# Last Time

# Last Time

- Does value iteration always converge? $\quad$ Yes!
  Because
  $\quad$ B is a
  $\quad$ contraction
- Is the value function unique?

any number of optimal policies

every optimal policy achieves $U^*$

$\pi^1 \quad U^{\pi^1} = U^*$

$\pi^2 \quad U^{\pi^2} = U^*$
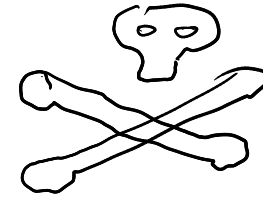
$A = \{ left, right+1, right+2 \}$

# Guiding Questions

# Guiding Questions

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

# Curse of Dimensionality

# Curse of Dimensionality

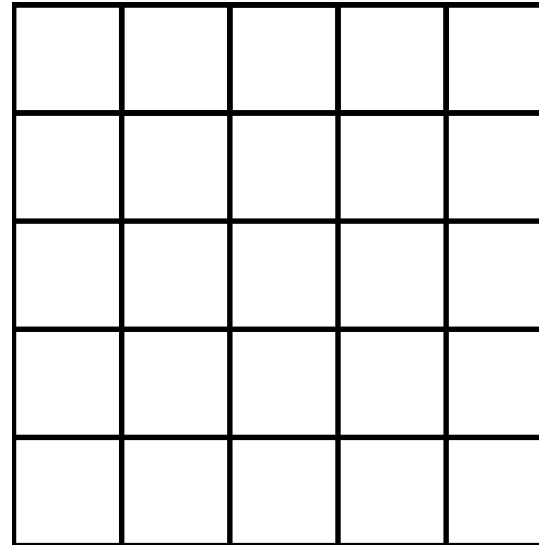1 dimension, 5 segments

$$|\mathcal{S}| = 5$$

# Curse of Dimensionality

1 dimension, 5 segments

$$|\mathcal{S}| = 5$$

2 dimensions, 5 segments

$$|\mathcal{S}| = 25$$

# Curse of Dimensionality
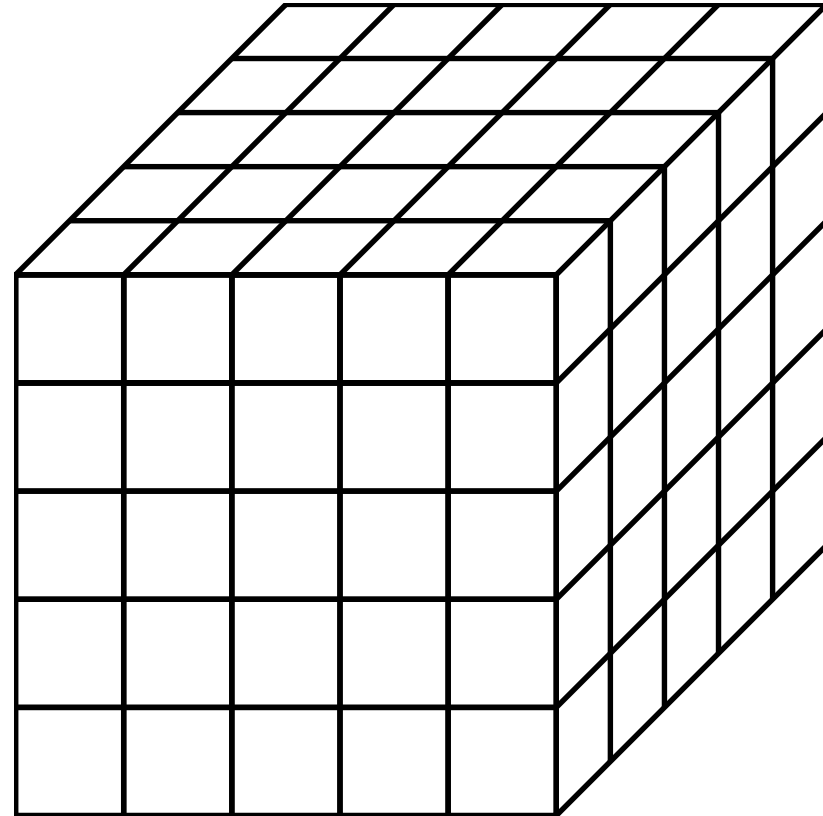
1 dimension, 5 segments

$$|\mathcal{S}| = 5$$

2 dimensions, 5 segments

$$|\mathcal{S}| = 25$$

3 dimensions, 5 segments

$$|\mathcal{S}| = 125$$

# Curse of Dimensionality
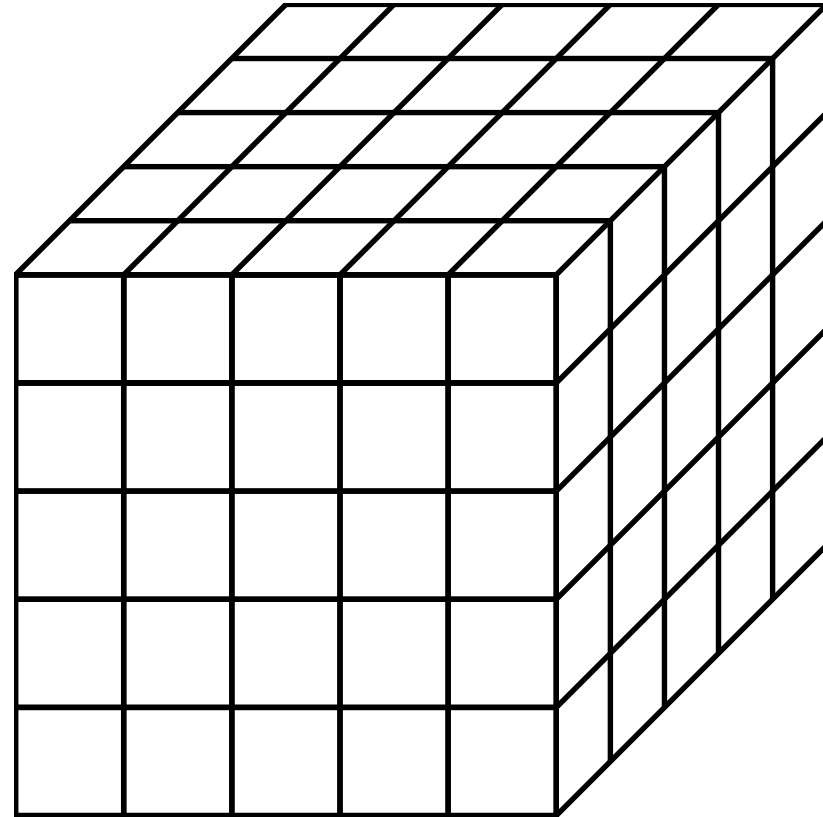
1 dimension, 5 segments

$$|\mathcal{S}| = 5$$

2 dimensions, 5 segments

$$|\mathcal{S}| = 25$$

3 dimensions, 5 segments

$$|\mathcal{S}| = 125$$

$$n \text{ dimensions, } k \text{ segments } \rightarrow |\mathcal{S}| = k^n$$

# Offline vs Online Solutions

Offline

