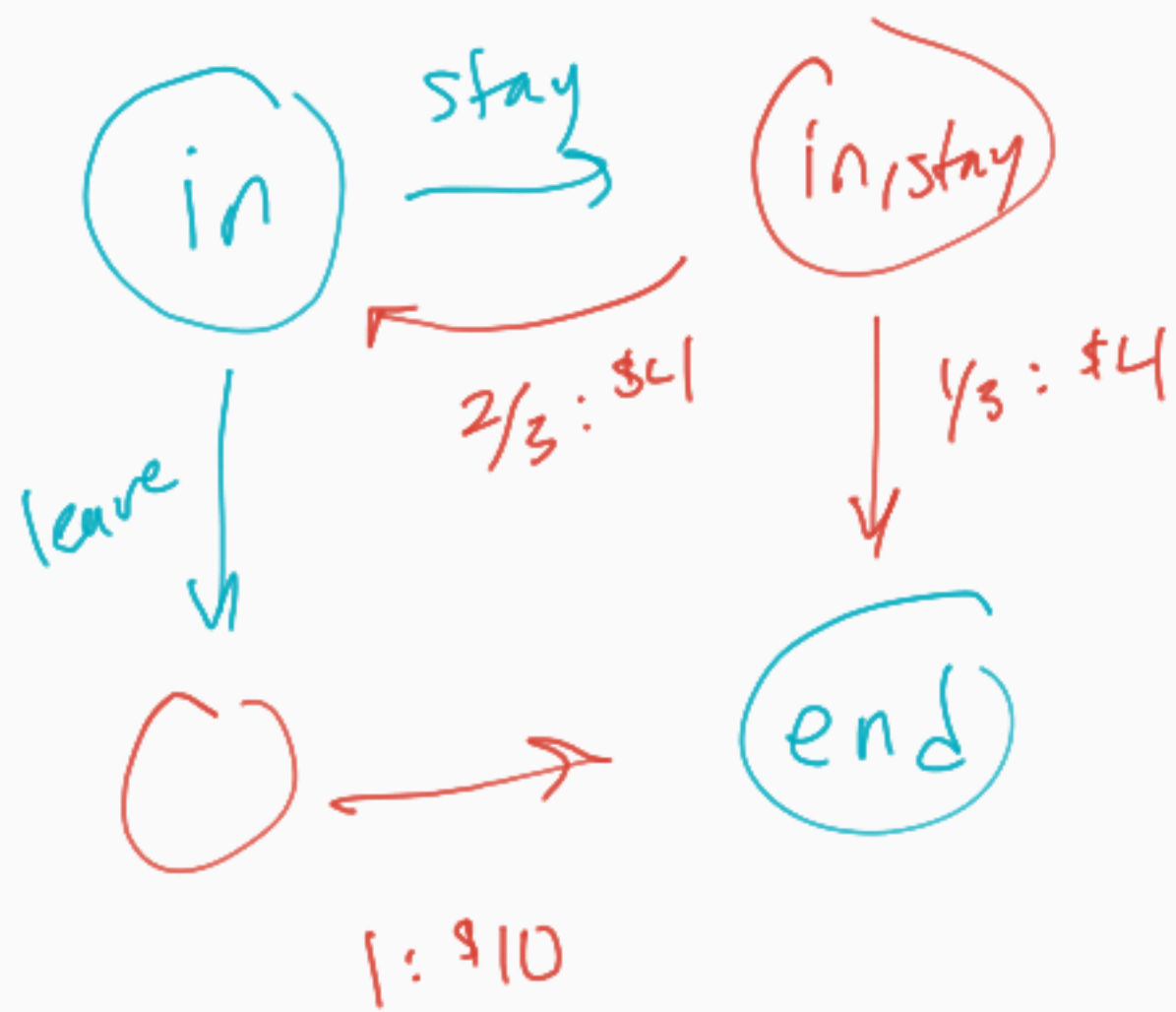


C221 Lecture 8

Modeling
MDP

Inference
policy evaluation
value iteration

Learning
→ RL



$\pi(s) \rightarrow a$

random path $s_1, a_1, r_1, s_2, a_2, r_2, \dots$

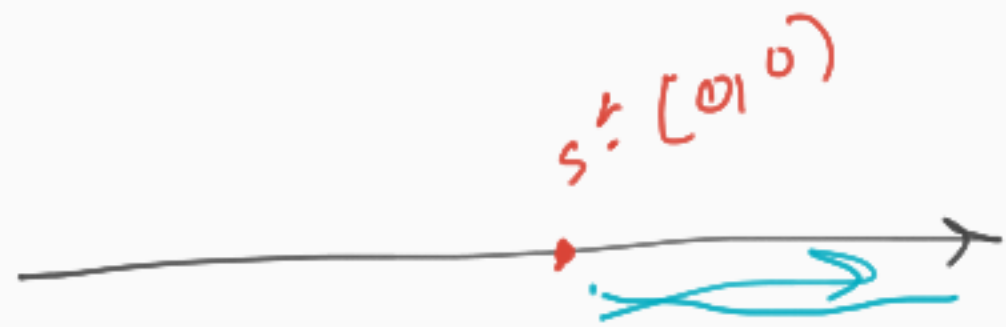
"episode" "roll-out"

utility $u = r_1 + \gamma r_2 + \gamma^2 r_3 + \dots$

"return" / "discounted return"

