

# Last Time

# Last Time

- Does value iteration always converge?
- Is the value function unique?

# Guiding Questions

# Guiding Questions

- What are the differences between online and offline solutions?
- Are there solution techniques that are *independent* of the state space size?

# Offline vs Online Solutions

Offline

Online

# Offline vs Online Solutions

## Offline

- Before Execution: find  $V^*/Q^*$

## Online

# Offline vs Online Solutions

## Offline

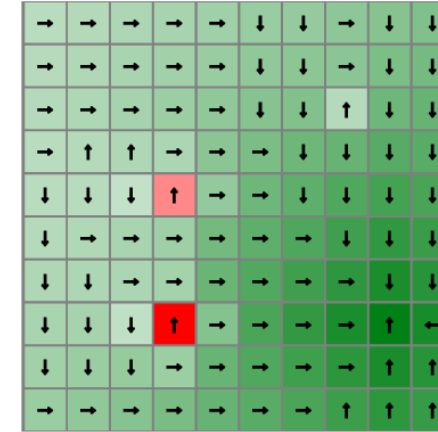
- Before Execution: find  $V^*/Q^*$
- During Execution:  $\pi^*(s) = \operatorname{argmax} Q^*(s, a)$

## Online

# Offline vs Online Solutions

## Offline

- Before Execution: find  $V^*/Q^*$
- During Execution:  $\pi^*(s) = \operatorname{argmax} Q^*(s, a)$



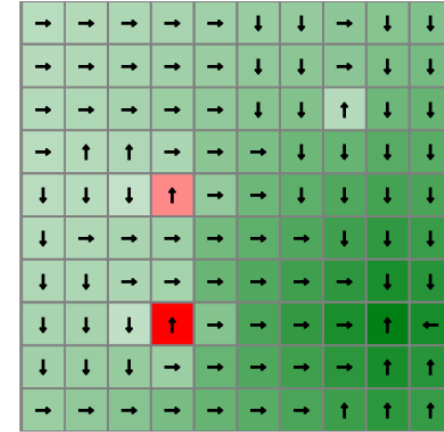
## Online



# Offline vs Online Solutions

## Offline

- Before Execution: find  $V^*/Q^*$
- During Execution:  $\pi^*(s) = \operatorname{argmax} Q^*(s, a)$



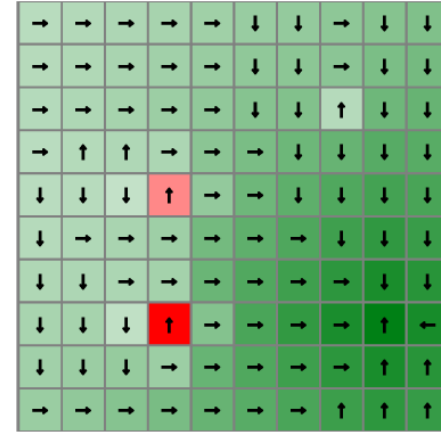
## Online

- Before Execution: <nothing>

# Offline vs Online Solutions

## Offline

- Before Execution: find  $V^*/Q^*$
- During Execution:  $\pi^*(s) = \operatorname{argmax} Q^*(s, a)$



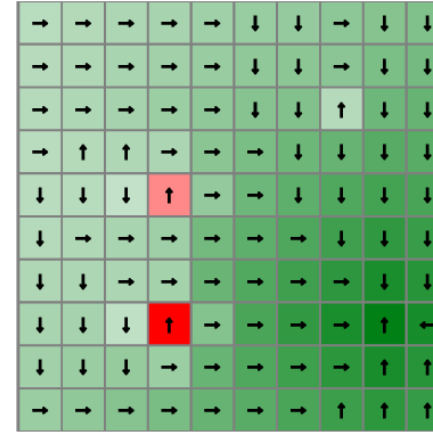
## Online

- Before Execution: <nothing>
- During Execution: Consider actions and their consequences (everything)

# Offline vs Online Solutions

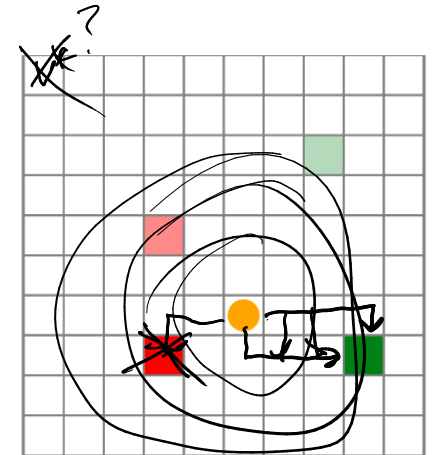
## Offline

- Before Execution: find  $V^*/Q^*$
- During Execution:  $\pi^*(s) = \operatorname{argmax} Q^*(s, a)$



## Online

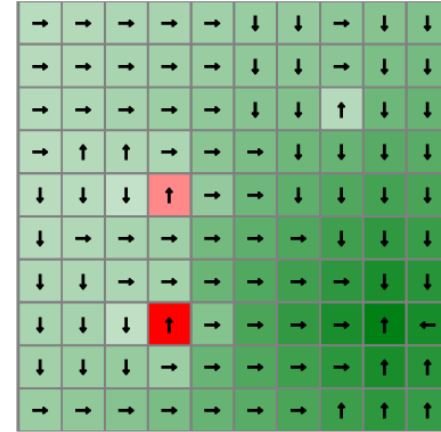
- Before Execution: <nothing>
- During Execution: Consider actions and their consequences (everything)



# Offline vs Online Solutions

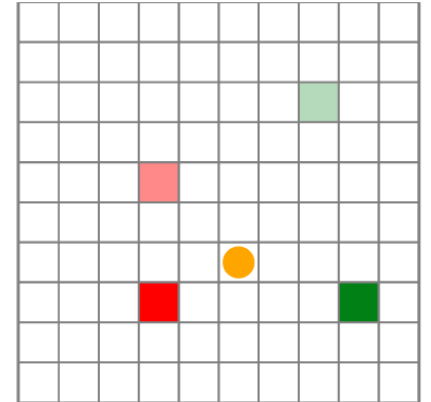
## Offline

- Before Execution: find  $V^*/Q^*$
- During Execution:  $\pi^*(s) = \operatorname{argmax} Q^*(s, a)$



## Online

- Before Execution: <nothing>
- During Execution: Consider actions and their consequences (everything)

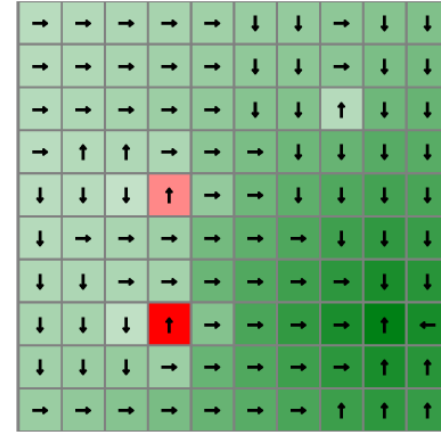


- Why?

