# Continuous Space MDPs

# 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 $\mathcal{S}$ and $\mathcal{A}$?

# Current Tool-Belt

# 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!

# Today: Four Tools

# 1. Linear Dynamics, Quadratic Reward

# 2. Value Function Approximation

# 2. Value Function Approximation

$V_\theta(s) = f_\theta(s)$     (e.g. neural network)

# 2. Value Function Approximation

$V_\theta(s) = f_\theta(s)$     (e.g. neural network)

$V_\theta(s) = \theta^\top \beta(s)$     (linear feature)

# 2. Value Function Approximation

$V_\theta(s) = f_\theta(s)$     (e.g. neural network)

$V_\theta(s) = \theta^\top \beta(s)$     (linear feature)

**Fitted Value Iteration**

while not converged

$\quad \theta \leftarrow \theta'$

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

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

# 2. Value Function Approximation

$V_\theta(s) = f_\theta(s)$     (e.g. neural network)

$V_\theta(s) = \theta^\top \beta(s)$     (linear feature)

**Fitted Value Iteration**

while not converged

$\quad \theta \leftarrow \theta'$

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

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

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

# 2. Value Function Approximation

$V_\theta(s) = f_\theta(s)$     (e.g. neural network)

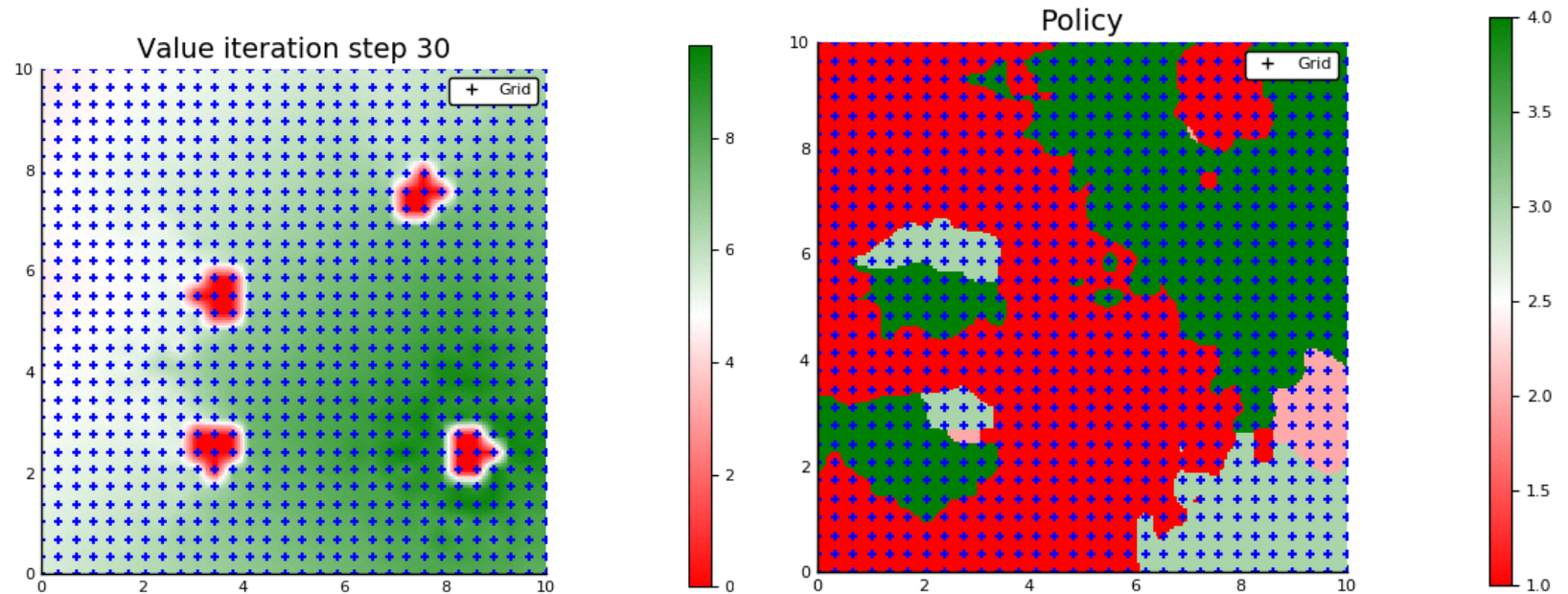$V_\theta(s) = \theta^\top \beta(s)$     (linear feature)

**Fitted Value Iteration**

while not converged

$\quad \theta \leftarrow \theta'$

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

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



Value iteration step 30

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

# 2. Value Function Approximation

$V_\theta(s) = f_\theta(s)$      (e.g. neural network)

$V_\theta(s) = \theta^\top \beta(s)$      (linear feature)

**Fitted Value Iteration**

while not converged

     $\theta \leftarrow \theta'$

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

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



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

# Function Approximation

Weighting of $2^d$ points

Weighting of only $d+1$ points!

