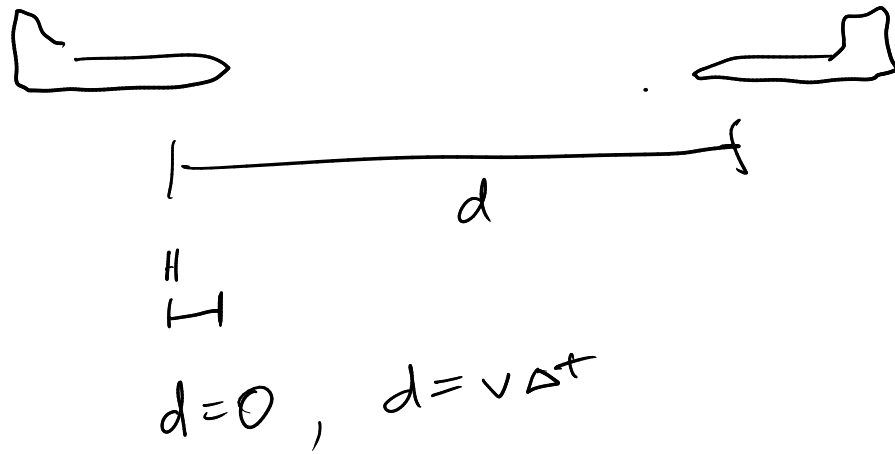


Last Time



Last Time

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

Guiding Questions

Guiding Questions

- What tools do we have to solve MDPs with continuous S and A ?

Continuous S and A

Continuous S and A

e.g. $S \subseteq \mathbb{R}^n, A \subseteq \mathbb{R}^m$

Continuous S and A

e.g. $S \subseteq \mathbb{R}^n, A \subseteq \mathbb{R}^m$

The old rules still work!

$$V^*(s) = \max_a \left(R(s, a) + \gamma \int_S T(s'|s, a) V^*(s') \right)$$
$$V^\pi(s) = R(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim T(s, \pi(s))} [V^\pi(s')]$$

Linear Dynamics, Quadratic Reward

Linear Dynamics, Quadratic Reward

$$s' = T_s s + T_a a + w$$

Linear Dynamics, Quadratic Reward

$$s' = T_s s + T_a a + w$$

w is zero-mean, finite variance R.V.

Linear Dynamics, Quadratic Reward

$$s' = T_s s + T_a a + w$$

w is zero-mean, finite variance

$$R(s, a) = s^\top R_s s + a^\top R_a a$$

LQR

Linear Dynamics, Quadratic Reward

$$s' = T_s s + T_a a + w$$

w is zero-mean, finite variance

$$R(s, a) = s^\top R_s s + a^\top R_a a$$

$$V_h(s) = s^\top P_h s + q_h$$

Linear Dynamics, Quadratic Reward

$$s' = T_s s + T_a a + w$$

w is zero-mean, finite variance

$$R(s, a) = s^\top R_s s + a^\top R_a a$$

$$V_h(s) = s^\top P_h s + q_h$$

$$P_{h+1} = R_s + \underbrace{T_s^\top P_h^\top T_s}_{\text{}} - \left(\underbrace{T_a^\top P_h T_s}_{\text{}} \right)^\top \left(R_a + \underbrace{T_a^\top P_h T_a}_{\text{}} \right)^{-1} \underbrace{\left(T_a^\top P_h T_s \right)}_{\text{}}$$

Linear Dynamics, Quadratic Reward

$$s' = T_s s + T_a a + w$$

w is zero-mean, finite variance

$$R(s, a) = s^\top R_s s + a^\top R_a a$$

$$V_h(s) = s^\top P_h s + q_h$$

$$P_{h+1} = R_s + T_s^\top P_h^\top T_s - (T_a^\top P_h T_s)^\top (R_a + T_a^\top P_h T_a)^{-1} (T_a^\top P_h T_s)$$

$$\pi_h^*(s) = - \underbrace{(R_a + T_a^\top P_h T_a)^{-1} T_a^\top P_h T_s}_K s = -K s$$

Linear Dynamics, Quadratic Reward

$$s' = T_s s + T_a a + w$$

w is zero-mean, finite variance

$$R(s, a) = s^\top R_s s + a^\top R_a a$$

$$V_h(s) = s^\top P_h s + \underline{\underline{q_h}}$$

$$P_{h+1} = R_s + T_s^\top P_h^\top T_s - (T_a^\top P_h T_s)^\top (R_a + T_a^\top P_h T_a)^{-1} (T_a^\top P_h T_s)$$

$$\pi_h^*(s) = - (R_a + T_a^\top P_h T_a)^{-1} T_a^\top P_h T_s s = -K s$$

if $w \sim \mathcal{N}(0, \Sigma)$, then $q_{h+1} = \Sigma_{i=1}^h \text{tr}(\Sigma P_i)$

Linear Dynamics, Quadratic Reward

$$s' = T_s s + T_a a + w$$

w is zero-mean, finite variance

$$R(s, a) = s^\top R_s s + a^\top R_a a$$

$$V_h(s) = s^\top P_h s + q_h$$

$$P_{h+1} = R_s + T_s^\top P_h^\top T_s - (T_a^\top P_h T_s)^\top (R_a + T_a^\top P_h T_a)^{-1} (T_a^\top P_h T_s)$$

$$\pi_h^*(s) = - (R_a + T_a^\top P_h T_a)^{-1} T_a^\top P_h T_s s = -K s$$

if $w \sim \mathcal{N}(0, \Sigma)$, then $q_{h+1} = \Sigma_{i=1}^h \text{tr}(\Sigma P_i)$

"Linear Quadratic Regulator" / "LQR"

If not Linear Quadratic...

Offline:

Online:

If not Linear Quadratic...

Offline:

- Approximate Dynamic Programming (ADP)

Online:

If not Linear Quadratic...

Offline:

- Approximate Dynamic Programming (ADP)

V_θ

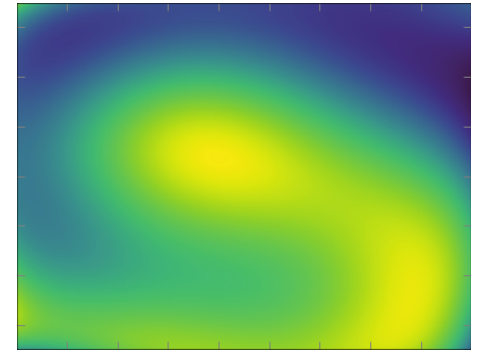
Online:

If not Linear Quadratic...