# Offline vs Online Solutions

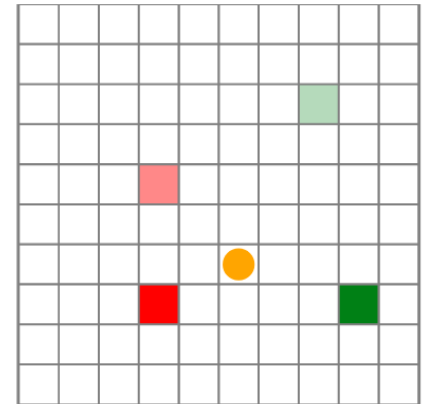
## Offline

- Before Execution: find  $V^*/Q^*$
- During Execution:  $\pi^*(s) = \operatorname{argmax} Q^*(s, a)$



## Online

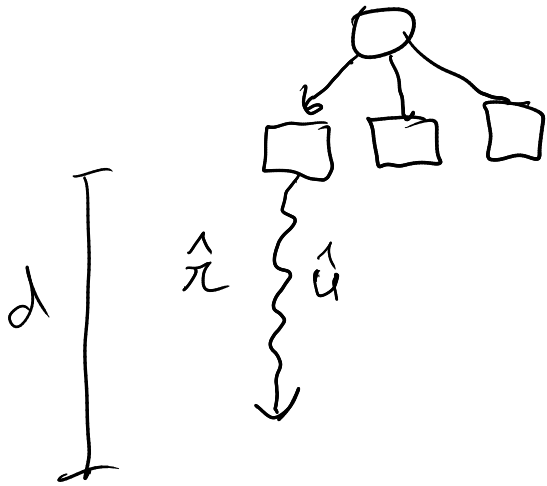
- Before Execution: <nothing>
- During Execution: Consider actions and their consequences (everything)



- Why?
- Online methods are insensitive to the size of  $S$  !

# Lookahead

# Lookahead



```
randstep( $\mathcal{P}::\text{MDP}$ ,  $s$ ,  $a$ ) =  $\mathcal{P}.\text{TR}(s, a)$ 

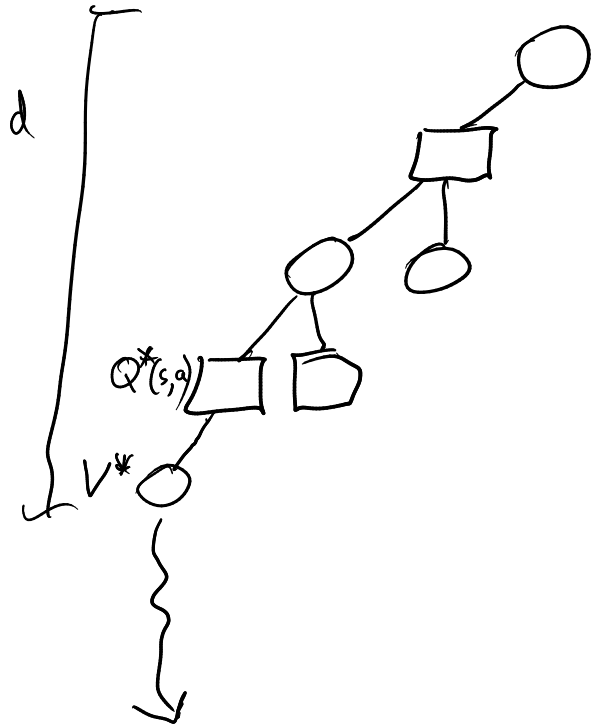
function rollout( $\mathcal{P}$ ,  $s$ ,  $\pi$ ,  $d$ )
    ret = 0.0
    for t in 1:d
         $a = \pi(s)$ 
         $s, r$  = randstep( $\mathcal{P}$ ,  $s$ ,  $a$ )
        ret +=  $\mathcal{P}.\gamma^{(t-1)} * r$ 
    end
    return ret
end

function ( $\pi::\text{RolloutLookahead}$ )( $s$ )
     $U(s)$  = rollout( $\pi.\mathcal{P}$ ,  $s$ ,  $\pi.\pi$ ,  $\pi.d$ )
    return greedy( $\pi.\mathcal{P}$ ,  $U$ ,  $s$ ).a
end
```

# Forward Search



# Forward Search



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

```
function forward_search( $\mathcal{P}$ , s, d, U)
    if  $d \leq 0$ 
        return (a=nothing, u=U(s))
    end
    best = (a=nothing, u=-Inf)
     $U'(s) = \text{forward\_search}(\mathcal{P}, s, d-1, U).u$ 
    for a in  $\mathcal{P}.A$ 
        u = lookahead( $\mathcal{P}$ ,  $U'$ , s, a)
        if u > best.u
            best = (a=a, u=u)
        end
    end
    return best
end
```

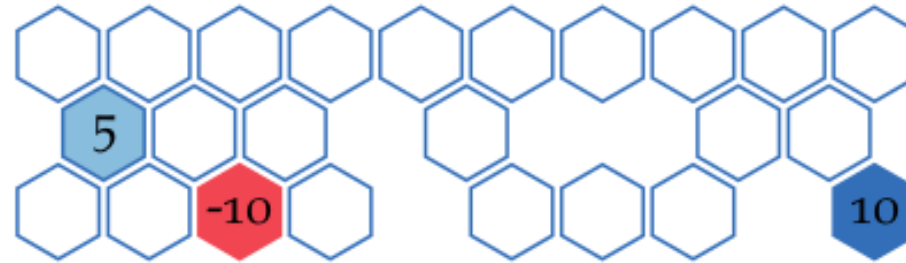
# Forward Search

```
function forward_search( $\mathcal{P}$ , s, d, U)
  if d ≤ 0
    return (a=nothing, u=U(s))
  end
  best = (a=nothing, u=-Inf)
  U'(s) = forward_search( $\mathcal{P}$ , s, d-1, U).u
  for a in  $\mathcal{P}.A$ 
    u = lookahead( $\mathcal{P}$ , U', s, a)
    if u > best.u
      best = (a=a, u=u)
    end
  end
  return best
end
```

