

①

② Perceptual rewards Perceptual Rewards

IRL

$s_t \rightarrow$ visual feature activation at time $t \Rightarrow [s_{1t}, s_{2t}, \dots]$

$\tau \rightarrow \{s_1, \dots, s_T\}$ // sequence of trajectory.

$$P(\tau) = P(s_1, \dots, s_T) = \frac{1}{Z} \exp\left(\sum_{i=1}^T R(s_{it})\right); \text{max Ent model}$$

↑ unknown rewards (target)

{ dynamic programming }

← challenge: How to compute? Boltzmann distribution.

Now, next state, $s_{t+1} = \begin{cases} f(a_t, s_t) \text{ deterministic} \\ \text{or } P(s_{t+1} | a_t, s_t) \text{ probabilistic.} \end{cases}$

Simplifying Assumption,

$$P(\tau) = \prod_{t=1}^T \prod_{i=1}^N P(s_{it}) = \prod_{t=1}^T \prod_{i=1}^N \frac{1}{Z_{it}} \exp(R_i(s_{it}))$$

where $\Rightarrow R_t(s_t) = \sum_{i=1}^N R_i(s_{it})$

Intermediate stage discovery:

$$P(\tau) = \prod_{t=1}^T \prod_{i=1}^N \frac{1}{Z_{it}} \exp(R_{i_g}(s_{it}))$$

↓
index of good step at time t