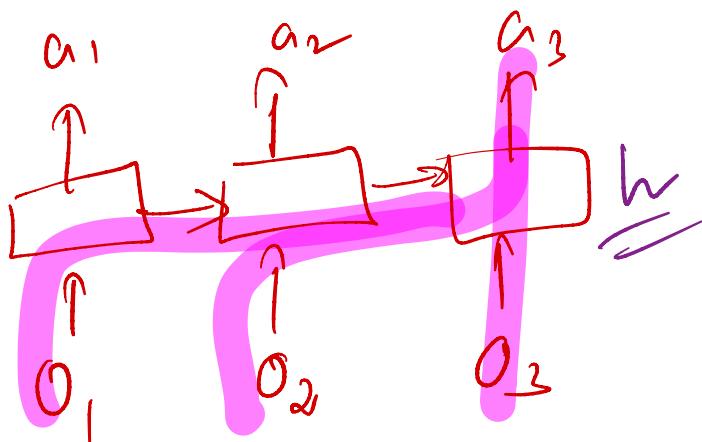


$$h_t = f(o_t) = \sigma(w o_t + b)$$

$$a_t = f(h_t) = v h_t + c$$

$$h_t = f(o_t, h_{t-1}) = \sigma(w q_t + p h_{t-1} + b)$$

$$a_t = f(h_t) = v h_t + c$$

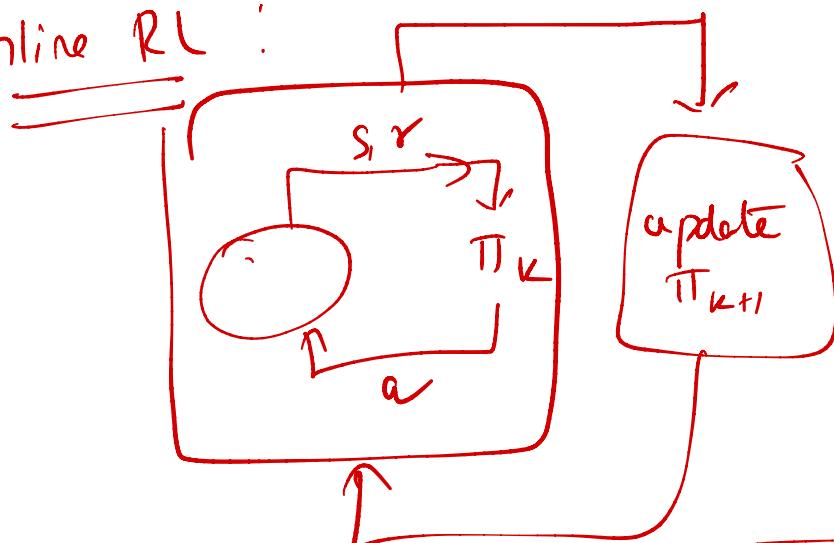


→ offline RL

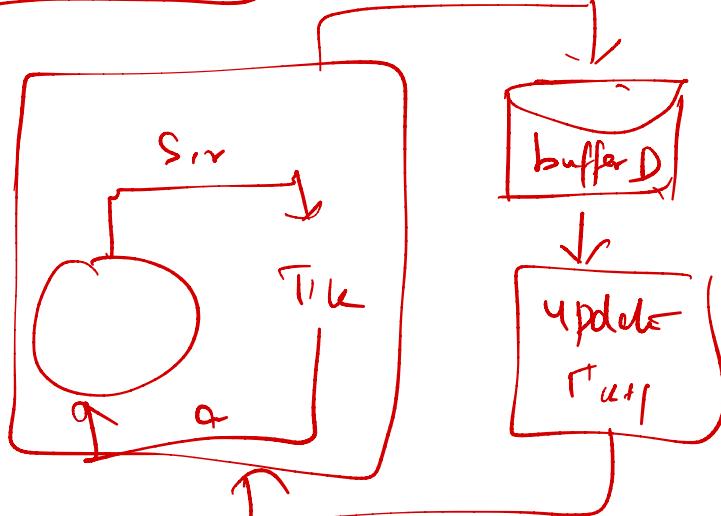
→ learning through interaction.

on-policy
off-Policy

Online RL :