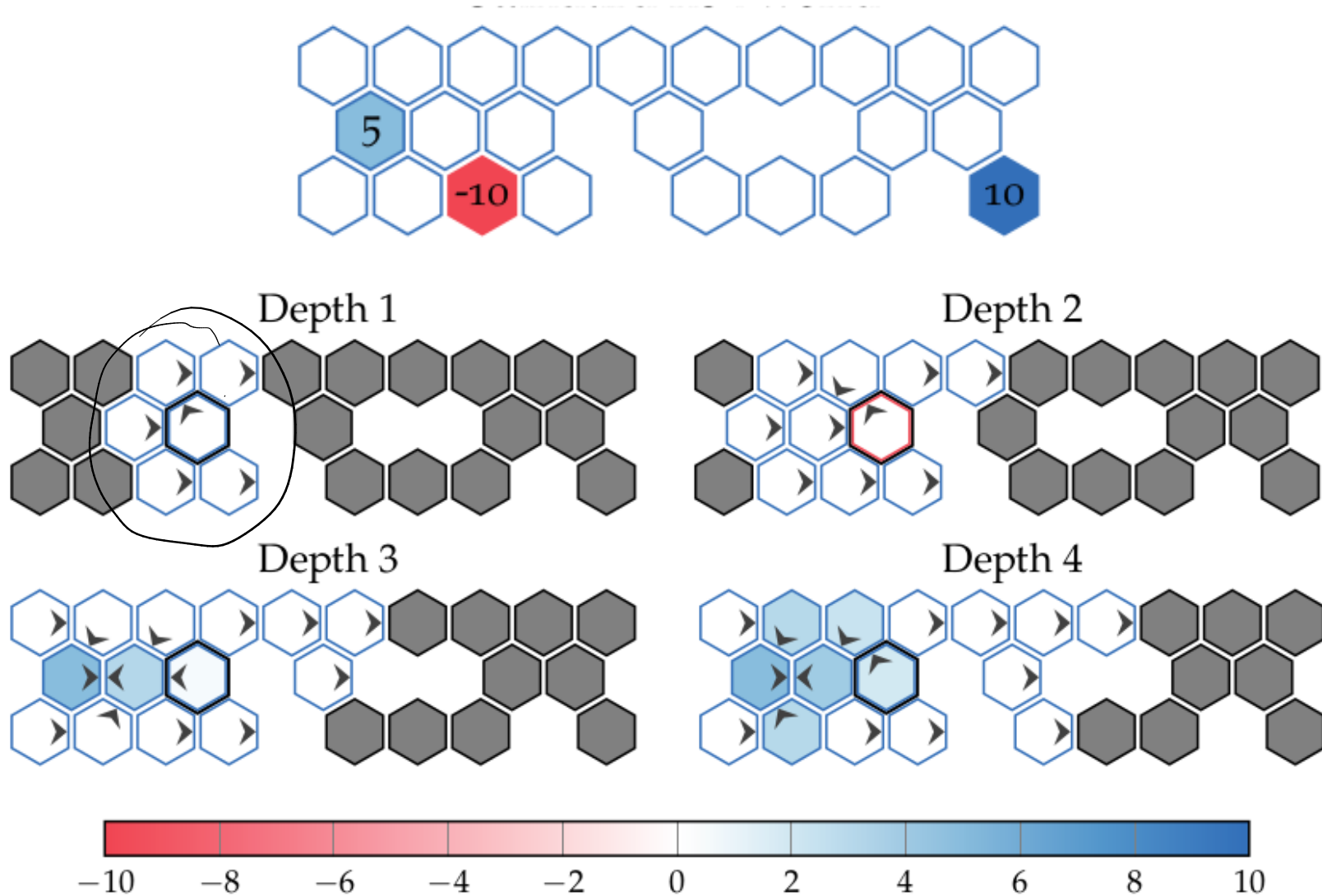
$$O((|S| \times |A|)^d)$$

# Forward Search depth

# Forward Search depth



# Forward Search depth



# Branch and Bound

# Branch and Bound

Assume you have  $\underline{V}(s)$  and  $\bar{Q}(s, a)$

# Branch and Bound

Assume you have  $\underline{V}(s)$  and  $\bar{Q}(s, a)$

```
function branch_and_bound( $\mathcal{P}$ , s, d, Ulo, Qhi)
    if d ≤ 0
        return (a=nothing, u=Ulo(s))
    end
    U'(s) = branch_and_bound( $\mathcal{P}$ , s, d-1, Ulo, Qhi).u
    best = (a=nothing, u=-Inf)
    for a in sort( $\mathcal{P}.\mathcal{A}$ , by=a→Qhi(s,a), rev=true)
        if Qhi(s, a) < best.u
            return best # safe to prune
        end
        u = lookahead( $\mathcal{P}$ , U', s, a)
        if u > best.u
            best = (a=a, u=u)
        end
    end
    return best
end
```



# Branch and Bound

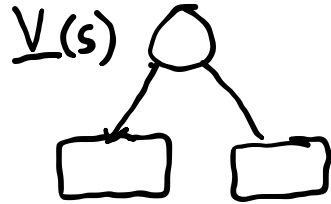
Assume you have  $\underline{V}(s)$  and  $\bar{Q}(s, a)$

$\underline{V}(s)$  ○

```
function branch_and_bound( $\mathcal{P}$ , s, d, Ulo, Qhi)
    if d ≤ 0
        return (a=nothing, u=Ulo(s))
    end
    U'(s) = branch_and_bound( $\mathcal{P}$ , s, d-1, Ulo, Qhi).u
    best = (a=nothing, u=-Inf)
    for a in sort( $\mathcal{P}.A$ , by=a→Qhi(s,a), rev=true)
        if Qhi(s, a) < best.u
            return best # safe to prune
        end
        u = lookahead( $\mathcal{P}$ , U', s, a)
        if u > best.u
            best = (a=a, u=u)
        end
    end
    return best
end
```

# Branch and Bound

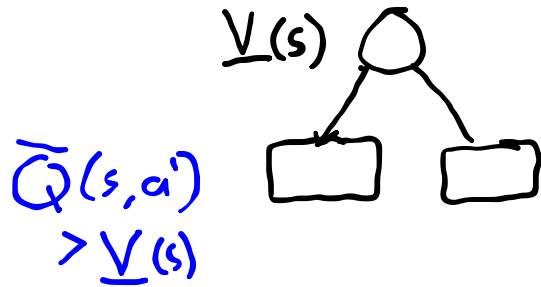
Assume you have  $\underline{V}(s)$  and  $\bar{Q}(s, a)$



```
function branch_and_bound( $\mathcal{P}$ , s, d, Ulo, Qhi)
  if d ≤ 0
    return (a=nothing, u=Ulo(s))
  end
  U'(s) = branch_and_bound( $\mathcal{P}$ , s, d-1, Ulo, Qhi).u
  best = (a=nothing, u=-Inf)
  for a in sort( $\mathcal{P}.\mathcal{A}$ , by=a→Qhi(s,a), rev=true)
    if Qhi(s, a) < best.u
      return best # safe to prune
    end
    u = lookahead( $\mathcal{P}$ , U', s, a)
    if u > best.u
      best = (a=a, u=u)
    end
  end
  return best
end
```

# Branch and Bound

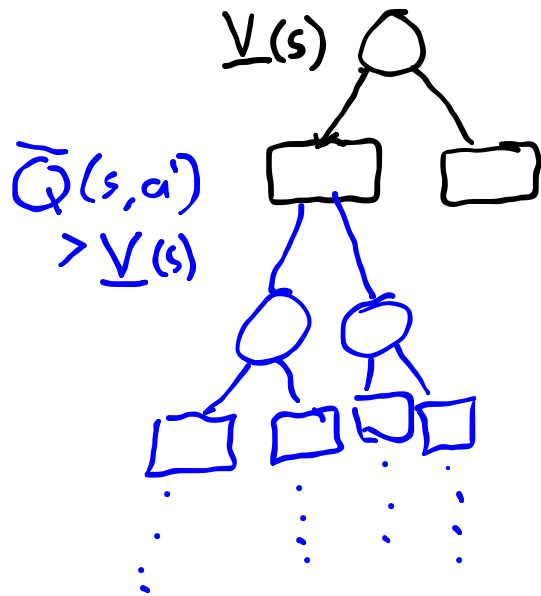
Assume you have  $\underline{V}(s)$  and  $\bar{Q}(s, a)$



```
function branch_and_bound( $\mathcal{P}$ , s, d, Ulo, Qhi)
  if d ≤ 0
    return (a=nothing, u=Ulo(s))
  end
   $U'(s)$  = branch_and_bound( $\mathcal{P}$ , s, d-1, Ulo, Qhi).u
  best = (a=nothing, u=-Inf)
  for a in sort( $\mathcal{P}.A$ , by=a→Qhi(s,a), rev=true)
    if Qhi(s, a) < best.u
      return best # safe to prune
    end
    u = lookahead( $\mathcal{P}$ ,  $U'$ , s, a)
    if u > best.u
      best = (a=a, u=u)
    end
  end
  return best
end
```