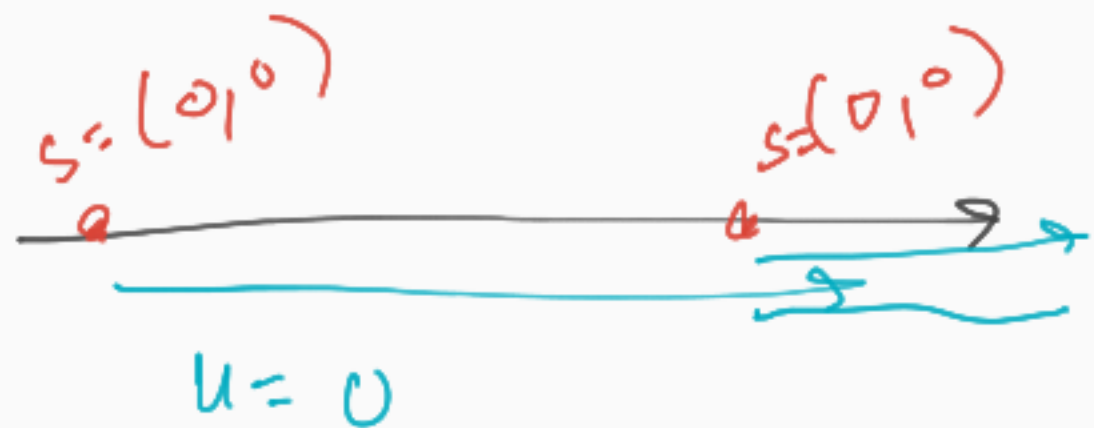


$$\hat{Q}^s((0, 0), a = 0) = \frac{5 + 2 + 0}{3}$$



$$u = 5$$

$$\gamma = 0.9$$



$$u = 2$$

$$u_1 = 4 + 0 + 0$$

$$\hat{Q}_{t+1} = \frac{u_1 + u_2 + u}{3}$$

$$u_2 = 4 + 84$$

$$u_3 = 4 + 84 + 8^2 4 + 8^3 4$$

Model-based value iteration

Model-free Monte Carlo

batch average, online convex combination, stochastic gradient.