Online

# Offline vs Online Solutions

Offline

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

Online
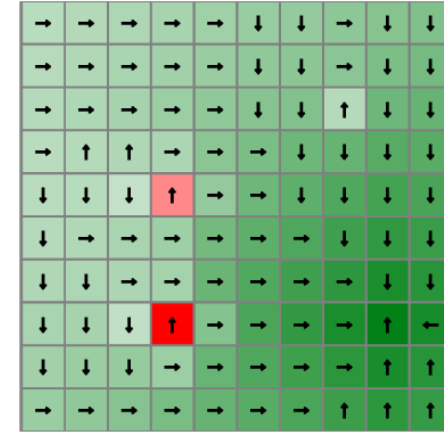
# Offline vs Online Solutions

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: $\pi^*(s) = \operatorname*{argmax} Q^*(s, a)$
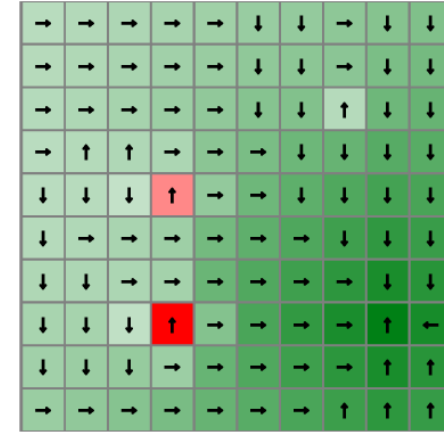
## Online

# Offline vs Online Solutions



## Offline

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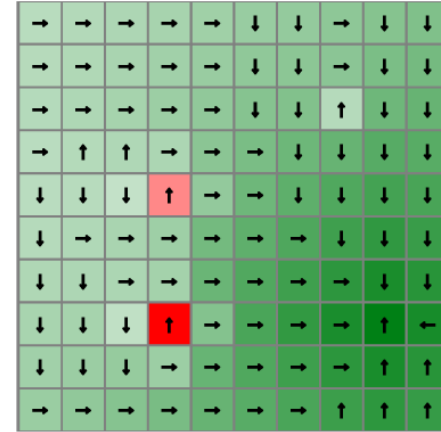- During Execution: $\pi^*(s) = \operatorname{argmax} Q^*(s, a)$

## Online

- Before Execution: <nothing>

# Offline vs Online Solutions

## Offline

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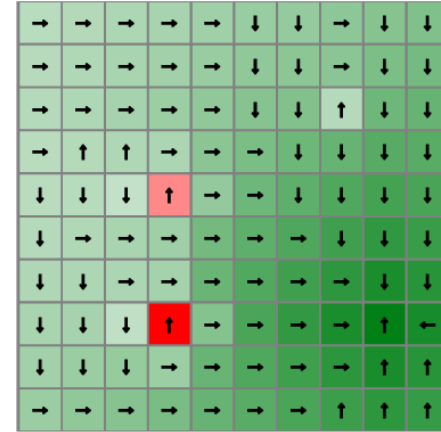


## Online

- Before Execution: \<nothing\>
- During Execution: Consider actions and their consequences (everything)
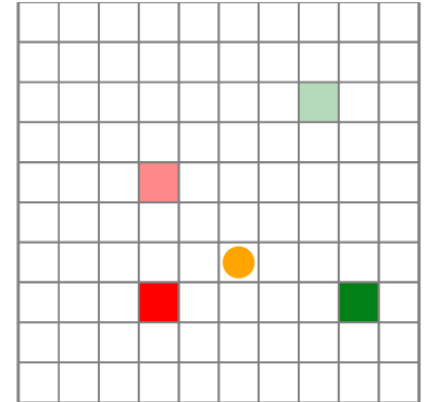
# Offline vs Online Solutions

## Offline

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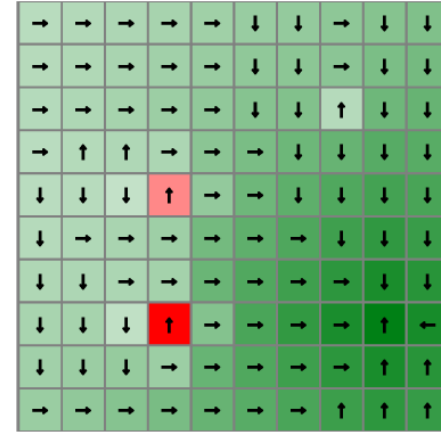
## Online

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

# Offline vs Online Solutions

<u>Offline</u>

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



<u>Online</u>

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

- Why?