# Branch and Bound

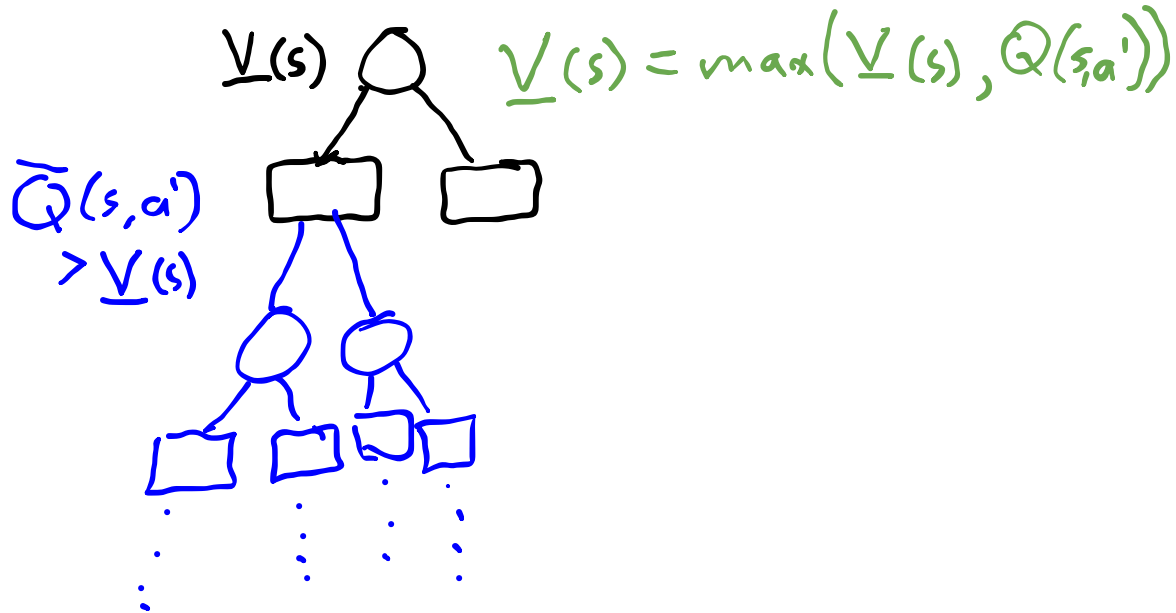
Assume you have  $\underline{V}(s)$  and  $\bar{Q}(s, a)$



```
function branch_and_bound( $\mathcal{P}$ , s, d, Ulo, Qhi)
    if d ≤ 0
        return (a=nothing, u=Ulo(s))
    end
    U'(s) = branch_and_bound( $\mathcal{P}$ , s, d-1, Ulo, Qhi).u
    best = (a=nothing, u=-Inf)
    for a in sort( $\mathcal{P}.A$ , by=a→Qhi(s,a), rev=true)
        if Qhi(s, a) < best.u
            return best # safe to prune
        end
        u = lookahead( $\mathcal{P}$ , U', s, a)
        if u > best.u
            best = (a=a, u=u)
        end
    end
    return best
end
```

# Branch and Bound

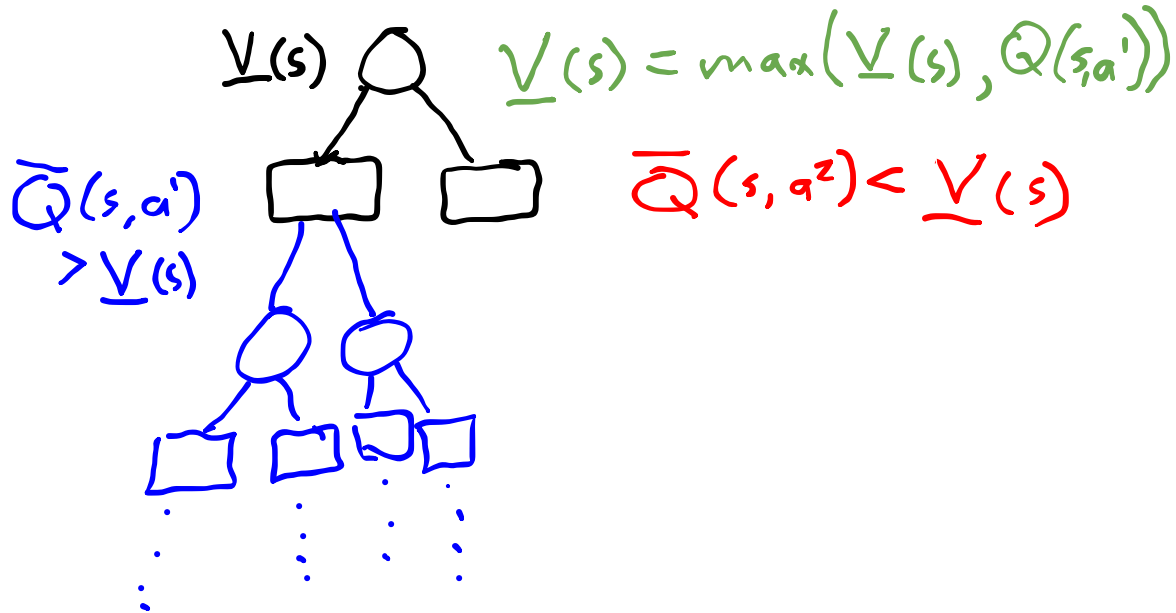
Assume you have  $\underline{V}(s)$  and  $\bar{Q}(s, a)$



```
function branch_and_bound( $\mathcal{P}$ , s, d, Ulo, Qhi)
    if d ≤ 0
        return (a=nothing, u=Ulo(s))
    end
    U'(s) = branch_and_bound( $\mathcal{P}$ , s, d-1, Ulo, Qhi).u
    best = (a=nothing, u=-Inf)
    for a in sort( $\mathcal{P}.A$ , by=a→Qhi(s,a), rev=true)
        if Qhi(s, a) < best.u
            return best # safe to prune
        end
        u = lookahead( $\mathcal{P}$ , U', s, a)
        if u > best.u
            best = (a=a, u=u)
        end
    end
    return best
end
```