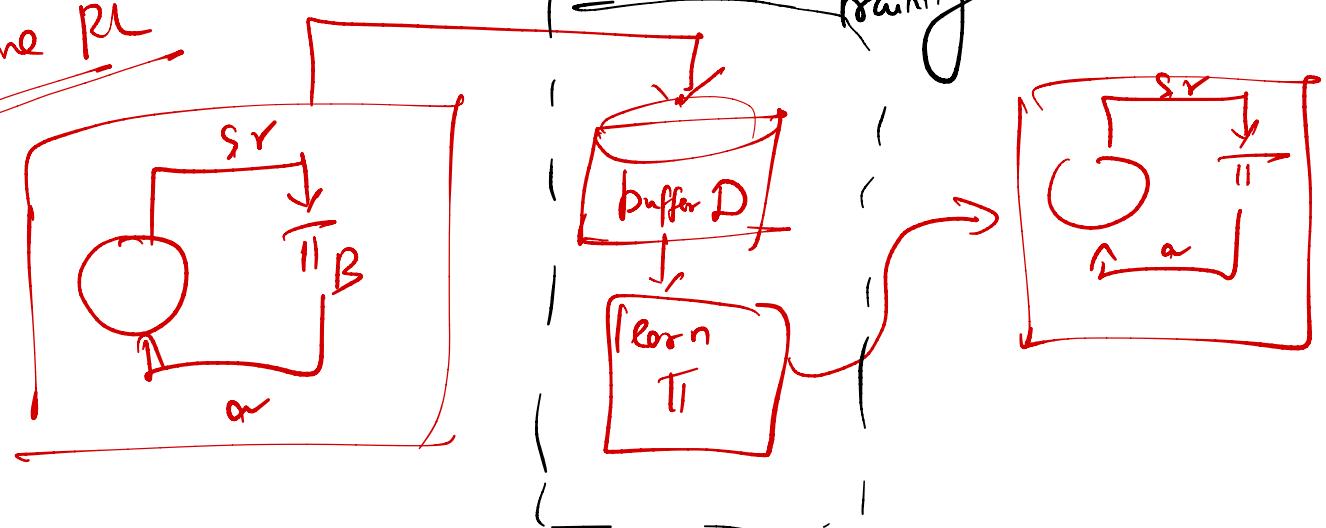


off-Policy RL



offline RL

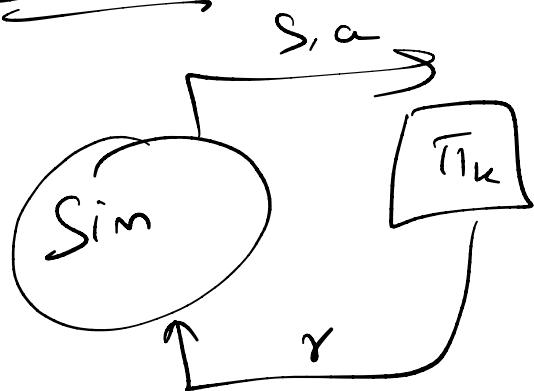
training



why offline RL ?

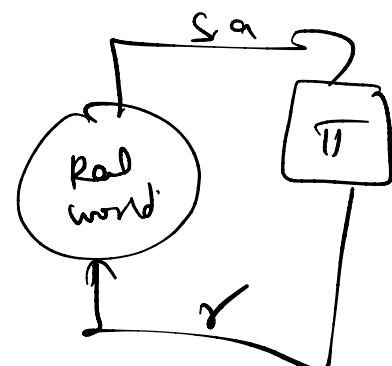
- Safe RL
- offline RL \rightarrow pretraining for online RL
- data efficiency
- multi-tasking.

Sim2Real



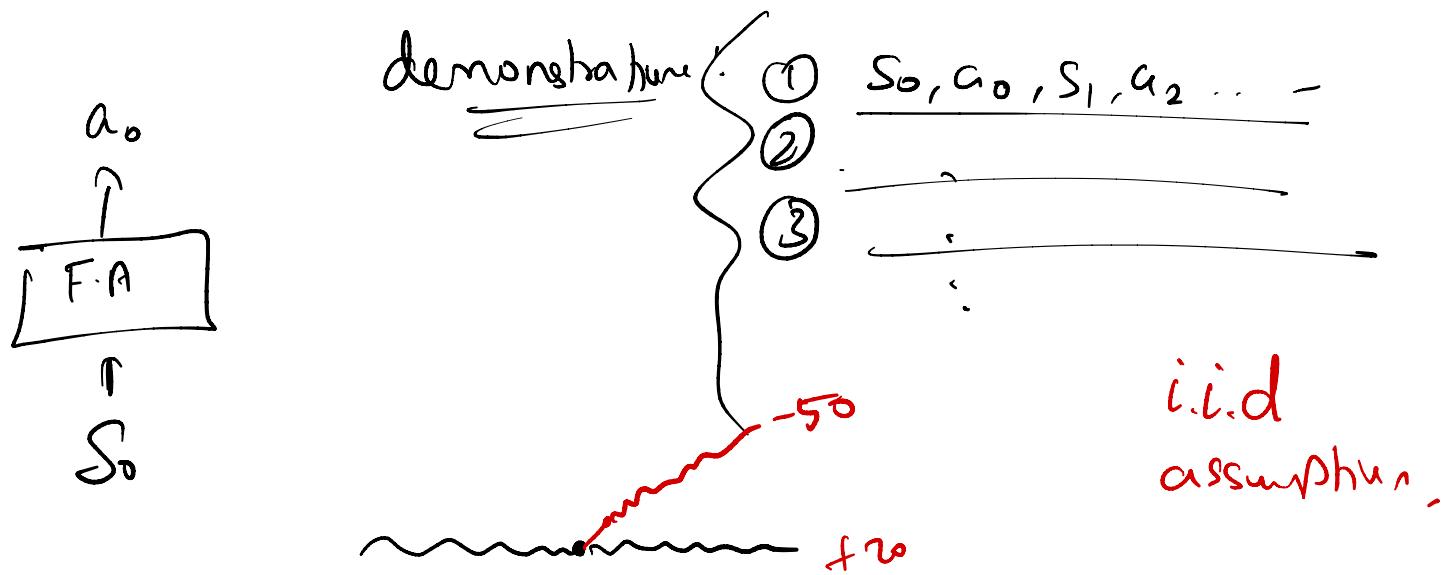
Simulator

CARLA



Behavior cloning:

$$S, A, P(S'|S, a)$$



Generalization

Exploration.

DAGGER: Dataset Aggregation.

$$D \leftarrow \emptyset$$

initialize $\hat{\pi}_i$ to any poly in Π .

for $i = 1 \rightarrow N$:

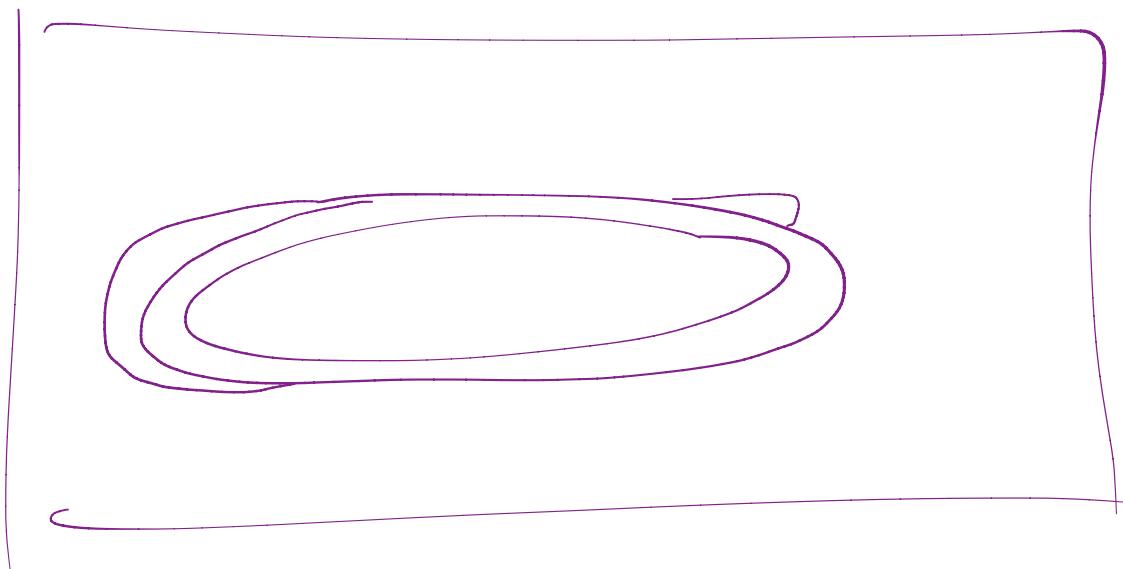
$$\text{let } \pi_i = \beta_i \pi^* + (1 - \beta_i) \hat{\pi}_i$$

sample T-step traj. using π_i

$\xrightarrow{} D_i = \{ (s, \pi^*(s)) \}$ of visited states and actions given by expert.

Aggr. $D \leftarrow D \cup D_i$

Train π_{i+1} on D .



Pure offline RL

Policy constraint methods

$$\pi(a' | s') \sim T\Gamma_B(a' | s')$$

→ Policy constraint

 Policy penalty -

Policy Constraints:

$$\hat{Q}_{k+1}^{\pi} \leftarrow \arg\min_{Q} E_{(s,a,s') \sim D} \left[Q(s,a) - (\gamma r_{s,a}) + \gamma E_{a' \sim \pi_k(a'|s)} [Q_k^{\pi}(s', a')] \right]$$

$$\pi_{k+1} \leftarrow \arg\max_{\pi} E_{\text{SND}} \left[E_{\text{anti}(a|s)} \left[\hat{Q}_{\pi_{k+1}}^{\pi}(s+1) \right] \right]$$

s.t. $D(\pi_i, \pi_{i'}) \leq \epsilon$

~~Policy~~ Policy penalty methods:-

$$\hat{Q}_{k+1}^{\pi} \leftarrow \arg \min_{Q} E_{(s,a,s')} [Q(s,a) - \left(\tilde{r}(s,a) + \gamma \mathbb{E}_{a'} [\hat{Q}_{k+1}^{\pi}(s',a')] \right)^2]$$
$$\tilde{r}(s,a) = r(s,a) - \alpha D(\pi(.|s), \pi_B(.|s))$$

$$\pi_{k+1} \leftarrow \arg \max_{\pi} E_s [\mathbb{E}_a [\hat{Q}_{k+1}^{\pi}(s,a) - \alpha D[\cdot]]]$$