# Offline vs Online Solutions
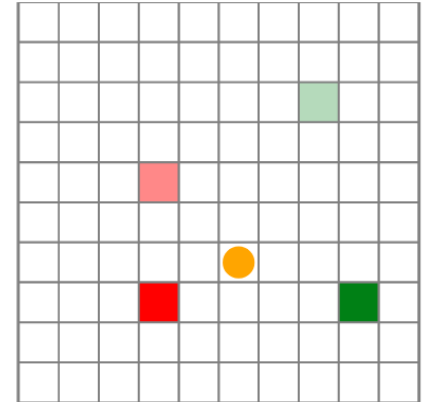
## Offline

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

## Online

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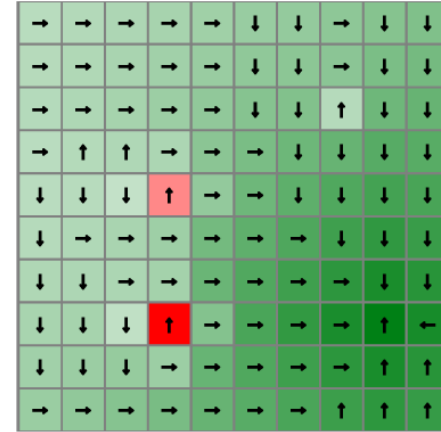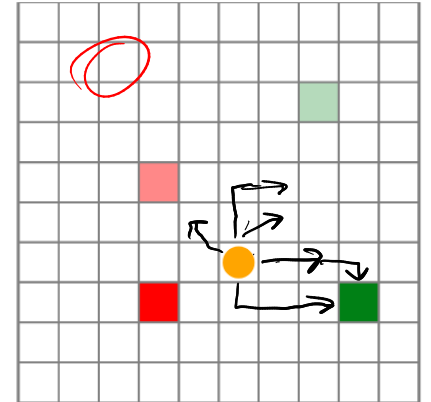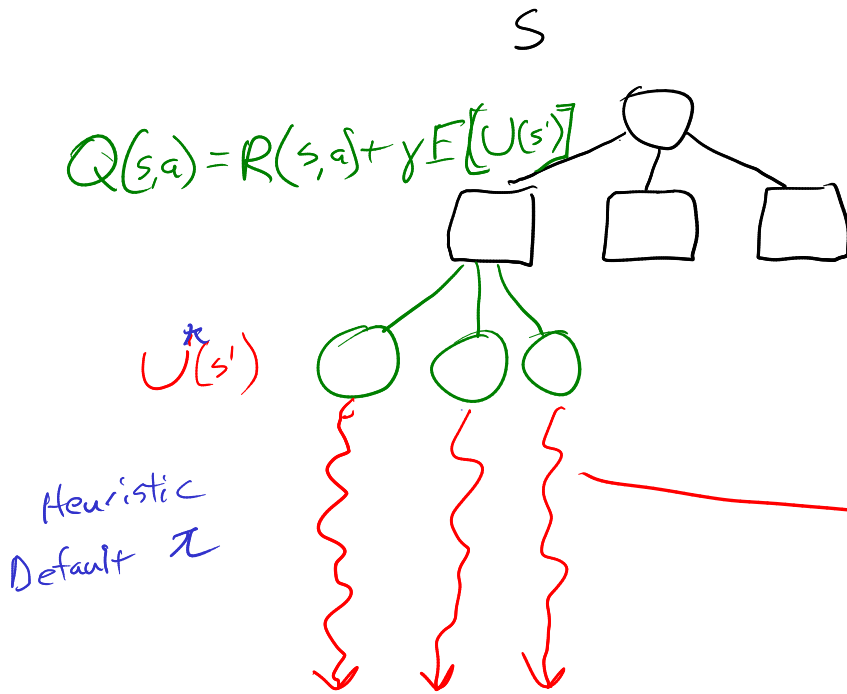- During Execution: Consider actions and their consequences (everything)

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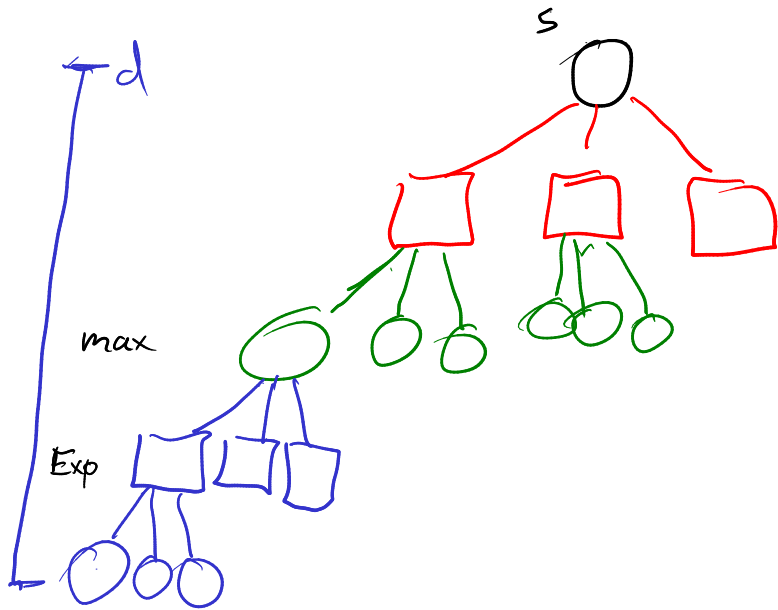- Online methods are insensitive to the size of $S$ !

# One Step Lookahead

$S$

$$Q(s,a) = R(s,a) + \gamma E[U(s')]$$

$\overset{\pi}{U(s')}$

Heuristic
Default $\pi$



```
randstep(𝒫::MDP, s, a) = 𝒫.TR(s, a)

function rollout(𝒫, s, π, d)
    ret = 0.0
    for t in 1:d
        a = π(s)
        s, r = randstep(𝒫, s, a)
        ret += 𝒫.γ^(t-1) * r
    end
    return ret
end


function (π::RolloutLookahead)(s)
    U(s) = rollout(π.𝒫, s, π.π, π.d)
    return greedy(π.𝒫, U, s).a
end


function greedy(𝒫::MDP, U, s)
    u, a = findmax(a→lookahead(𝒫, U, s, a), 𝒫.𝒜)
    return (a=a, u=u)
end
```

5

# Forward Search



$Q(s,a)$

$U = \max_a Q(s,a)$

$Q = R(s,a) + \gamma E[U(s')]$

$U$

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

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

sad!

6

# Forward Search depth

# Forward Search depth

# Forward Search depth



Depth 1

Depth 2

Depth 3

Depth 4

# Sparse Sampling

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

# Sparse Sampling

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

$$O\left((m|A|)^d\right)$$

# Sparse Sampling

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

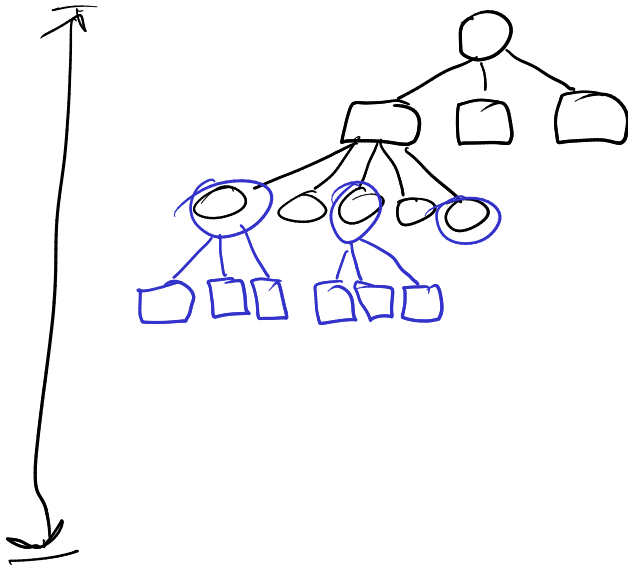$$O\left((m|A|)^d\right) \qquad |V^{\text{SS}}(s) - V^*(s)| \le \epsilon$$