# Branch and Bound

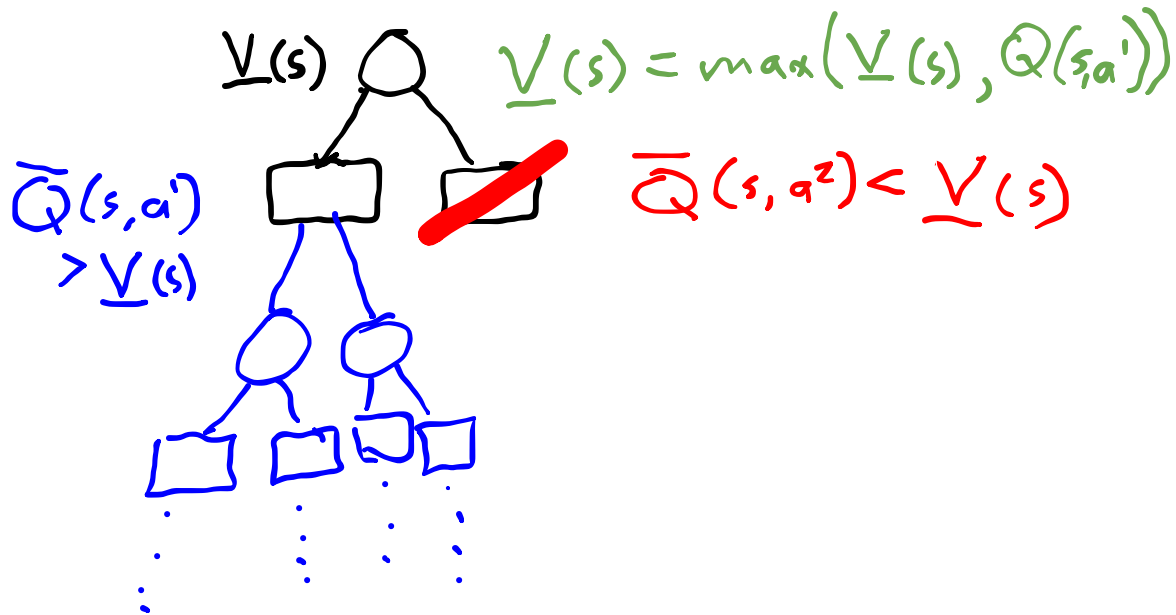
Assume you have  $\underline{V}(s)$  and  $\bar{Q}(s, a)$



```
function branch_and_bound( $\mathcal{P}$ , s, d, Ulo, Qhi)
    if d ≤ 0
        return (a=nothing, u=Ulo(s))
    end
    U'(s) = branch_and_bound( $\mathcal{P}$ , s, d-1, Ulo, Qhi).u
    best = (a=nothing, u=-Inf)
    for a in sort( $\mathcal{P}.\mathcal{A}$ , by=a→Qhi(s,a), rev=true)
        if Qhi(s, a) < best.u
            return best # safe to prune
        end
        u = lookahead( $\mathcal{P}$ , U', s, a)
        if u > best.u
            best = (a=a, u=u)
        end
    end
    return best
end
```

# Branch and Bound

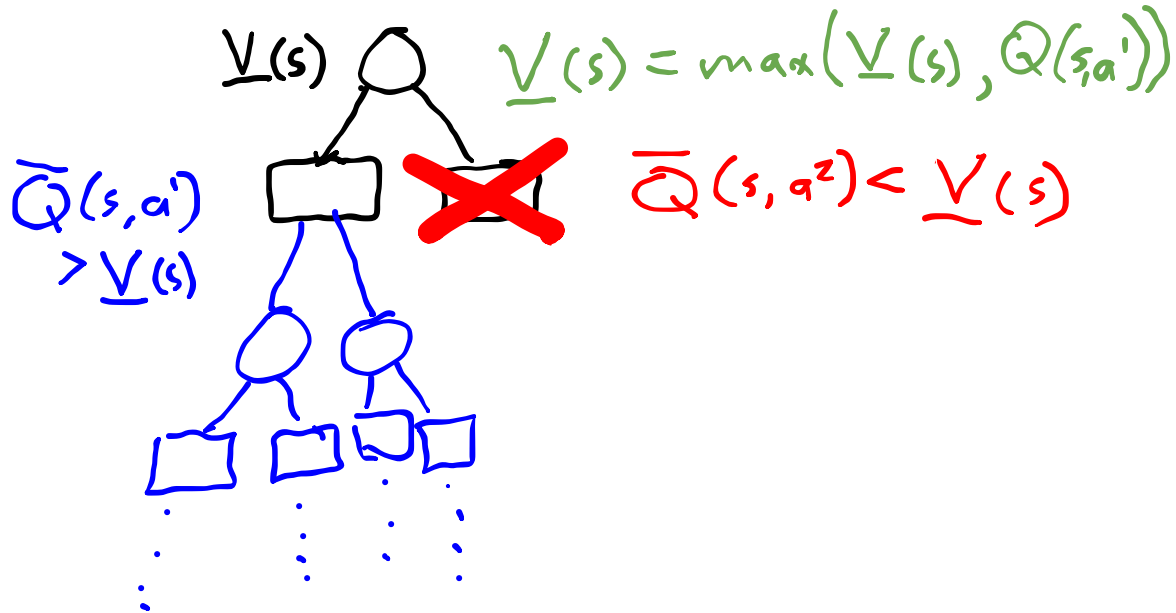
Assume you have  $\underline{V}(s)$  and  $\bar{Q}(s, a)$



```
function branch_and_bound( $\mathcal{P}$ , s, d, Ulo, Qhi)
    if d ≤ 0
        return (a=nothing, u=Ulo(s))
    end
    U'(s) = branch_and_bound( $\mathcal{P}$ , s, d-1, Ulo, Qhi).u
    best = (a=nothing, u=-Inf)
    for a in sort( $\mathcal{P}.A$ , by=a→Qhi(s,a), rev=true)
        if Qhi(s, a) < best.u
            return best # safe to prune
        end
        u = lookahead( $\mathcal{P}$ , U', s, a)
        if u > best.u
            best = (a=a, u=u)
        end
    end
    return best
end
```

# Branch and Bound

Assume you have  $\underline{V}(s)$  and  $\bar{Q}(s, a)$

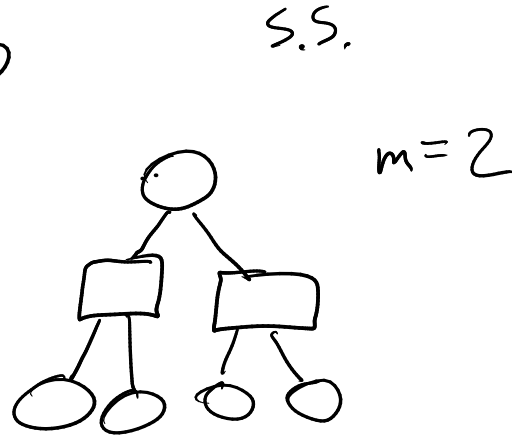
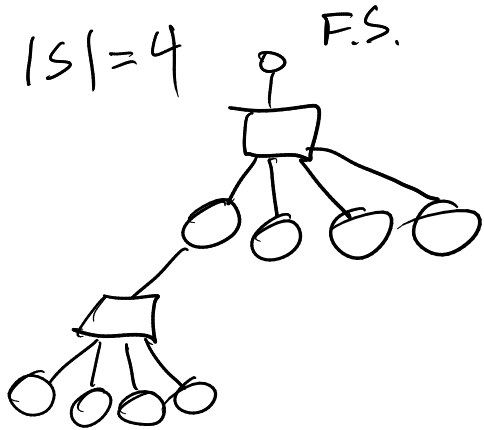


```
function branch_and_bound( $\mathcal{P}$ , s, d, Ulo, Qhi)
    if d ≤ 0
        return (a=nothing, u=Ulo(s))
    end
    U'(s) = branch_and_bound( $\mathcal{P}$ , s, d-1, Ulo, Qhi).u
    best = (a=nothing, u=-Inf)
    for a in sort( $\mathcal{P}.A$ , by=a→Qhi(s,a), rev=true)
        if Qhi(s, a) < best.u
            return best # safe to prune
        end
        u = lookahead( $\mathcal{P}$ , U', s, a)
        if u > best.u
            best = (a=a, u=u)
        end
    end
    return best
end
```