Offline:

- Approximate Dynamic Programming (ADP)

V_θ



Online:

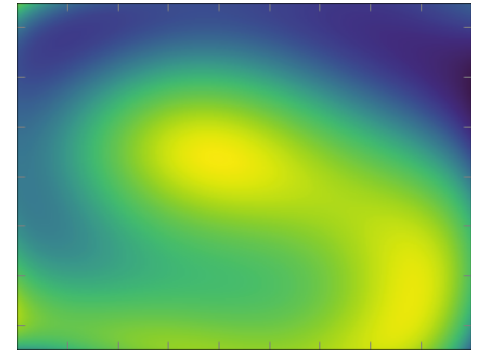
If not Linear Quadratic...

Offline:

- Approximate Dynamic Programming (ADP)
- Policy Search

Online:

V_θ



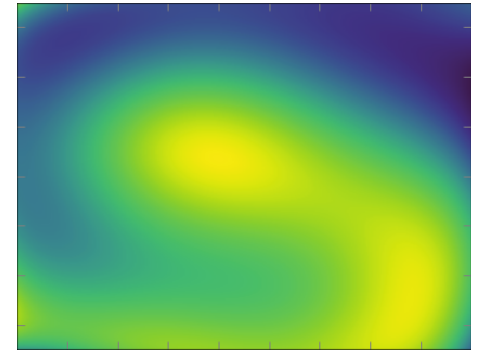
If not Linear Quadratic...

Offline:

- Approximate Dynamic Programming (ADP)
- Policy Search

Online:

V_θ



π_θ

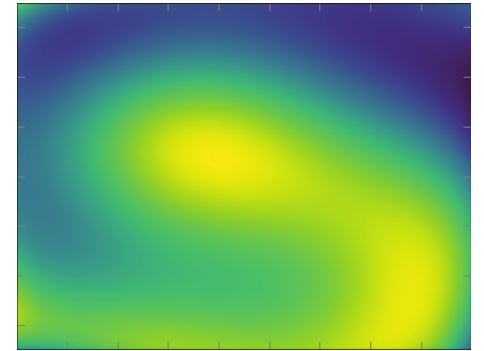
If not Linear Quadratic...

Offline:

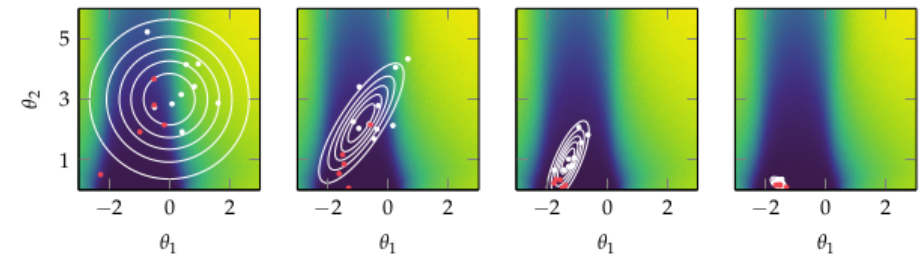
- Approximate Dynamic Programming (ADP)
- Policy Search

Online:

V_θ



π_θ



If not Linear Quadratic...

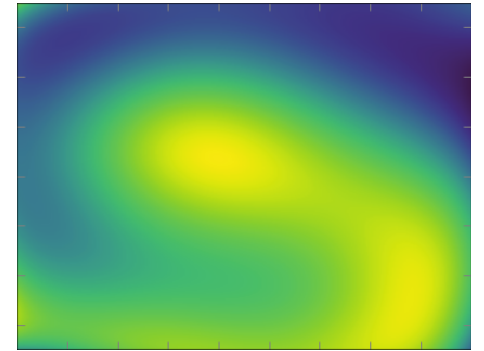
Offline:

- Approximate Dynamic Programming (ADP)
- Policy Search

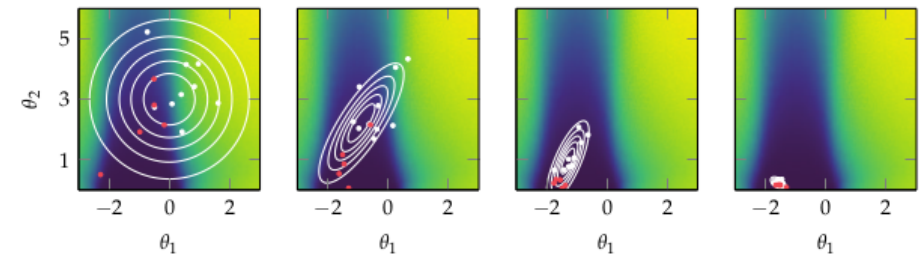
Online:

- Model Predictive Control (MPC)

V_θ



π_θ



If not Linear Quadratic...

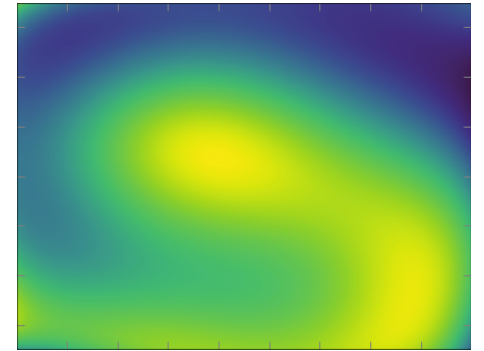
Offline:

- Approximate Dynamic Programming (ADP)
- Policy Search

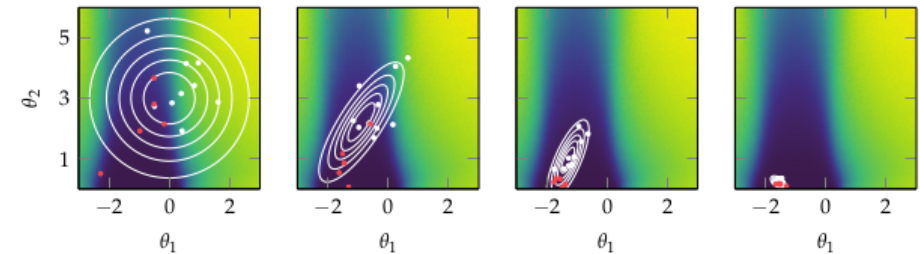
Online:

- Model Predictive Control (MPC)

V_θ



π_θ



If not Linear Quadratic...