# Sparse Sampling

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

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

# Sparse Sampling



$m = 3$

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

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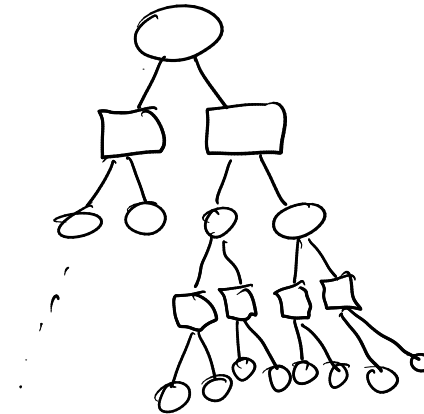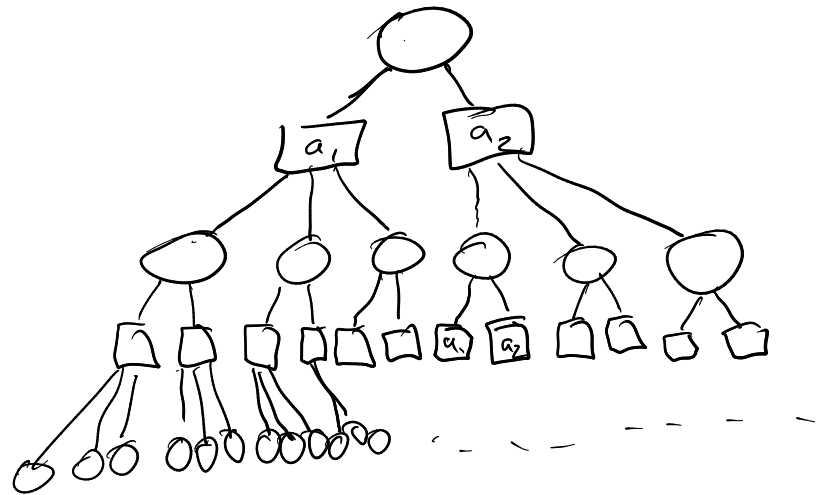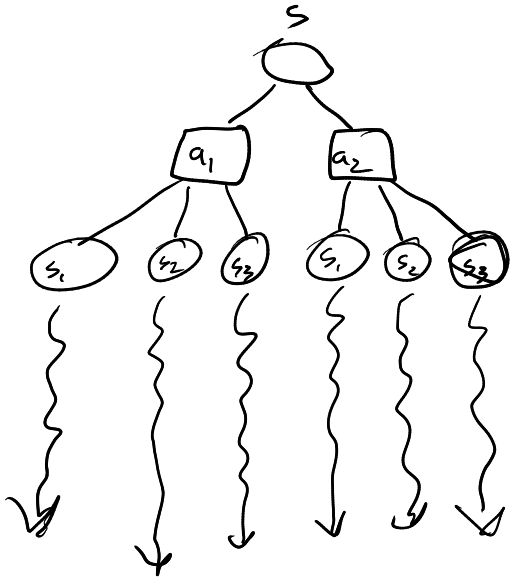
mad!

not on exam

# Break

Draw the trees produced by the following algorithms for a problem with 2 actions and 3 states:

1. One-step lookahead with rollout
2. Forward search (d=2)
3. Sparse sampling (d=2, m=2)

# Branch and Bound

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

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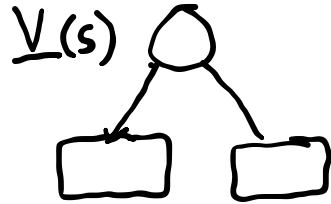
# Branch and Bound

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

$\underline{V}(s)$ ◯

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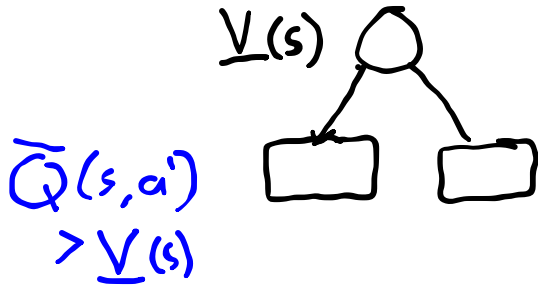
# Branch and Bound

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

$\underline{V}(s)$

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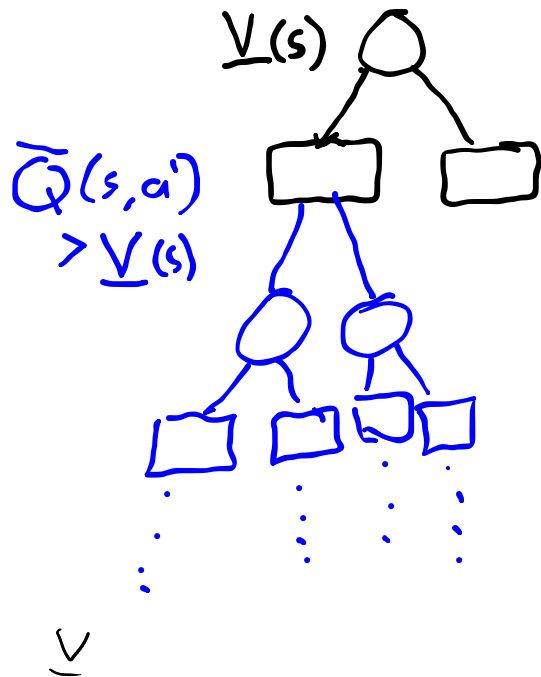
# Branch and Bound