# Sparse Sampling

# Sparse Sampling



```
function sparse_sampling( $\mathcal{P}$ , s, d, m, U)
    if d ≤ 0
        return (a=nothing, u=U(s))
    end
    best = (a=nothing, u=-Inf)
    for a in  $\mathcal{P}.A$ 
        u = 0.0
        for i in 1:m
            s', r = randstep( $\mathcal{P}$ , s, a)
            a', u' = sparse_sampling( $\mathcal{P}$ , s', d-1, m, U)
            u += (r +  $\mathcal{P}.\gamma * u'$ ) / m
        end
        if u > best.u
            best = (a=a, u=u)
        end
    end
    return best
end
```

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

$$|V^{ss}(s) - V^*(s)| \leq \epsilon$$



# Sparse Sampling

```
function sparse_sampling( $\mathcal{P}$ , s, d, m, U)
    if d ≤ 0
        return (a=nothing, u=U(s))
    end
    best = (a=nothing, u=-Inf)
    for a in  $\mathcal{P}.A$ 
        u = 0.0
        for i in 1:m
            s', r = randstep( $\mathcal{P}$ , s, a)
            a', u' = sparse_sampling( $\mathcal{P}$ , s', d-1, m, U)
            u += (r +  $\mathcal{P}.\gamma*u'$ ) / m
        end
        if u > best.u
            best = (a=a, u=u)
        end
    end
    return best
end
```

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

# Sparse Sampling

```
function sparse_sampling( $\mathcal{P}$ , s, d, m, U)
    if d ≤ 0
        return (a=nothing, u=U(s))
    end
    best = (a=nothing, u=-Inf)
    for a in  $\mathcal{P}.A$ 
        u = 0.0
        for i in 1:m
            s', r = randstep( $\mathcal{P}$ , s, a)
            a', u' = sparse_sampling( $\mathcal{P}$ , s', d-1, m, U)
            u += (r +  $\mathcal{P}.\gamma * u'$ ) / m
        end
        if u > best.u
            best = (a=a, u=u)
        end
    end
    return best
end
```

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

$$|V^{\text{SS}}(s) - V^*(s)| \leq \epsilon$$

# Sparse Sampling

```
function sparse_sampling( $\mathcal{P}$ , s, d, m, U)
    if d ≤ 0
        return (a=nothing, u=U(s))
    end
    best = (a=nothing, u=-Inf)
    for a in  $\mathcal{P}.A$ 
        u = 0.0
        for i in 1:m
            s', r = randstep( $\mathcal{P}$ , s, a)
            a', u' = sparse_sampling( $\mathcal{P}$ , s', d-1, m, U)
            u += (r +  $\mathcal{P}.\gamma*u'$ ) / m
        end
        if u > best.u
            best = (a=a, u=u)
        end
    end
    return best
end
```

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

$$|V^{\text{SS}}(s) - V^*(s)| \leq \epsilon$$

$m$ ,  $\epsilon$ , and  $d$  related, but independent of  $|S|$

# Sparse Sampling

```
function sparse_sampling( $\mathcal{P}$ , s, d, m, U)
    if d ≤ 0
        return (a=nothing, u=U(s))
    end
    best = (a=nothing, u=-Inf)
    for a in  $\mathcal{P}.A$ 
        u = 0.0
        for i in 1:m
            s', r = randstep( $\mathcal{P}$ , s, a)
            a', u' = sparse_sampling( $\mathcal{P}$ , s', d-1, m, U)
            u += (r +  $\mathcal{P}.\gamma * u'$ ) / m
        end
        if u > best.u
            best = (a=a, u=u)
        end
    end
    return best
end
```

$O((m|A|)^d)$        $|V^{\text{SS}}(s) - V^*(s)| \leq \epsilon$        $m, \epsilon$ , and  $d$  related, but independent of  $|S|$

# Monte Carlo Tree Search (MCTS/UCT)

# Monte Carlo Tree Search (MCTS/UCT)

Search



# Monte Carlo Tree Search (MCTS/UCT)

Search

Expansion

# Monte Carlo Tree Search (MCTS/UCT)

Search

Expansion

Rollout

# Monte Carlo Tree Search (MCTS/UCT)

Search

Expansion

Rollout

Backup

# Monte Carlo Tree Search (MCTS/UCT)

Search



Expansion

Rollout

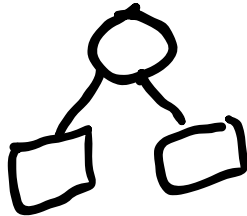
Backup

# Monte Carlo Tree Search (MCTS/UCT)

Search



Expansion

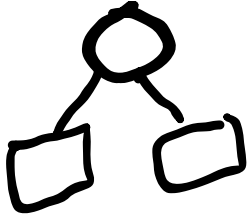


Rollout

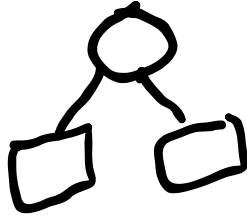
Backup

# Monte Carlo Tree Search (MCTS/UCT)

Search



Expansion

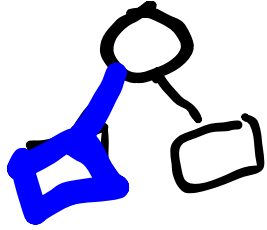


Rollout

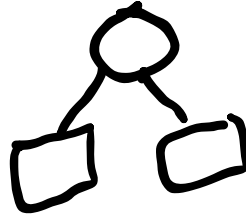
Backup

# Monte Carlo Tree Search (MCTS/UCT)

Search



Expansion



Rollout

Backup

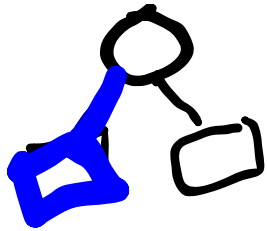
$$Q(s, a) + c \sqrt{\frac{\log N(s)}{N(s, a)}}$$

low  $N(s, a)/N(s)$  = high bonus

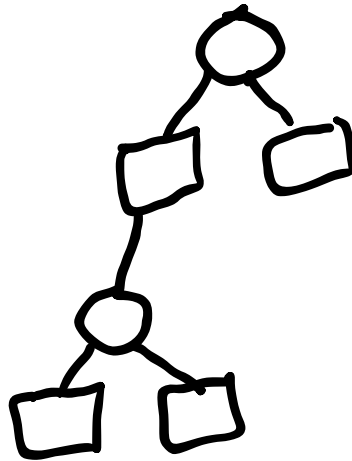
start with  $c = 2(\bar{V} - \underline{V})$