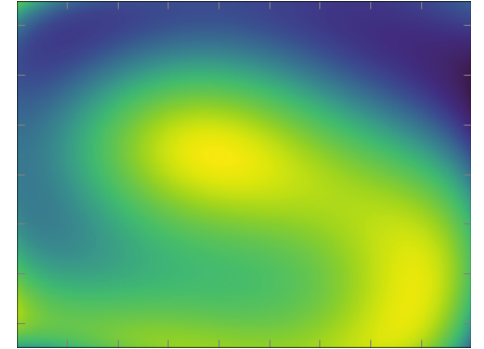
Offline:

- Approximate Dynamic Programming (ADP)
- Policy Search

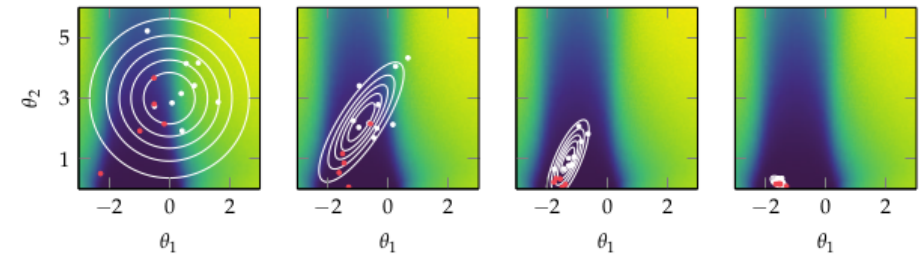
Online:

- Model Predictive Control (MPC)
- Sparse Tree Search/Progressive Widening

V_θ



π_θ

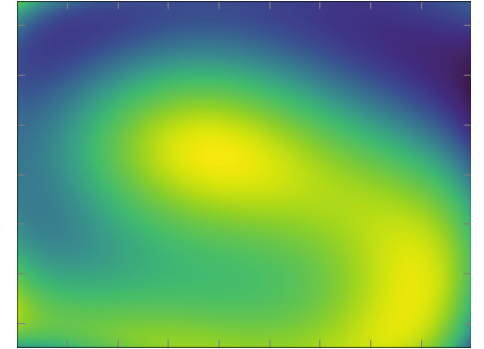


If not Linear Quadratic...

Offline:

- Approximate Dynamic Programming (ADP)
- Policy Search

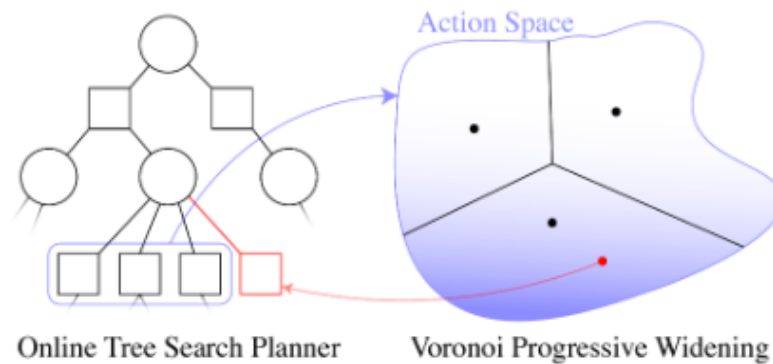
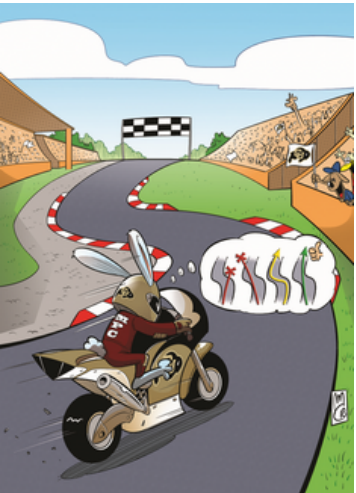
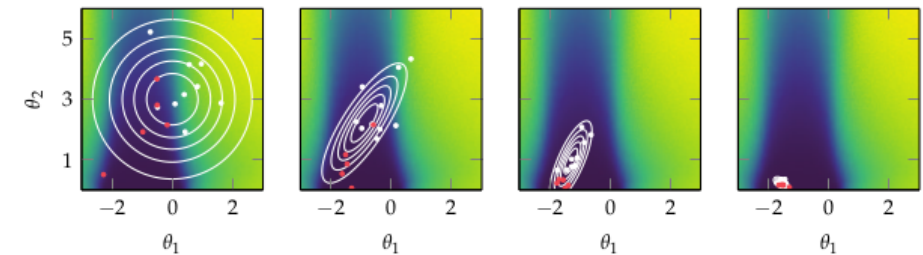
V_θ



Online:

- Model Predictive Control (MPC)
- Sparse Tree Search/Progressive Widening

π_θ



Online Tree Search Planner

Voronoi Progressive Widening

Approximate Dynamic Programming

Approximate Dynamic Programming

$$V_{\theta}(s) = f_{\theta}(s) \quad (\text{e.g. neural network})$$

Approximate Dynamic Programming

$$V_{\theta}(s) = f_{\theta}(s) \quad (\text{e.g. neural network})$$

$$V_{\theta}(s) = \theta^{\top} \beta(s) \quad (\text{linear feature})$$

Approximate Dynamic Programming

$$V_{\theta}(s) = f_{\theta}(s) \quad (\text{e.g. neural network})$$

$$V_{\theta}(s) = \theta^{\top} \beta(s) \quad (\text{linear feature})$$

while not converged

$$\theta \leftarrow \theta'$$

$$\hat{V}' \leftarrow \underbrace{B_{\text{approx}}[V_{\theta}]}$$

$$\theta' \leftarrow \text{fit}(\hat{V}')$$

Approximate Dynamic Programming

$$V_{\theta}(s) = f_{\theta}(s) \quad (\text{e.g. neural network})$$

$$V_{\theta}(s) = \theta^{\top} \beta(s) \quad (\text{linear feature})$$

$$\left. \begin{array}{l} \text{while not converged} \\ \theta \leftarrow \theta' \\ \hat{V}' \leftarrow B_{\text{approx}}[V_{\theta}] \\ \theta' \leftarrow \text{fit}(\hat{V}') \end{array} \right\}$$

$$B_{\text{MC}(N)}[V_{\theta}](s) = \max_a \left(R(s, a) + \gamma \underbrace{\sum_{i=1}^N V_{\theta}(G(s, a, \downarrow w_i))} \right)$$