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

$\underline{V}(s)$

$\bar{Q}(s,a')$
$> \underline{V}(s)$

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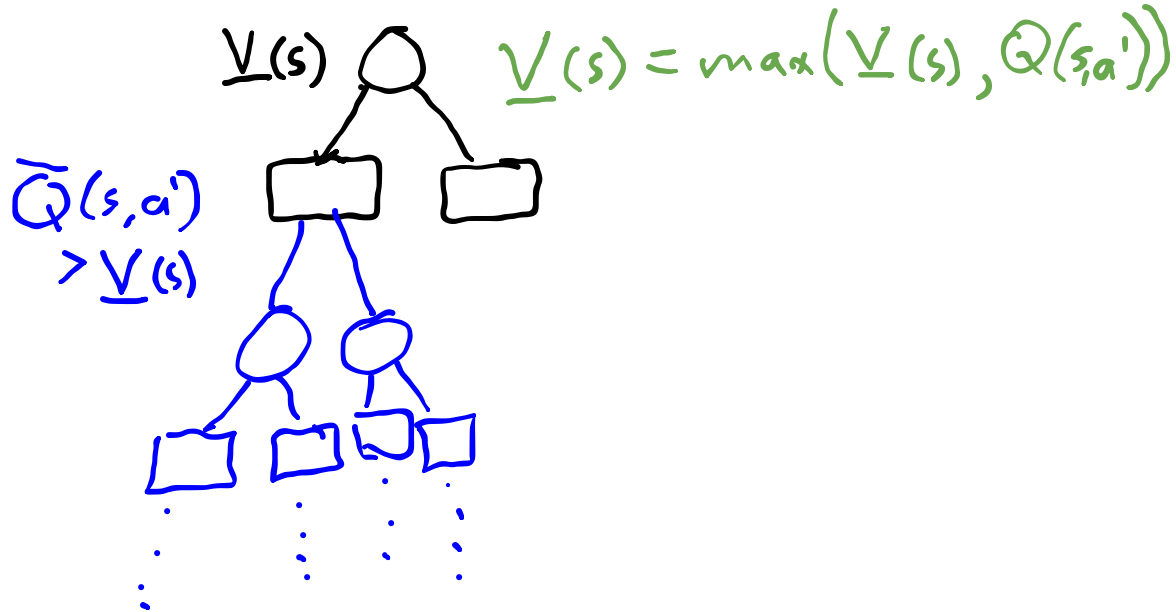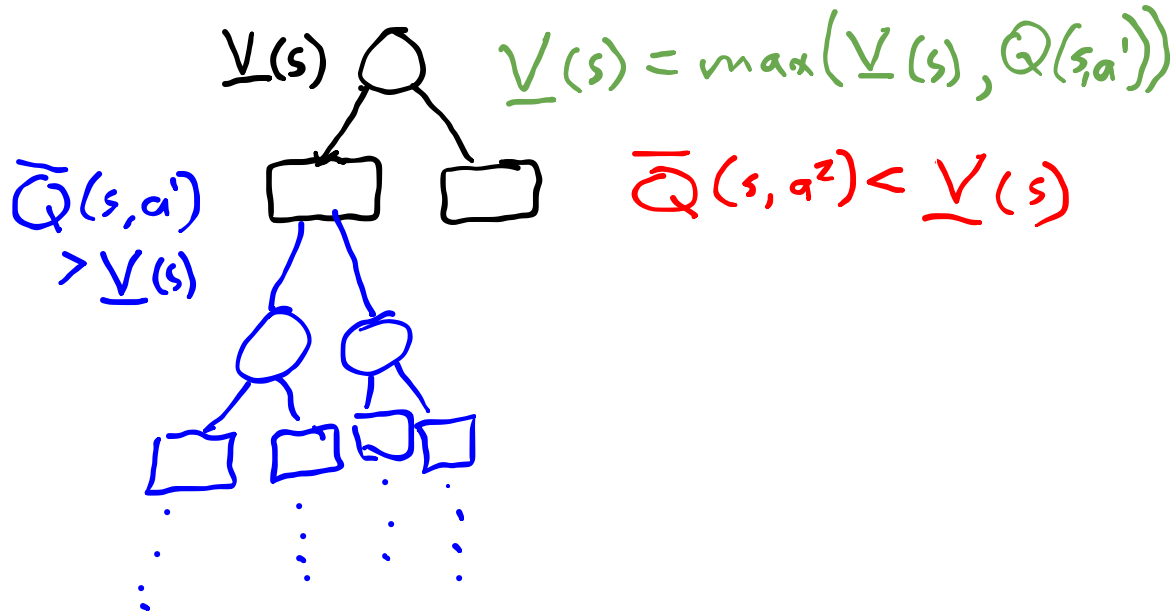
# Branch and Bound

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



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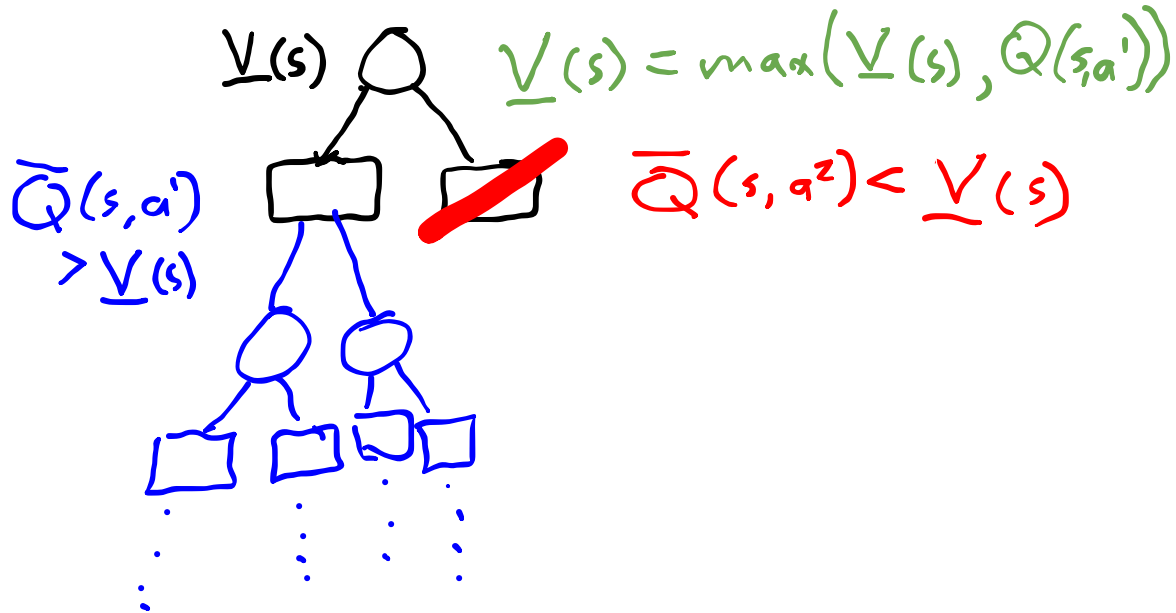
# Branch and Bound