# Monte Carlo Tree Search (MCTS/UCT)

Search



Expansion



Rollout

Backup

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

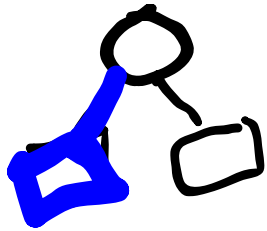
low  $N(s, a)/N(s)$  = high bonus

start with  $c = 2(\bar{V} - \underline{V})$

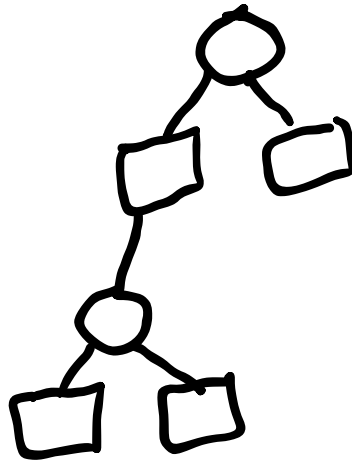


# Monte Carlo Tree Search (MCTS/UCT)

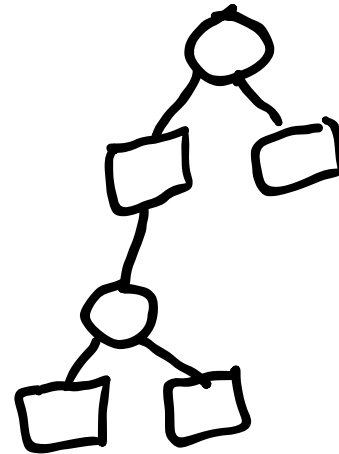
Search



Expansion



Rollout



Backup

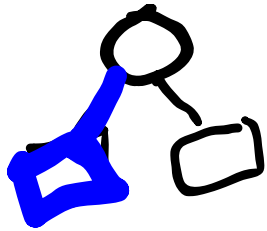
$$Q(s, a) + c \sqrt{\frac{\log N(s)}{N(s, a)}}$$

low  $N(s, a)/N(s)$  = high bonus

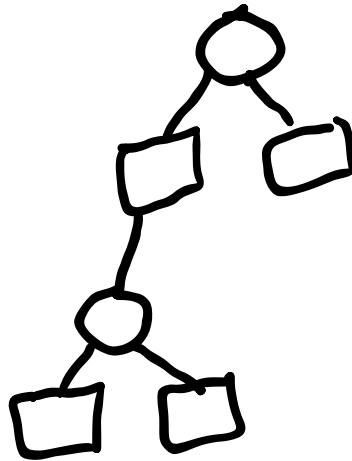
start with  $c = 2(\bar{V} - \underline{V})$

# Monte Carlo Tree Search (MCTS/UCT)

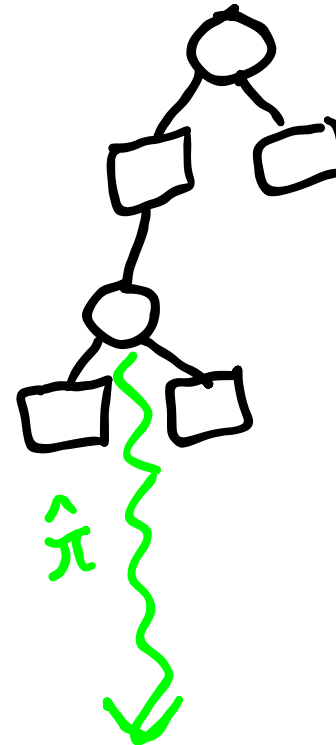
Search



Expansion



Rollout



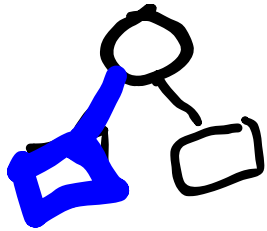
Backup

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

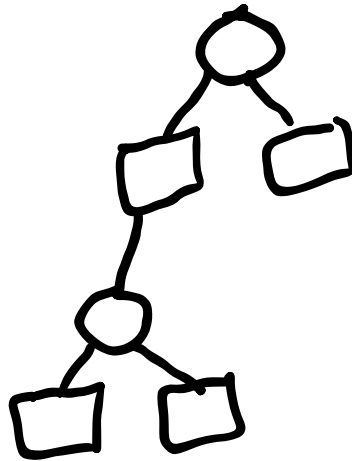
low  $N(s, a)/N(s)$  = high bonus  
start with  $c = 2(\bar{V} - \underline{V})$

# Monte Carlo Tree Search (MCTS/UCT)

Search



Expansion



Rollout



Backup

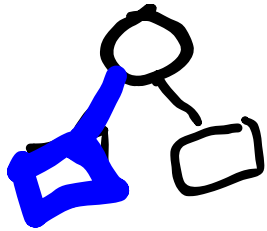
$$Q(s, a) + c \sqrt{\frac{\log N(s)}{N(s, a)}}$$

low  $N(s, a)/N(s)$  = high bonus

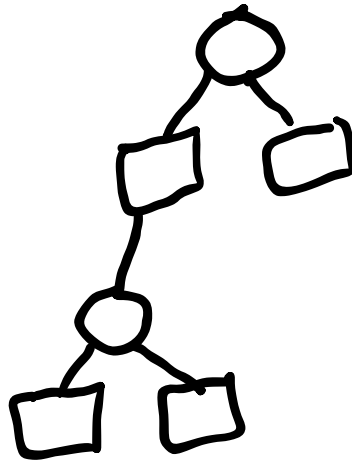
start with  $c = 2(\bar{V} - \underline{V})$

# Monte Carlo Tree Search (MCTS/UCT)

Search



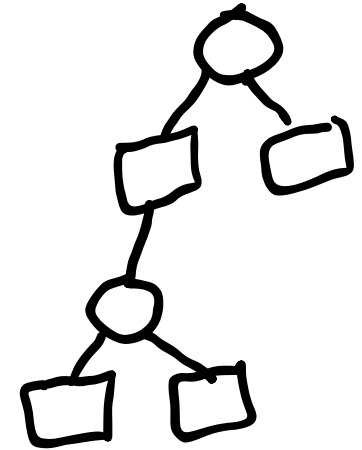
Expansion



Rollout



Backup

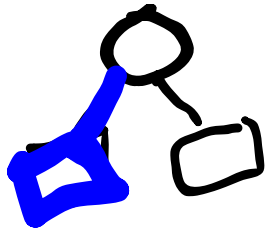


$$Q(s, a) + c \sqrt{\frac{\log N(s)}{N(s, a)}}$$

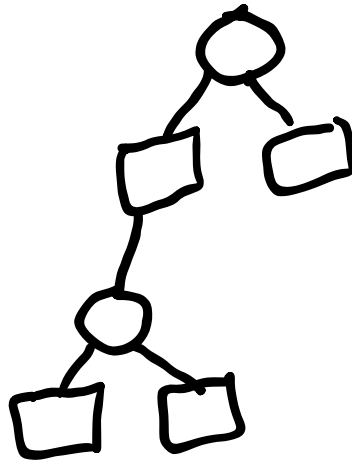
low  $N(s, a)/N(s)$  = high bonus  
start with  $c = 2(\bar{V} - \underline{V})$

