

Offline Planning

- Policy learnt before runtime
- Used when state space is small
- Value iteration, Q-learning etc

Online planning

- Policy for current state learnt during runtime
- Used for large state space problems
- MCTS

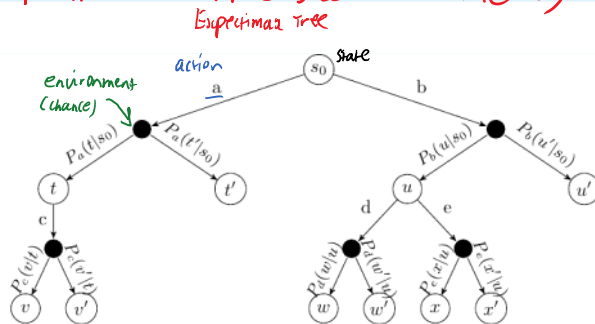
General Online planning algorithm.

For each available action in current state:

- 1) Approximate Q-value through multiple simulations
- 2) Choose action with highest Q-value

Requires simulator or model to approximate/get transition probabilities

Monte Carlo Tree Search (MCTS)



Like a search tree, the Expectimax tree is built incrementally

- Terminate search if complete tree is expanded, or some set time/node limit is reached
- Return action from current state with highest value.

4 steps to MCTS

While limit is not reached, repeat:

Select

Expand

Simulation

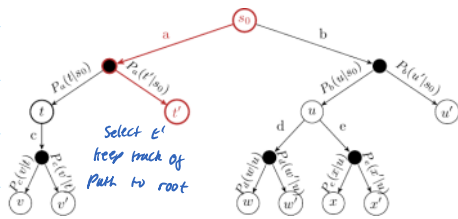
Backpropagate

Return best action

Selection

Select node that has not been fully expanded

(at least one child not yet explored)

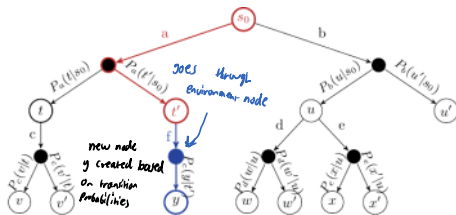


While S is fully expanded:

- 1) Choose action to perform (UCT look at later)
- 2) Use simulator/transition probabilities to determine new state
- 3) $S \leftarrow$ new state

Expansion

Unless terminal/goal state, expand node by choosing an action and creating new child nodes for every possible new state

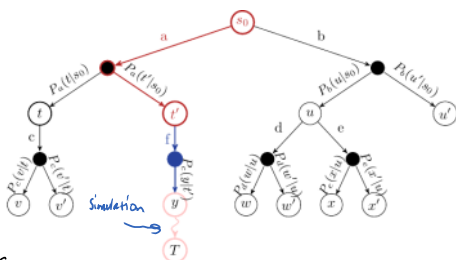


Choose action a to apply

Use simulator/transition probabilities to get new state s' and reward r

Simulation

Perform random actions in simulator until reaching a terminal state



$G = 0$

$t = 0$

While s is not terminal:

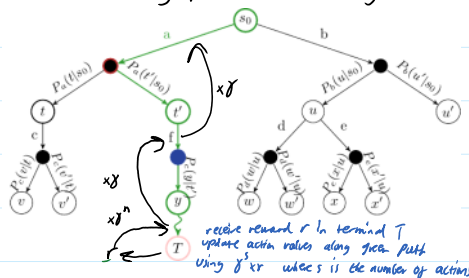
- 1) Apply random action a
- 2) Use simulation/transition probabilities to get new state s' and reward r
- 3) $G \leftarrow G + \gamma^t r$
- 4) $t \leftarrow t + 1$
- 5) $s \leftarrow s'$

After simulation, discard all simulated info except G and the first action taken

Backpropagation

Observe reward r at terminal state. Update value $V(s)$

for each node along path to root using discounted reward



$$\text{Update Rule: } Q(s, a) = Q(s, a) + \frac{1}{N(s, a)} * [r + \gamma G - Q(s, a)]$$

While s exists: (repeat until after we apply to root node)

$N(s, a) \leftarrow N(s, a) + 1$ Increase # times we see action (s, a)

$G \leftarrow r + \gamma G$ Add current reward and discount future reward

$Q(s, a) \leftarrow Q(s, a) + \frac{1}{N(s, a)} [G - Q(s, a)]$ Update Q value

$s \leftarrow \text{Parent state}$ Go up the tree 1 step

$a \leftarrow \text{parent action}$

Upper Confidence Trees (UCT)

Strategy for selecting actions during selection step

Choose action a maximising:

$$Q(s, a) + 2 C_p \sqrt{\frac{2 \ln N(s)}{N(s, a)}}$$

$C_p > 0$ determines weight of exploration

$N(s)$: # times visiting state s

$N(s, a)$: # times performing action (s, a)

Higher value of C_p encourages more exploration