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

$\underline{V}(s)$

$$\underline{V}(s) = \max\left(\underline{V}(s), \bar{Q}(s,a')\right)$$

$\bar{Q}(s,a')$
$> \underline{V}(s)$

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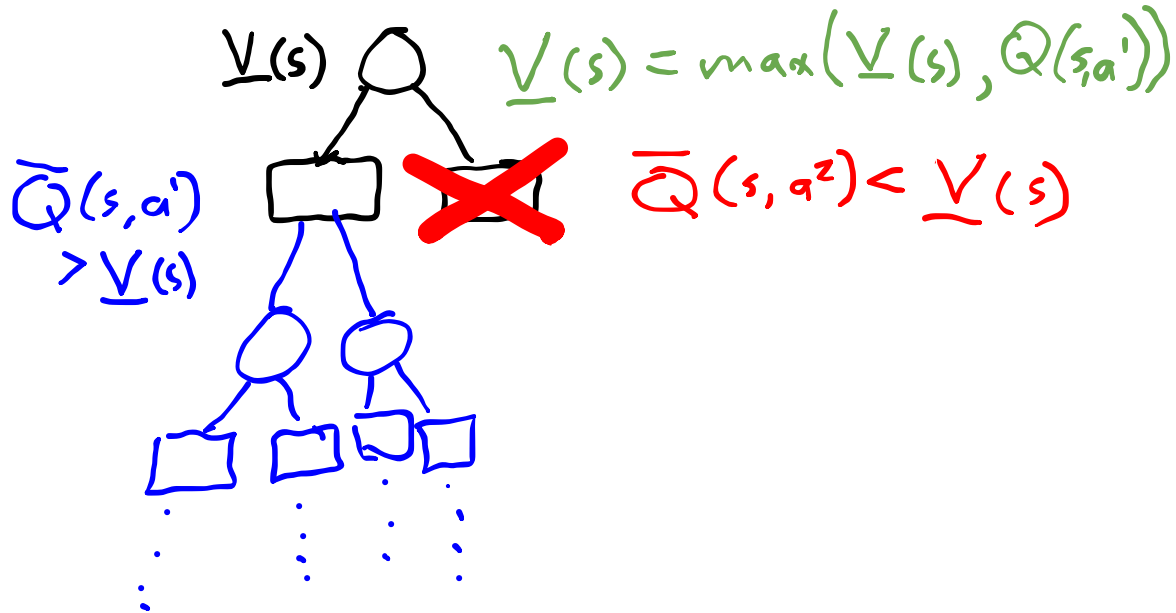
# Branch and Bound

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



$\underline{V}(s) = \max(\underline{V}(s), \underline{Q}(s,a'))$

$\bar{Q}(s,a^2) < \underline{V}(s)$

$\bar{Q}(s,a') > \underline{V}(s)$

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



$$\underline{V}(s) = \max\left(\underline{V}(s), \underline{Q}(s, a')\right)$$

$$\bar{Q}(s, a^2) < \underline{V}(s)$$

$\underline{V}(s)$

$\bar{Q}(s, a')$
$> \underline{V}(s)$

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



$\underline{V}(s) = \max\left(\underline{V}(s), \underline{Q}(s,a')\right)$

$\underline{V}(s)$

$\bar{Q}(s,a') > \underline{V}(s)$

$\bar{Q}(s, a^2) < \underline{V}(s)$

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

# 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)

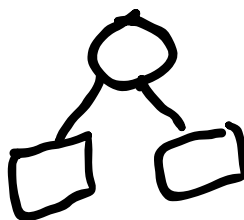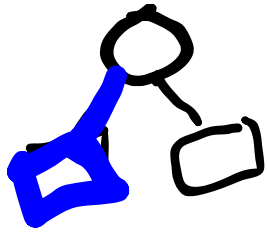Search                Expansion                Rollout                Backup

# 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
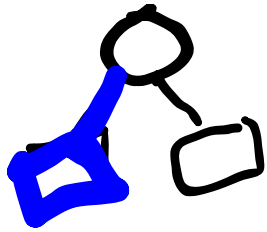
Rollout

Backup

more Common

more theoretically justified

$$Q(s,a) + c\sqrt{\frac{\log N(s)}{N(s,a)}} \quad \text{or} \quad Q(s,a) + c\frac{N(s)^\beta}{\sqrt{N(s,a)}}$$

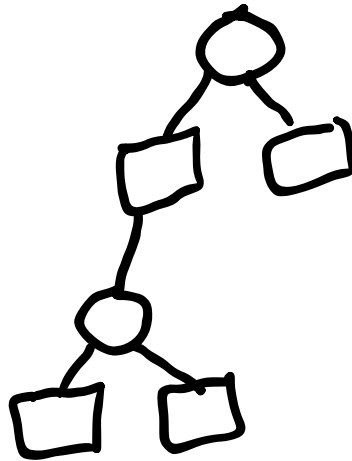low $N(s,a)/N(s)$ = high bonus

start with $c = 2(\bar{V} - \underline{V})$, $\beta = 1/4$

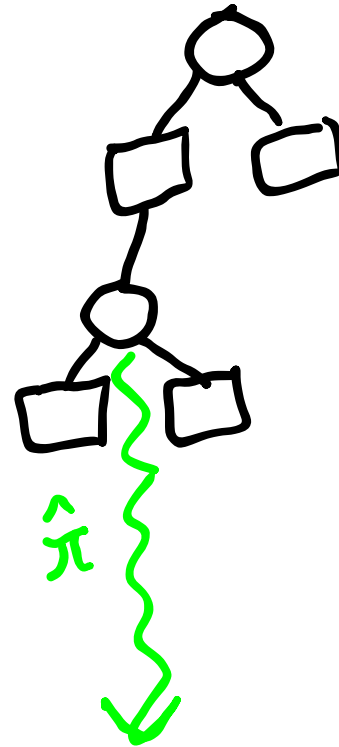# Monte Carlo Tree Search (MCTS/UCT)