Approximate Dynamic Programming

$$V_{\theta}(s) = f_{\theta}(s) \quad (\text{e.g. neural network})$$

$$V_{\theta}(s) = \theta^{\top} \beta(s) \quad (\text{linear feature})$$

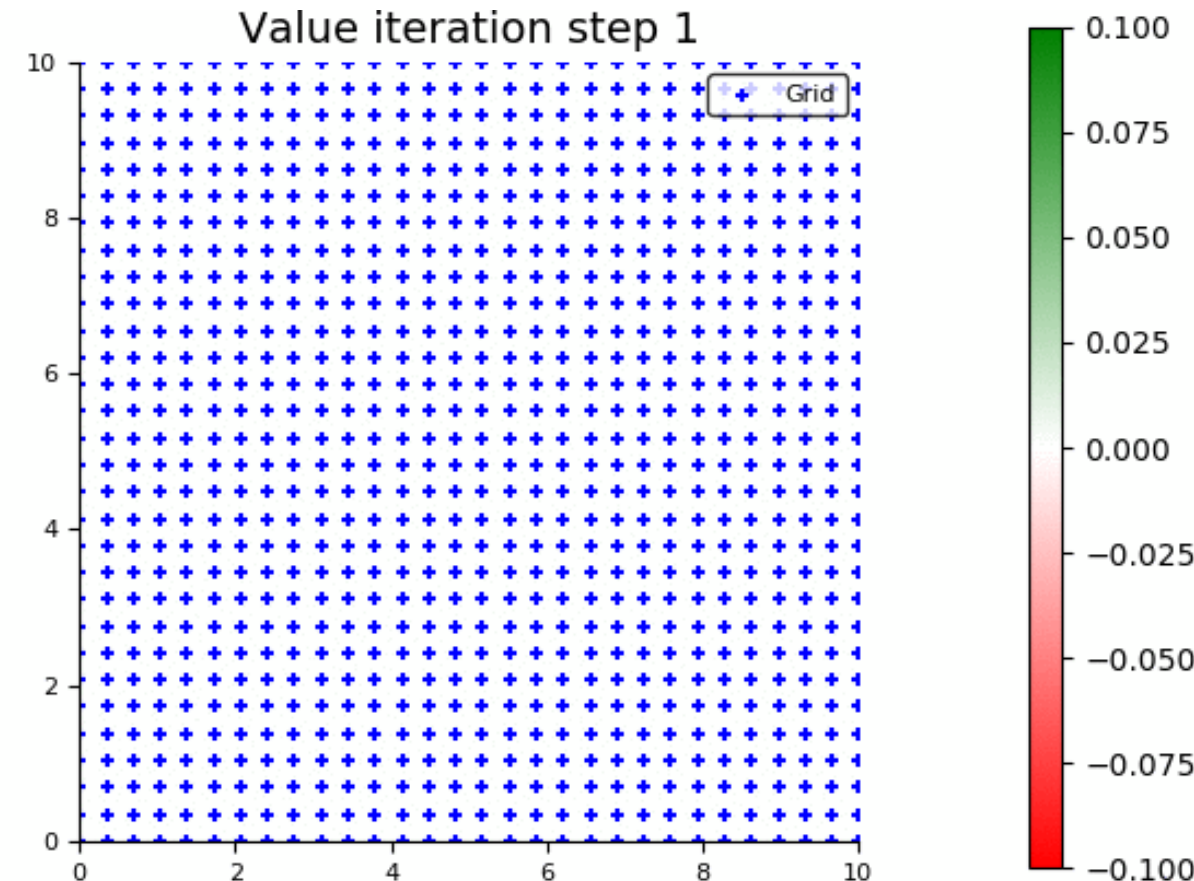
while not converged

$$\theta \leftarrow \theta'$$

$$\hat{V}' \leftarrow B_{\text{approx}}[V_{\theta}]$$

$$\theta' \leftarrow \text{fit}(\hat{V}')$$

LSVI



$$B_{\text{MC}(N)}[V_{\theta}](s) = \max_a \left(R(s, a) + \gamma \sum_{i=1}^N V_{\theta}(G(s, a, w_i)) \right)$$

Function Approximation

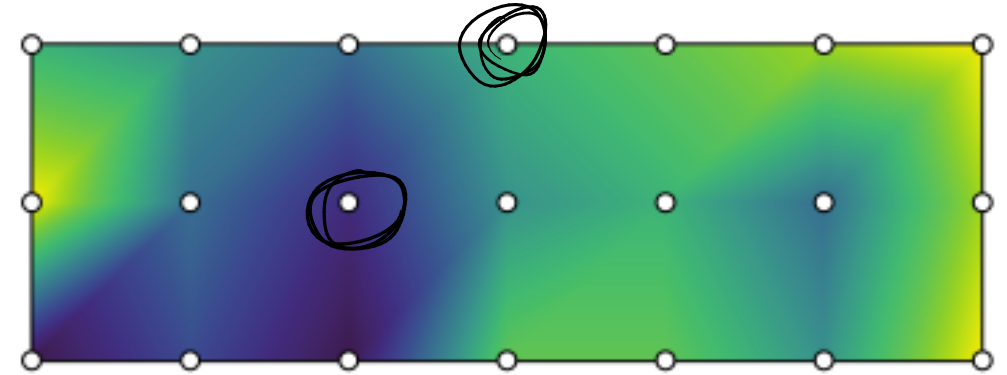
Function Approximation

- Local: (e.g. simplex interpolation)

