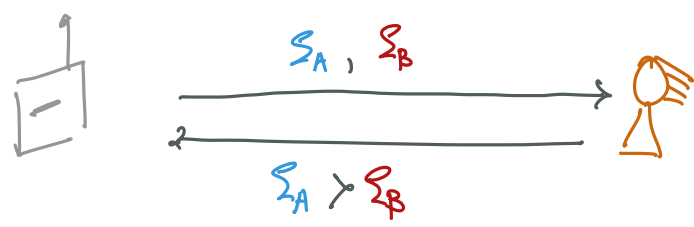


TRAJECTORY $\Sigma = \{s_0, a_0, r_1, a_1, s_1, a_2, \dots\}$



PREFERENCE DATASET $\left\{ \begin{array}{l} \Sigma_A^1, \Sigma_B^1, > \\ \Sigma_A^2, \Sigma_B^2, < \\ \vdots \\ \Sigma_A^N, \Sigma_B^N, \sim \end{array} \right\}$

Q: How CAN I LEARN A REWARD FUNCTION $R_\theta(s, a)$ FROM PREFERENCE?

BRADLEY TERRY MODEL

* PREFERENCES ARE PROBABILISTIC ↑ POSITIVE, REAL-VALUED

$$P(A > B) = \frac{\text{SCORE}(A)}{\text{SCORE}(A) + \text{SCORE}(B)}$$

$$P(\Sigma_A > \Sigma_B) = \frac{\exp(R(\Sigma_A))}{\exp(R(\Sigma_A)) + \exp(R(\Sigma_B))}$$

$$= \frac{1}{1 + \exp(\underbrace{R(\Sigma_B) - R(\Sigma_A)}_{\text{Difference of rewards}})}$$

Difference of rewards.

$$= \frac{1}{1 + \exp(\cdot)}$$

$$P(\Sigma_A > \Sigma_B) = \sigma(R(\Sigma_A) - R(\Sigma_B))$$

GOAL: LEARN $R_\theta(s, a) \rightarrow \sum_{t=0}^{T-1} R_\theta(s_t, a_t)$

$$\begin{aligned} \mathcal{L}(\theta) &= -\log P(\Sigma_A > \Sigma_B) \\ &= -\log \sigma\left(R_\theta(\Sigma_A) - R_\theta(\Sigma_B)\right) \\ &\quad \downarrow \\ &\quad \sum_{t=0}^{T-1} R_\theta(s_t, a_t) \end{aligned}$$

RLHF FROM PREFERENCES

- ① COLLECT PREF DATASET
- ② TRAIN REWARD FUNCTION R_θ
- ③ CALL RL ALGORITHM R_θ