

Last Time
Continuous

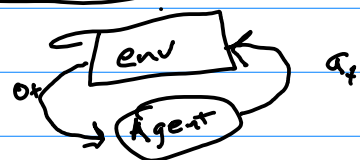
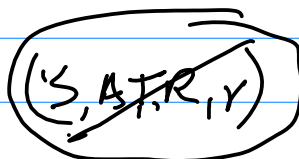
ADP

MPC

Sparse
Tree search
Progressive
widening

Reinforcement Learning

Max-Likelihood Model-Based Tabular RL



This time

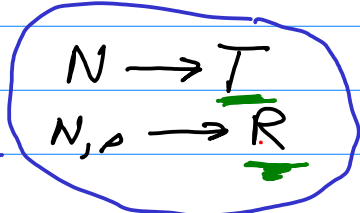
Julia tips

Exploration vs Exploitation: Bandits

Max-Likelihood Model-Based Tabular RL

$N \leftarrow 0$
 $p \leftarrow 0$
 $s \leftarrow s_0$
loop

$N(s, a, s')$
 $p(s, a)$



choose a with π
 $r = \text{act!}(\text{env}, a)$
 $s' = \text{observe}(\text{env})$

Solve with VI

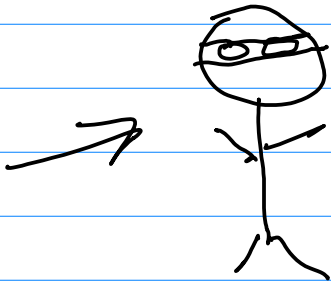
$N(s, a, s') += 1$
 $p(s, a) += r$
estimate \hat{T} from N
estimate \hat{R} from p

$$T(s'|sa) = \frac{N(s, a, s')}{\sum_{s'} N(s, a, s')}$$

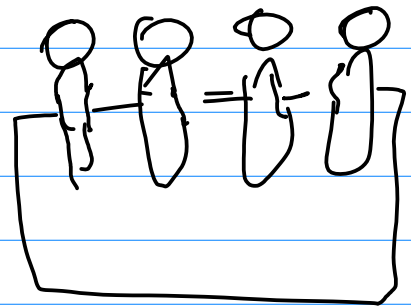
$\pi^* \text{solve}(S, A, \hat{T}, \hat{R}, \gamma)$
 $s \leftarrow s'$

every time
you get to a
terminal state

Bandits



Multi-Armed Bandit



which arm to choose
based on
experience

$A = \{1, \dots, n\}$ Bernoulli Bandit
 θ_a = probability of $r_t = 1$

Breakout Rooms

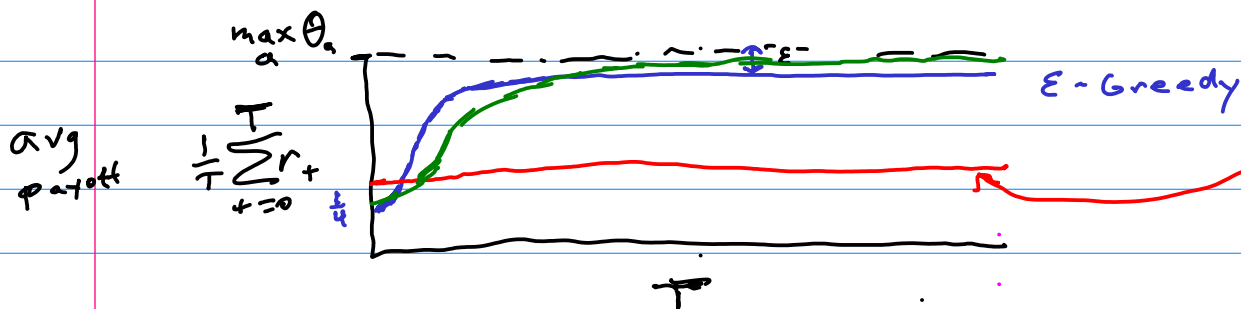
4-arm Bernoulli Bandit: w_i, l_i

How do you choose next action.

B.

Team 2

if $r_t = 1$
 $a_{t+1} = a_t$
 else
 $a_{t+1} = a_t + 1$



Bandit Algos

$$\hat{\mu}_a = \text{estimate of } \theta_a = \frac{w_a}{w_a + l_a}$$

- Greedy $\arg\max_a \hat{\mu}_a$ \emptyset no exploration

Pure
Exploitation

= Explore then commit: random for first k steps
 $\arg\max_a \hat{\mu}_a$ after that

- ϵ -Greedy $\arg\max_a \hat{\mu}_a$ w.p. $(1-\epsilon)$, random otherwise
Decay $\epsilon_{t+1} \leftarrow \alpha \epsilon_t$ where $\alpha \in [0, 1]$

less
regret

- softmax: choose a w.p. $\frac{e^{\lambda \hat{\mu}_a}}{\sum_i e^{\lambda \hat{\mu}_i}}$

- Interval Selection α



- UCB $\arg\max_a \hat{\mu}_a + c \sqrt{\frac{\log N}{N(a)}}$ logarithmic regret

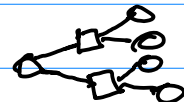
- Thompson Sampling

- keep track of $p(\theta_a)$
- sample $\hat{\mu}_a$
- choose $\arg\max_a \hat{\mu}_a$



- Optimal DP solution

- Usually not tractible
- Gittin's Index



regret $\equiv \Theta^* N - \sum_{t=1}^N r_t$
for Bernoulli

Bayesian

Optimal