Function Approximation

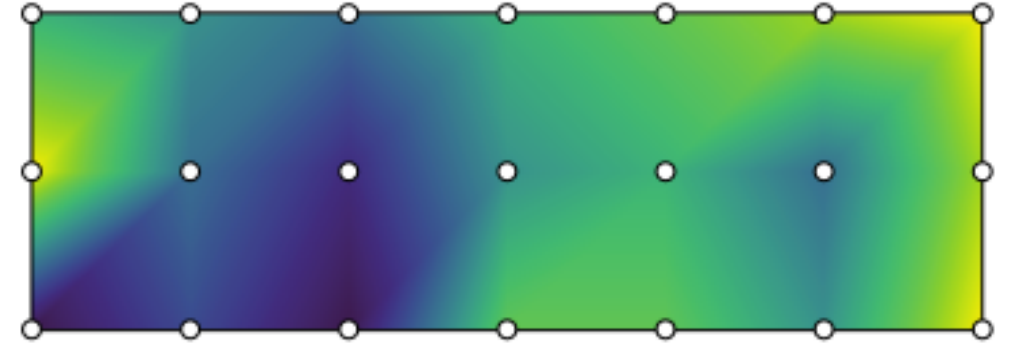


- Local: (e.g. simplex interpolation)



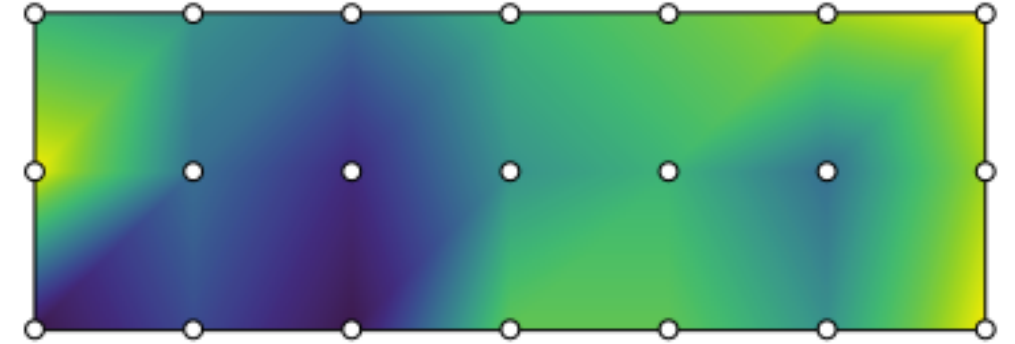
Function Approximation

- Local: (e.g. simplex interpolation)
- Global: (e.g. Fourier, neural network)



Function Approximation

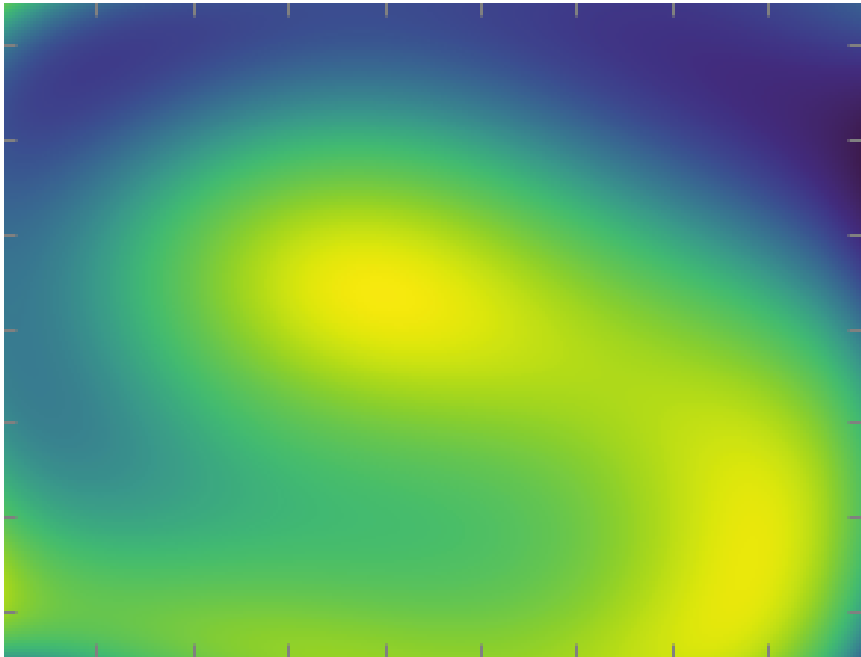
- Local: (e.g. simplex interpolation)
- Global: (e.g. Fourier, neural network)