# Monte Carlo Tree Search (MCTS/UCT)

Search



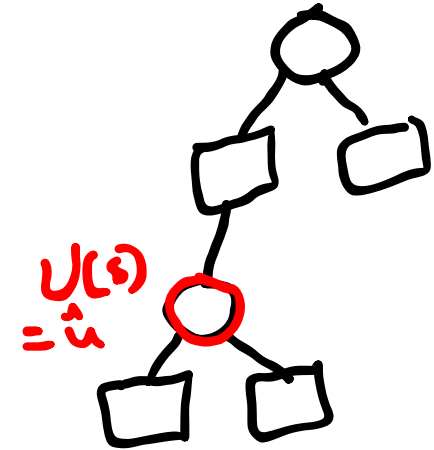
Expansion



Rollout



Backup

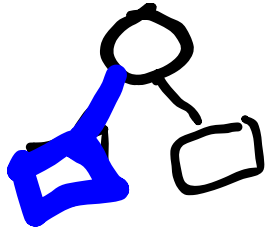


$$Q(s, a) + c \sqrt{\frac{\log N(s)}{N(s, a)}}$$

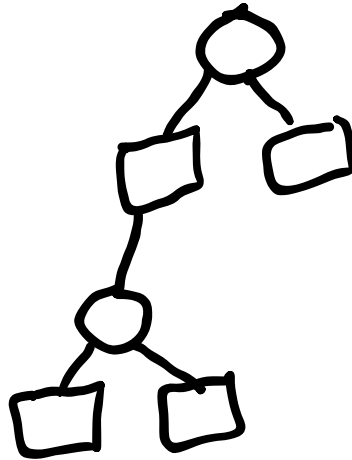
low  $N(s, a)/N(s)$  = high bonus  
start with  $c = 2(\bar{V} - \underline{V})$

# Monte Carlo Tree Search (MCTS/UCT)

Search



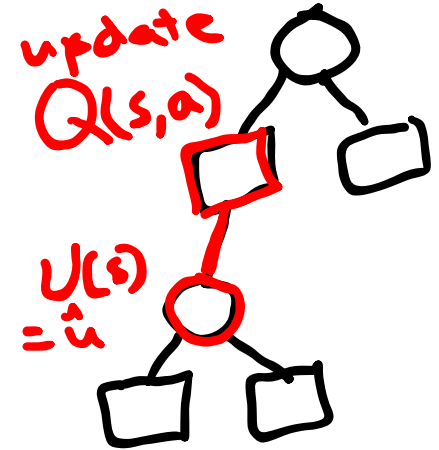
Expansion



Rollout

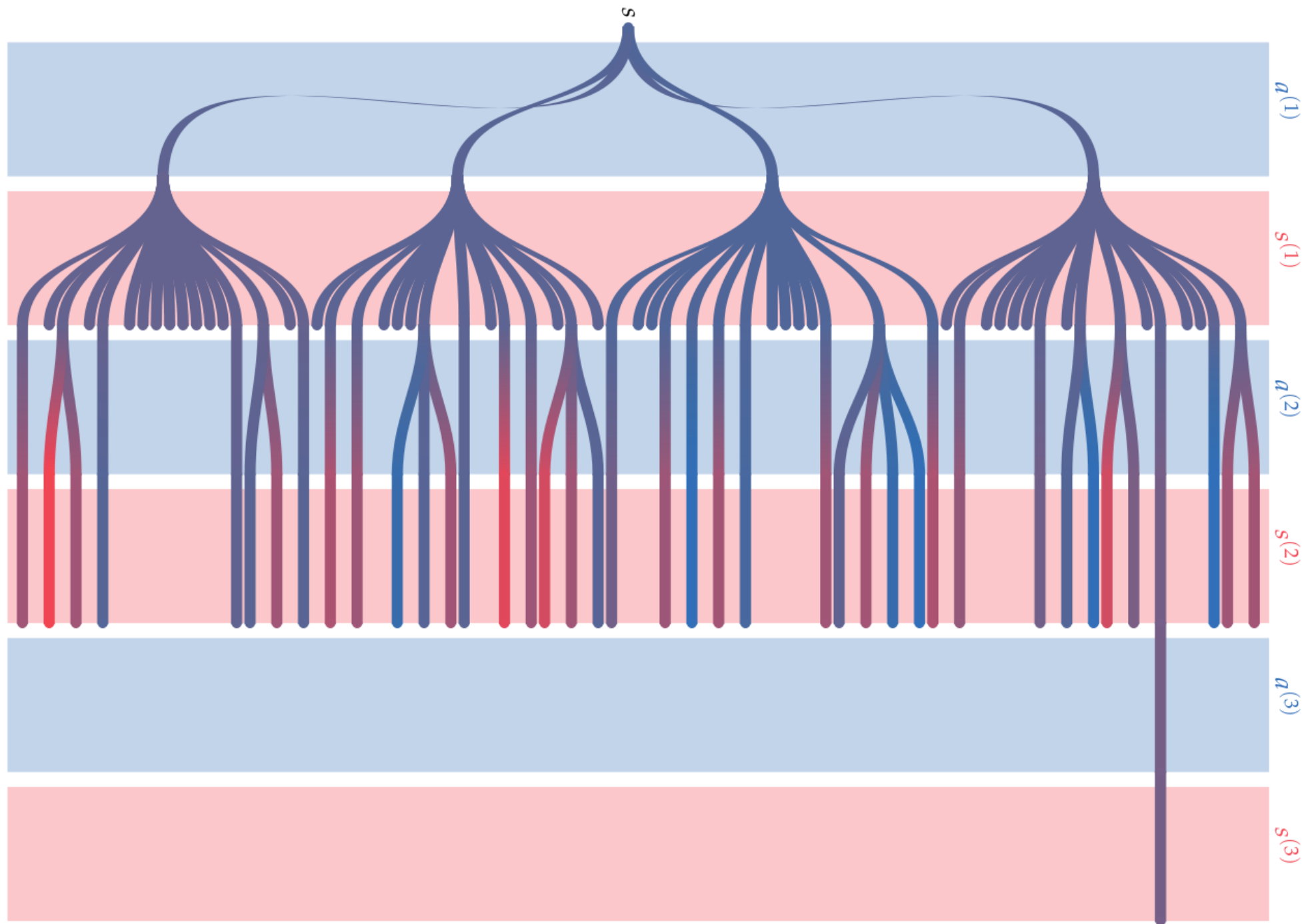


Backup



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

low  $N(s, a)/N(s)$  = high bonus  
start with  $c = 2(\bar{V} - \underline{V})$



# Guiding Questions



# Guiding Questions

- What are the differences between online and offline solutions?
- Are there solution techniques that are *independent* of the state space size?