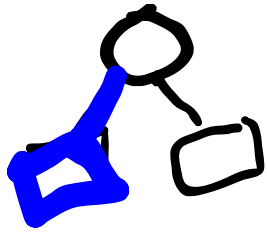Search          Expansion          Rollout          Backup

$$Q(s,a) + c\sqrt{\frac{\log N(s)}{N(s,a)}} \quad \text{or} \quad Q(s,a) + c\frac{N(s)^{\beta}}{\sqrt{N(s,a)}}$$

low $N(s,a)/N(s)$ = high bonus
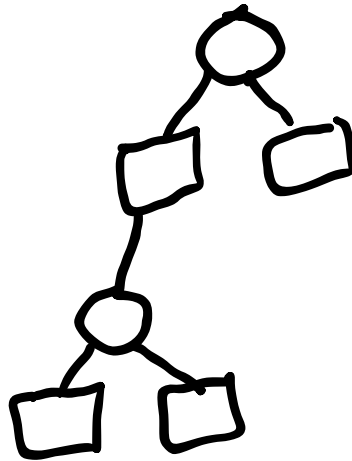
start with $c = 2(\bar{V} - \underline{V})$, $\beta = 1/4$

# Monte Carlo Tree Search (MCTS/UCT)

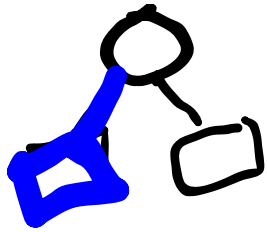Search                Expansion                Rollout                Backup



$$Q(s,a) + c\sqrt{\frac{\log N(s)}{N(s,a)}} \quad \text{or} \quad Q(s,a) + c\frac{N(s)^{\beta}}{\sqrt{N(s,a)}}$$

low $N(s,a)/N(s)$ = high bonus
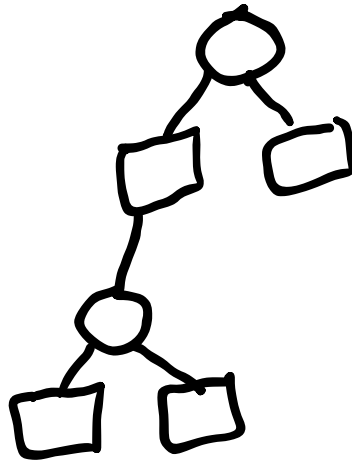
start with $c = 2(\bar{V} - \underline{V})$, $\beta = 1/4$

# Monte Carlo Tree Search (MCTS/UCT)

Search               Expansion               Rollout               Backup

$$Q(s,a) + c\sqrt{\frac{\log N(s)}{N(s,a)}} \quad \text{or} \quad Q(s,a) + c\frac{N(s)^{\beta}}{\sqrt{N(s,a)}}$$

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

start with $c = 2(\bar{V} - \underline{V})$, $\beta = 1/4$

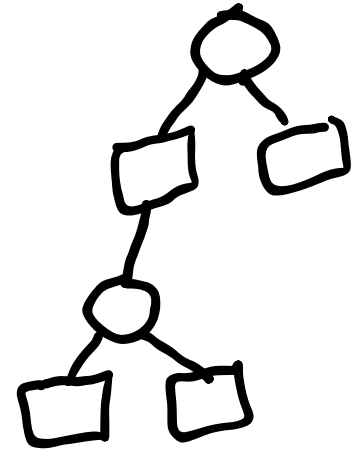# Monte Carlo Tree Search (MCTS/UCT)

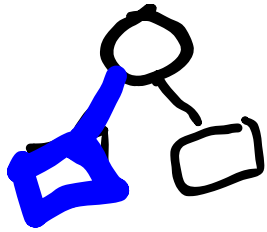Search              Expansion              Rollout              Backup



$$Q(s,a) + c\sqrt{\frac{\log N(s)}{N(s,a)}} \quad \text{or} \quad Q(s,a) + c\frac{N(s)^{\beta}}{\sqrt{N(s,a)}}$$

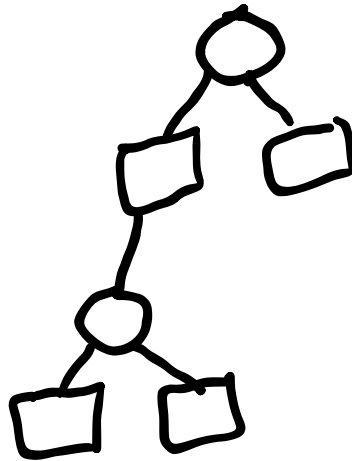low $N(s,a)/N(s)$ = high bonus

start with $c = 2(\bar{V} - \underline{V})$, $\beta = 1/4$

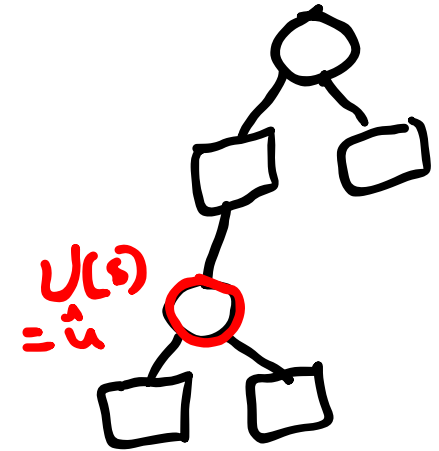# Monte Carlo Tree Search (MCTS/UCT)

Search         Expansion         Rollout         Backup



$$Q(s,a) + c\sqrt{\frac{\log N(s)}{N(s,a)}} \quad \text{or} \quad Q(s,a) + c\frac{N(s)^{\beta}}{\sqrt{N(s,a)}}$$

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

start with $c = 2(\bar{V} - \underline{V})$, $\beta = 1/4$

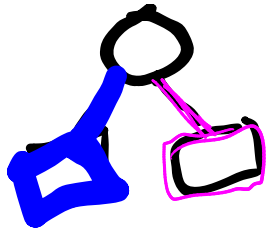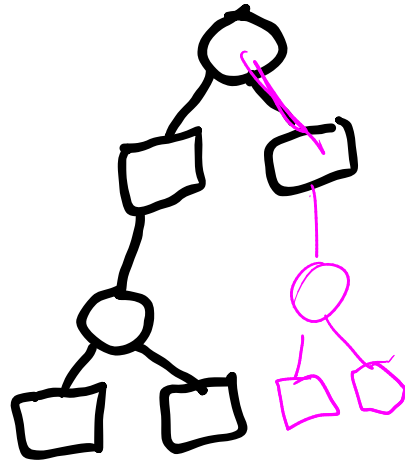# Monte Carlo Tree Search (MCTS/UCT)

Search

Expansion

Rollout

Backup



$$Q(s,a) + c\sqrt{\frac{\log N(s)}{N(s,a)}} \quad \text{or} \quad Q(s,a) + c\frac{N(s)^{\beta}}{\sqrt{N(s,a)}}$$