$$s = (x, v)$$

$$V(s) = \Theta^T \beta(s)$$

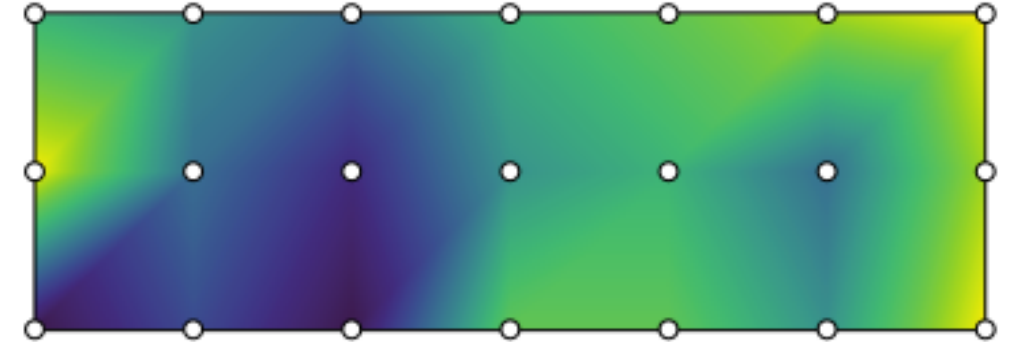
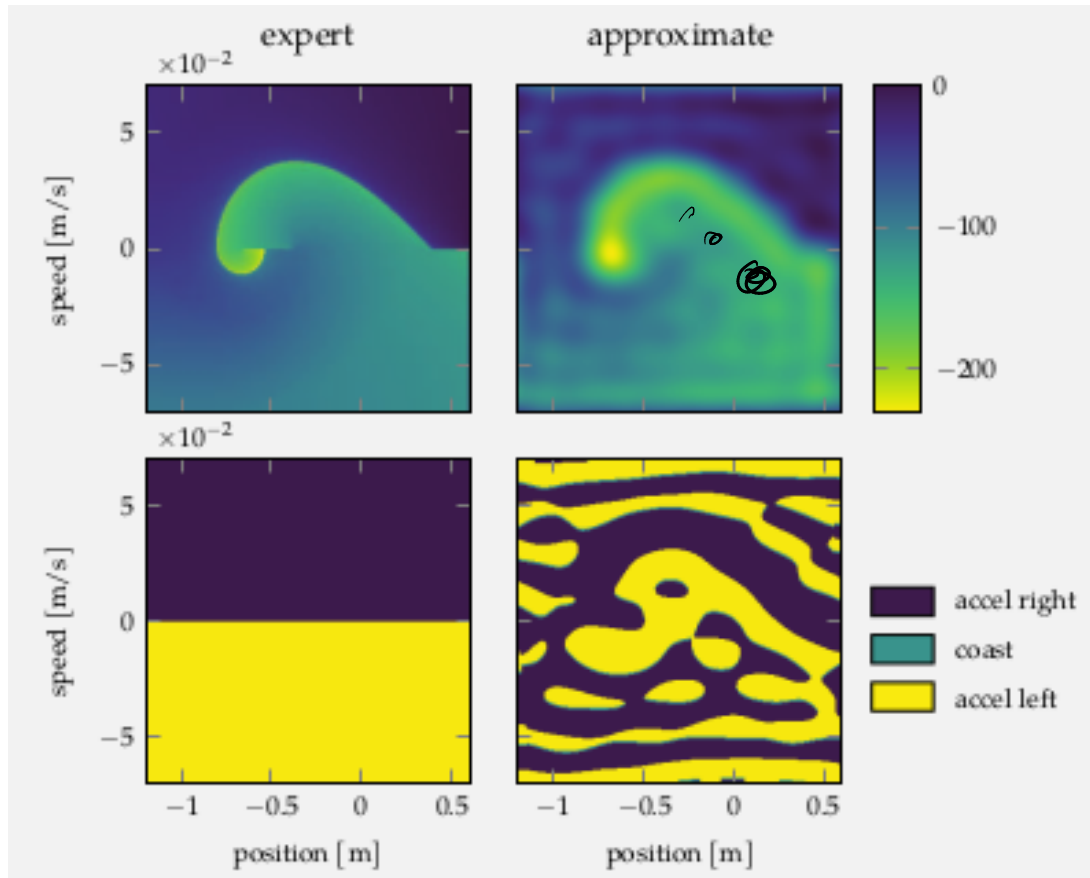
$$\beta(s) = \begin{bmatrix} 1, & \underline{x}, & \underline{v}, & & & & \\ x^2, & xv, & v^2, & & & & \\ x^3, & x^2v, & xv^2, & v^3, & & & \\ x^4, & x^3v, & x^2v^2, & xv^3, & v^4, & & \\ x^5, & x^4v, & x^3v^2, & x^2v^3, & xv^4, & v^5, & \\ x^6, & x^5v, & x^4v^2, & x^3v^3, & x^2v^4, & xv^5, & v^6 \end{bmatrix}$$



x

Function Approximation

- Local: (e.g. simplex interpolation)
- Global: (e.g. Fourier, neural network)



Policy Search

Policy Search

$$\underset{\theta}{\text{maximize}} \quad U(\pi_{\theta})$$

Policy Search

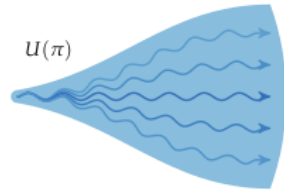
$$\underset{\theta}{\text{maximize}} \quad U(\pi_{\theta})$$

$$U(\pi) \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$$

Policy Search

$$\underset{\theta}{\text{maximize}} \quad U(\pi_{\theta})$$

$$U(\pi) \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$$

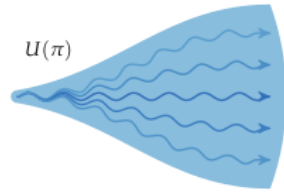


Policy Search

$$\underset{\theta}{\text{maximize}} \quad U(\pi_{\theta})$$

Common Approaches

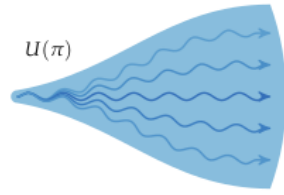
$$U(\pi) \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$$



Policy Search

$$\underset{\theta}{\text{maximize}} \quad U(\pi_{\theta})$$

$$U(\pi) \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$$



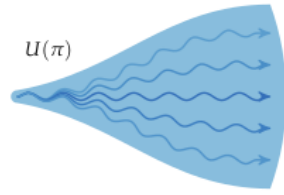
Common Approaches

- Evolutionary Algorithms

Policy Search

$$\underset{\theta}{\text{maximize}} \quad U(\pi_{\theta})$$

$$U(\pi) \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$$



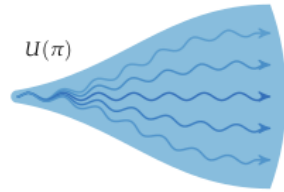
Common Approaches

- Evolutionary Algorithms
- Cross Entropy Method

Policy Search

$$\underset{\theta}{\text{maximize}} \quad U(\pi_{\theta})$$

$$U(\pi) \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$$



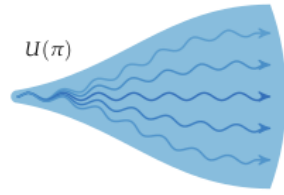
Common Approaches

- Evolutionary Algorithms
- Cross Entropy Method
- Policy Gradient (will cover in RL section)