SARSA

Q-learning

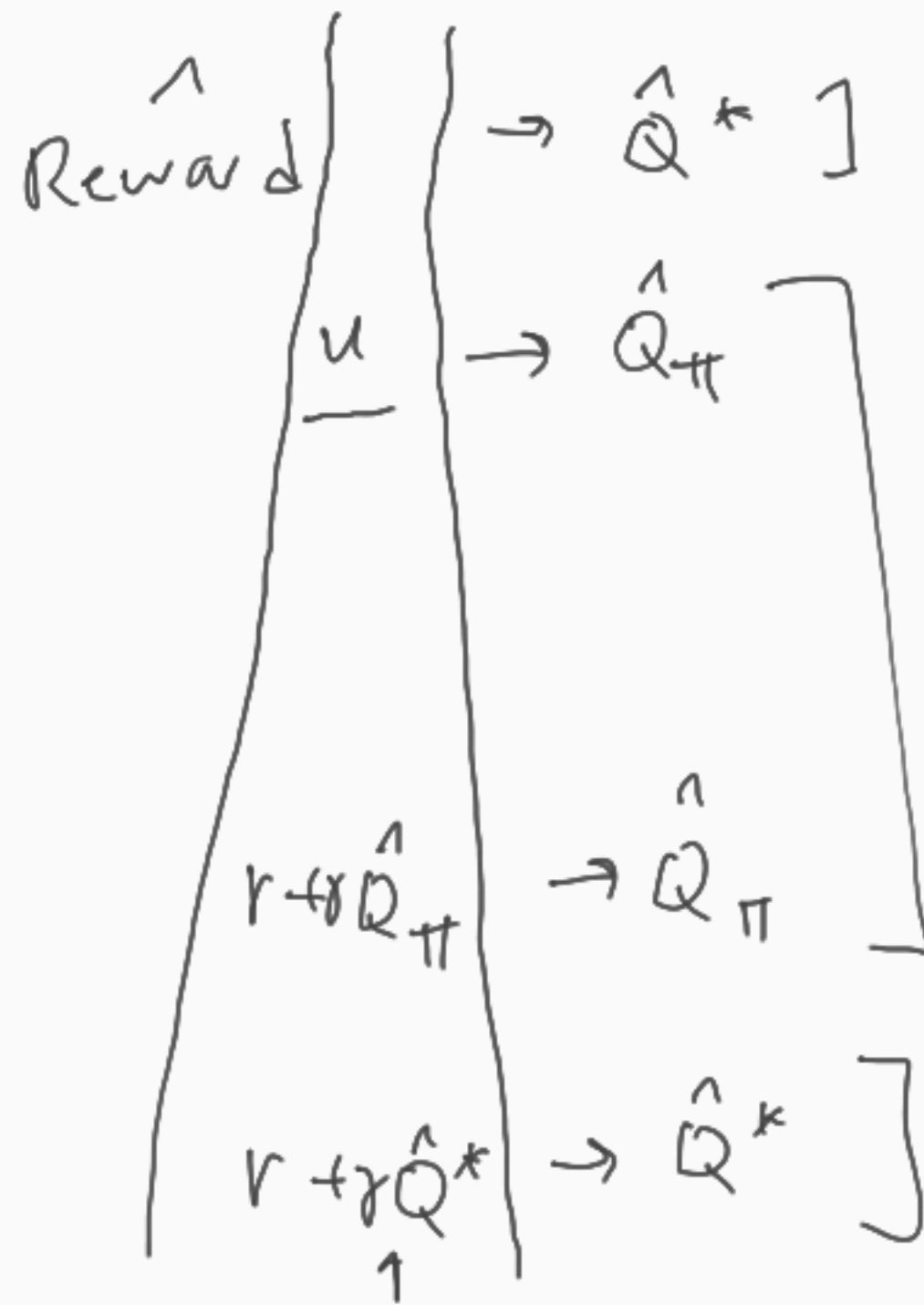
exploration \rightarrow epsilon-greedy

\hat{T}

\hat{Q}_π

\hat{Q}_π

\hat{Q}^*



$\max_{a'}$

$$\underline{u_1} = 4 \quad \underline{u_2} = 1 \quad \underline{u_3} = 1$$

$$\frac{4 + 1 + 1}{3} = \underline{2}$$

$$0 \cdot 0 + 1 \cdot 4 \quad \hat{Q}_\pi = 4$$

↑

$$\frac{1}{2} \cdot 4 + \frac{1}{2} \cdot 1 \quad \hat{Q} = 2 + 1/2$$

$$\frac{2}{3} (5/2) + \frac{1}{3} \cdot 1 = \frac{5}{3} + \frac{1}{3} = \frac{6}{3} = \underline{2}$$