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

start with $c = 2(\bar{V} - \underline{V})$, $\beta = 1/4$
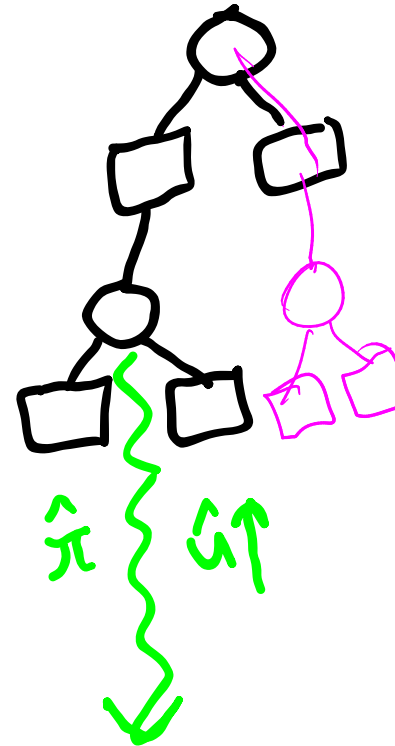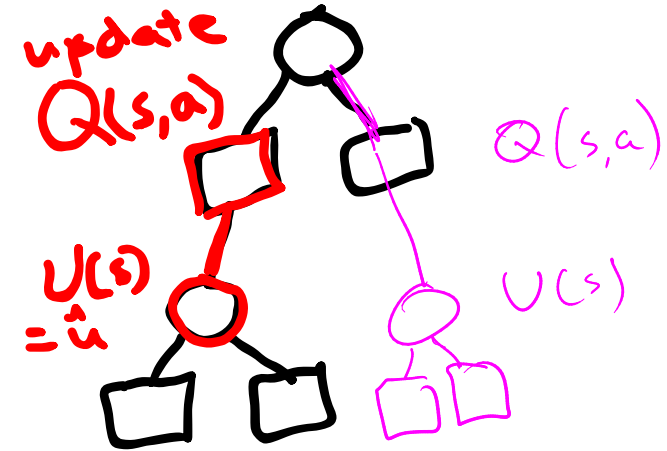
# Monte Carlo Tree Search (MCTS/UCT)
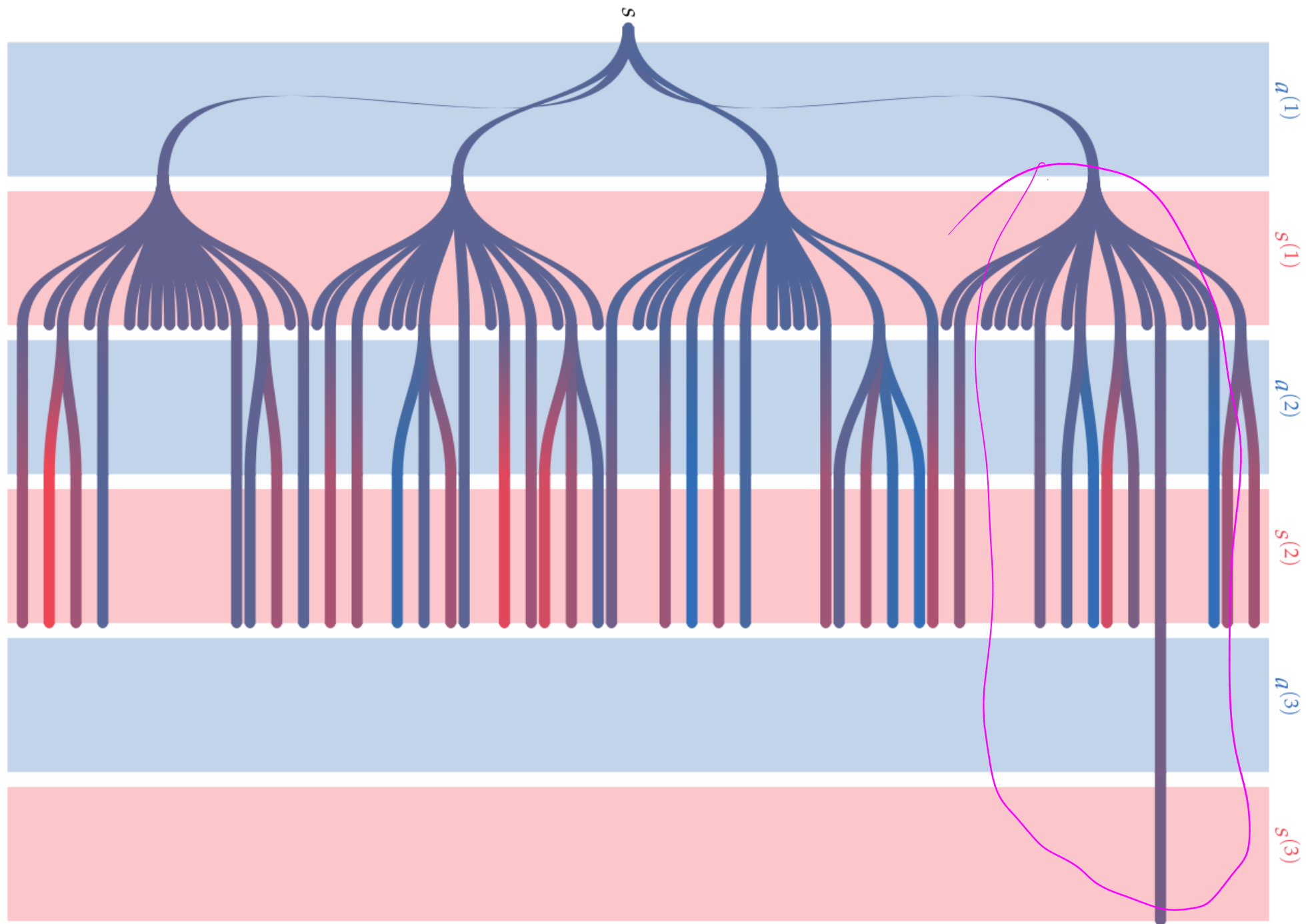


Search       Expansion       Rollout       Backup

$$Q(s,a) + c\sqrt{\frac{\log N(s)}{N(s,a)}} \quad \text{or} \quad Q(s,a) + c\frac{N(s)^{\beta}}{\sqrt{N(s,a)}}$$

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

start with $c = 2(\bar{V} - \underline{V})$, $\beta = 1/4$

11.15

# 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?