Policy Search

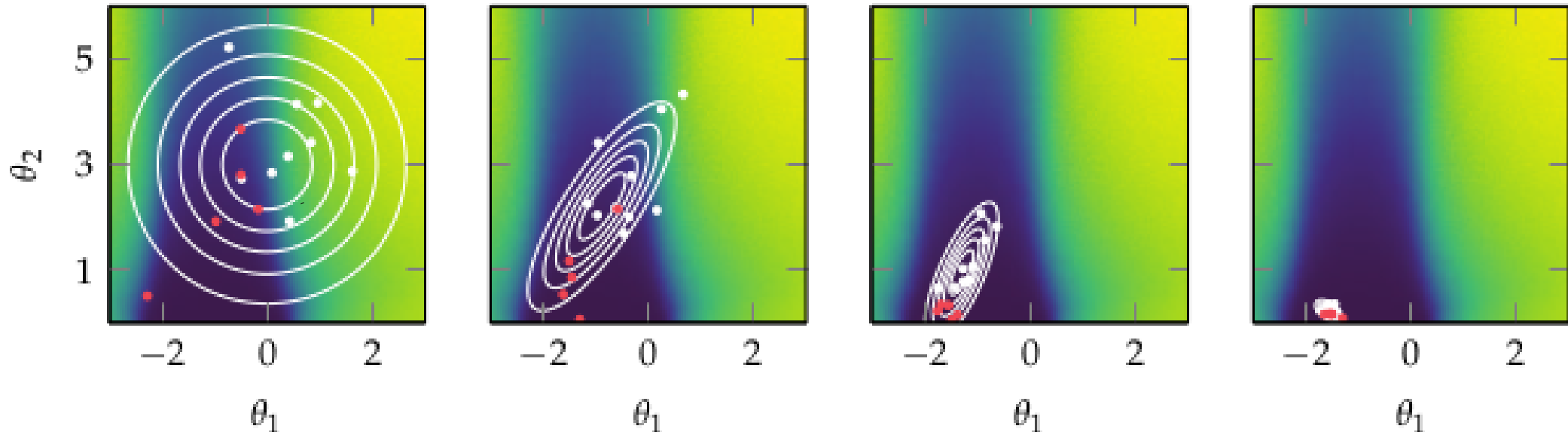
$$\underset{\theta}{\text{maximize}} \quad \underbrace{U(\pi_{\theta})}$$

$$U(\pi) \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$$



Common Approaches

- Evolutionary Algorithms
- Cross Entropy Method
- Policy Gradient (will cover in RL section)

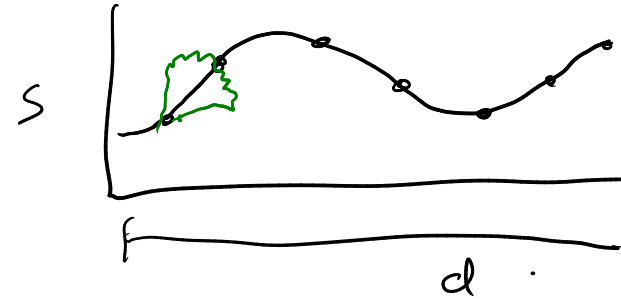


Model Predictive Control

Model Predictive Control

$$\begin{aligned}
 &\underset{a_{1:d}, s_{1:d}}{\text{maximize}} && \sum_{t=1}^d \gamma^t R(s_t, a_t) \\
 &\text{subject to} && \underline{s_{t+1} = \mathbb{E}[T(s_t, a_t)] \quad \forall t}
 \end{aligned}$$

mean



Model Predictive Control

$$\begin{aligned} & \underset{a_{1:d}, s_{1:d}}{\text{maximize}} && \sum_{t=1}^d \gamma^t R(s_t, a_t) \\ & \text{subject to} && s_{t+1} = \mathbb{E}[T(s_t, a_t)] \quad \forall t \end{aligned}$$

$$\begin{aligned} & \underset{a_{1:d}, s_{1:d}^{(1:m)}}{\text{maximize}} && \frac{1}{m} \sum_{i=1}^m \sum_{t=1}^d \gamma^t R(s_t^{(i)}, a_t) \\ & \text{subject to} && s_{t+1} = G(s_t^{(i)}, a_t, w_t^{(i)}) \quad \forall t, i \end{aligned}$$

Model Predictive Control

$$\underset{a_{1:d}, s_{1:d}}{\text{maximize}} \quad \sum_{t=1}^d \gamma^t R(s_t, a_t)$$

$$\text{subject to} \quad s_{t+1} = \mathbb{E}[T(s_t, a_t)] \quad \forall t$$

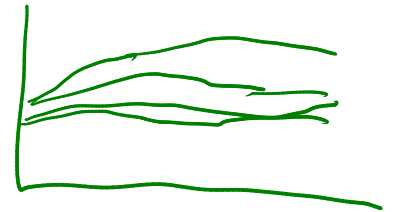
Optimistic

$$\underset{a_{1:d}, s_{1:d}^{(1:m)}}{\text{maximize}} \quad \frac{1}{m} \sum_{i=1}^m \sum_{t=1}^d \gamma^t R(s_t^{(i)}, a_t)$$

$$\text{subject to} \quad s_{t+1} = G(s_t^{(i)}, a_t, w_t^{(i)}) \quad \forall t, i$$

Open Loop

*Pessimistic
Conservative*



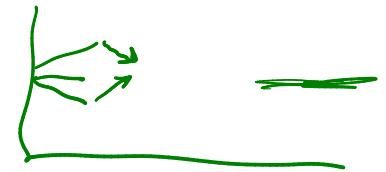
$$\underset{\boxed{a_{1:d}^{(1:m)}}, s_{1:d}^{(1:m)}}{\text{maximize}} \quad \frac{1}{m} \sum_{i=1}^m \sum_{t=1}^d \gamma^t R(s_t^{(i)}, a_t^{(i)})$$

$$\text{subject to} \quad s_{t+1} = G(s_t^{(i)}, a_t^{(i)}, w_t^{(i)}) \quad \forall t, i$$

