

$$S = \{1, 2\} \quad R(s, a) = s^2 \quad V^\pi(1) = 37$$

$$V^*(2)$$

$$V^*(2) = E \sum \gamma^t R(s, a) \quad R(s, a) \leq \frac{4}{1-\gamma}$$

$$\leq \frac{4}{1-\gamma} = 40$$

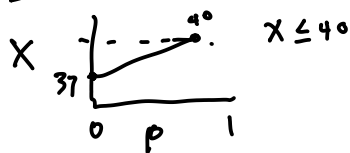
$$V^\pi(1) = 37 = R(1, \pi(1))$$

$$1 + \gamma E[V^\pi(s')]$$

$$E[V^\pi(s')] = 40$$

$$p = T(2|1, \pi(1))$$

$$E[V^\pi(s')] = 40 = pX + (1-p)37$$



Last time

What if  $s$  is continuous (not LQR)

- Approximate DP

- Direct Policy Opt.

Offline

What if problem too big for offline

Online

"Offline"

Before execution:

Find  $V^*$ , Find  $Q^*(s,a) = R(s,a) + \gamma E[V^*(s')]$

During execution:

$\pi^*(s) = \arg\max_a Q^*(s,a)$

"Online"

Before : do nothing

During Execution: Consider action and consequences

\*\* From the current state "

Cross Entropy

(Direct Policy Optimization)

(Offline)

$d$  = initial dist for  $\theta$

loop

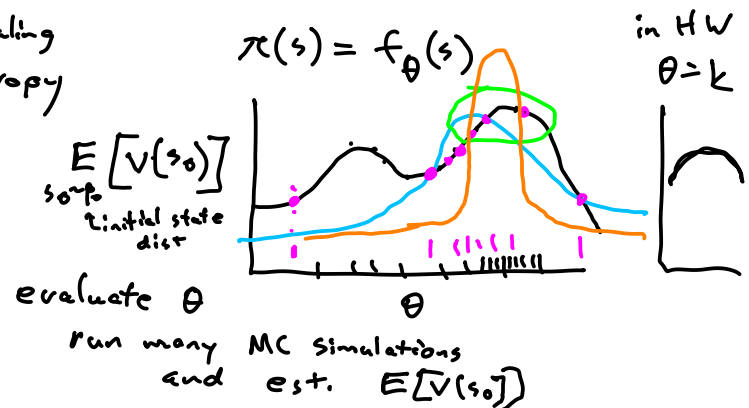
sample  $n$   $\theta$ 's from  $d$

evaluate

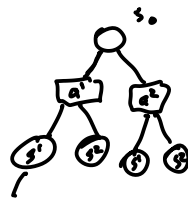
choose elite  $\theta_s$

$d = \text{fit}(\text{elite } \theta_s)$

[ Genetic  
Sim Annealing  
Cross Entropy



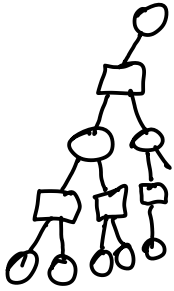
at state  $s$ , want to choose action "planning"

$$Q(s_0, a^2) = R(s_0, a^2) + \sum T(s' | s_0, a^2) V(s')$$


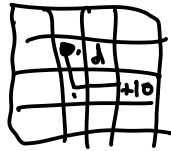
$$Q^\pi(s, a) \equiv R(s, a) + \gamma E[V^\pi(s')]$$

$$V^\pi(s) = \max_a Q^\pi(s, a)$$

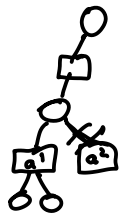
Forward Search



$$O(|A| \times |S|^d)$$



Branch and Bound

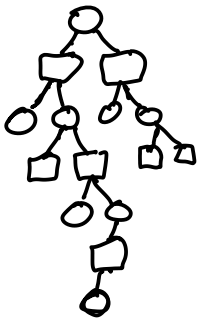


$$\bar{Q}(s, a) < Q(s, a')$$

$$\bar{Q}(s, a) \leq \frac{+10}{1-\gamma}$$

$$\underline{V}(s) \leq \gamma^d + 10$$

Heuristic Search: Not in book but should work well  
 Maintain  $\bar{V}, \underline{V}, \bar{Q}, \underline{Q}$   
 Expand actions with highest  $\bar{Q}$   
 states with biggest  $\bar{V} - \underline{V}$



at root argmax  $\underline{Q}$

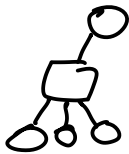
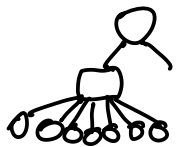
$$|S(s, a)| \ll |S|$$

Sparse Sampling

limit # children to  $n$

only need  $G$

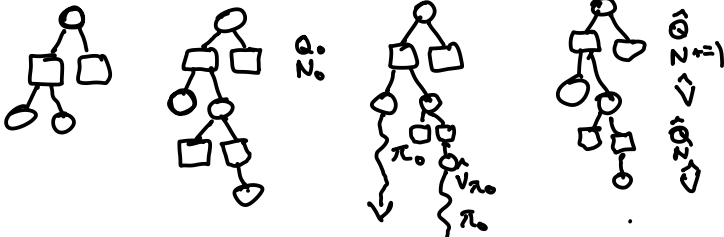
$$O((n|A|)^d)$$



MCTS

# MCTS

Search Expansion Rollout Backup



how to choose actions?

$$\arg \max \hat{Q}$$

UCB

$$\arg \max_a \hat{Q}(s,a) + c \sqrt{\frac{\log N(s)}{N(s,a)}}$$

tuning parameter

exploration bonus

$$MCTS + UCB = UCT$$

$S, A$  continuous

Double Progressive Widening (DPW)

limit |children| to  $k N^\alpha$   
 $\alpha = \frac{1}{2}$

$$\alpha < 1$$

$$\alpha = \frac{1}{2}$$

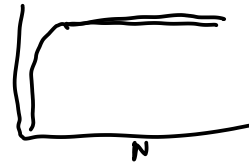


$A$  - only works for 1 dimension

$S$  -  $k = \infty$   $\alpha = \frac{1}{30}$

Sparse Sampling  
with  $n=8$

Sparse UCT



$\pi_\theta = MCTS$   
 $\theta = \{k, \alpha\}$   
 Cross Entropy