$$a_1^{(i)} = a_1^{(j)} \quad \forall i, j$$

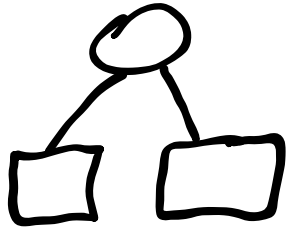
Hindsight Opt

Optimistic

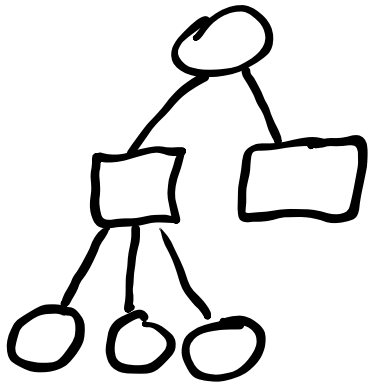


Sparse Tree Search/Progressive Widening

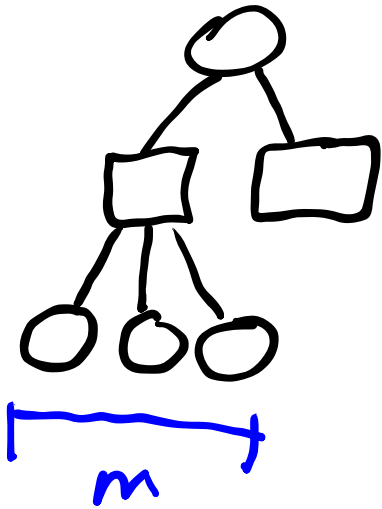
Sparse Tree Search/Progressive Widening



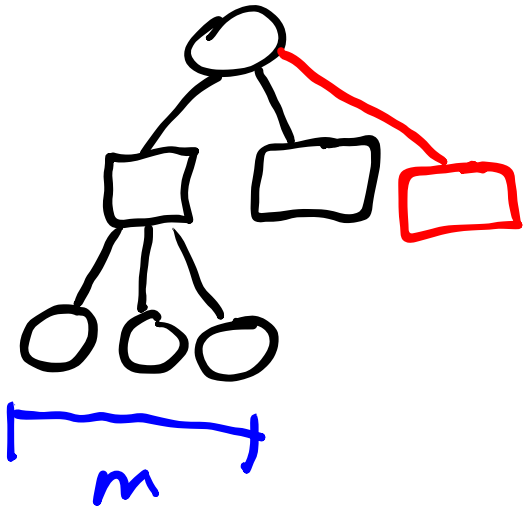
Sparse Tree Search/Progressive Widening



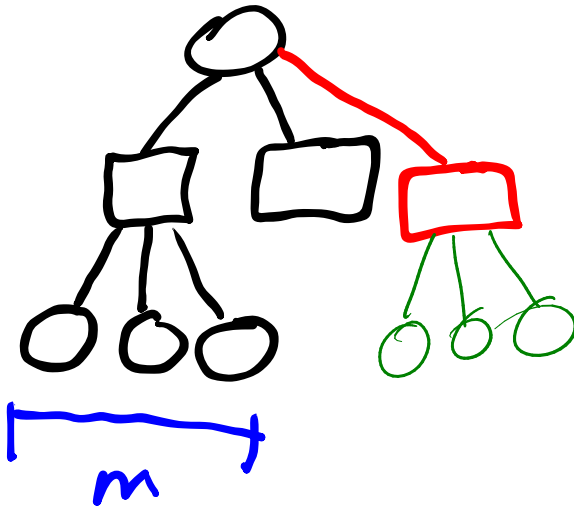
Sparse Tree Search/Progressive Widening



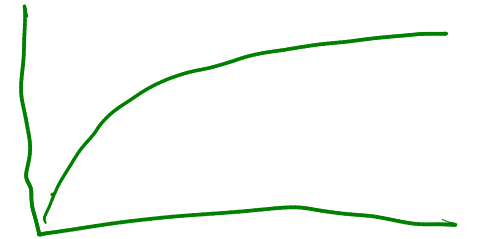
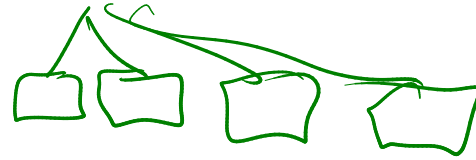
Sparse Tree Search / Progressive Widening



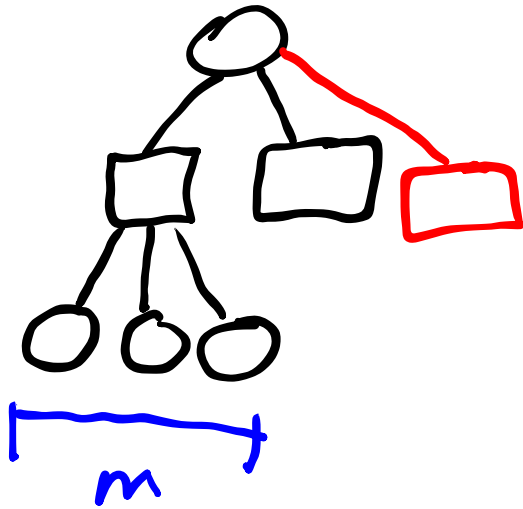
Sparse Tree Search/Progressive Widening



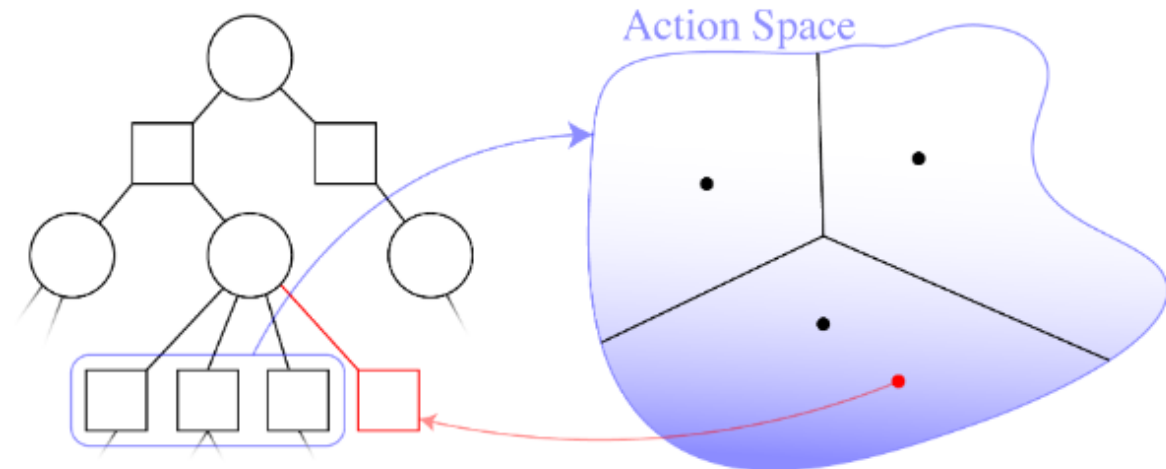
add new branch if $C < kN^\alpha$



Sparse Tree Search/Progressive Widening



add new branch if $C < kN^\alpha$



Online Tree Search Planner

Voronoi Progressive Widening

Guiding Questions

Guiding Questions

- What tools do we have to solve MDPs with continuous S and A ?