# Function Approximation

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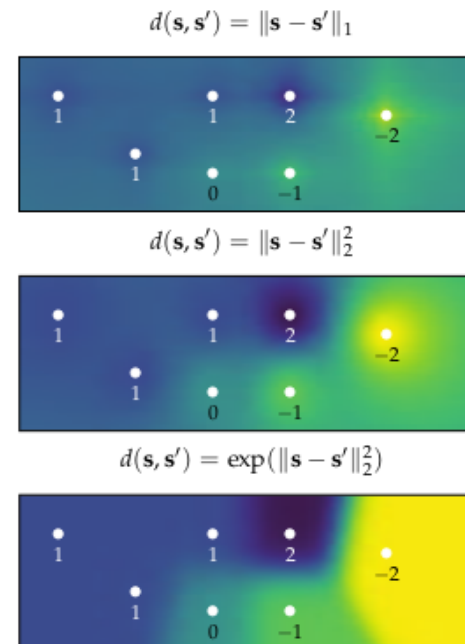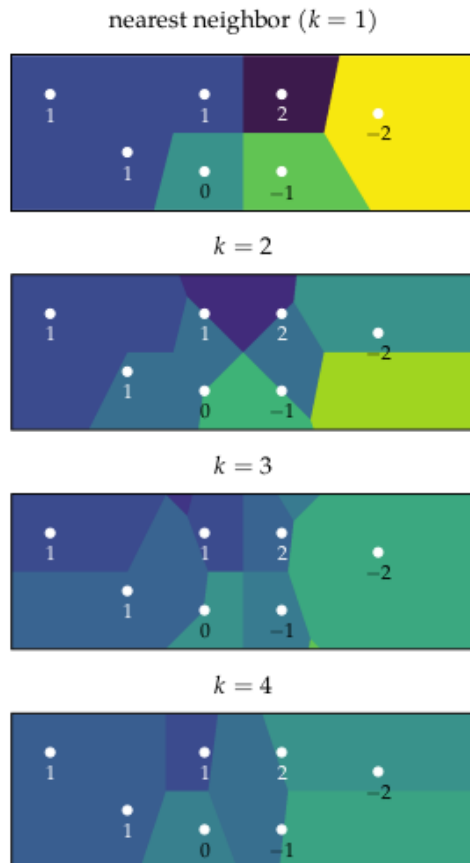- Local: (e.g. simplex interpolation)

Weighting of $2^d$ points

Weighting of only $d+1$ points!

# Function Approximation

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

nearest neighbor ($k = 1$)

$k = 2$

$k = 3$

$k = 4$

$d(\mathbf{s}, \mathbf{s}') = \|\mathbf{s} - \mathbf{s}'\|_1$

$d(\mathbf{s}, \mathbf{s}') = \|\mathbf{s} - \mathbf{s}'\|_2^2$

$d(\mathbf{s}, \mathbf{s}') = \exp(\|\mathbf{s} - \mathbf{s}'\|_2^2)$

Figure 8.9. Two-dimensional linear interpolation over a $3 \times 7$ grid.

Weighting of $2^d$ points

Figure 8.10. Two-dimensional simplex interpolation over a $3 \times 7$ grid.

Weighting of only $d + 1$ points!

# Function Approximation: Mountain Car

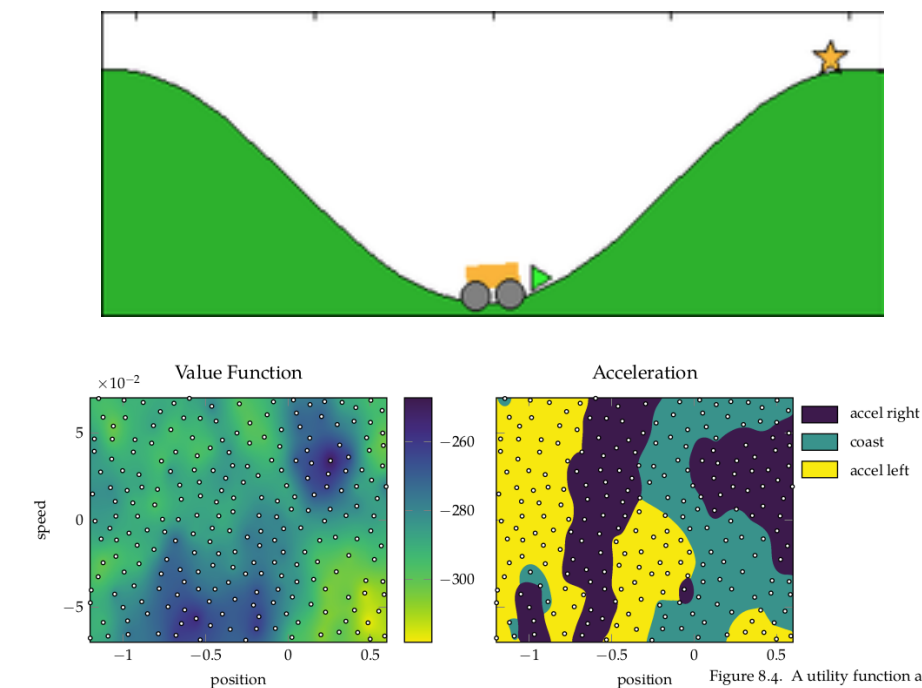# Function Approximation: Mountain Car
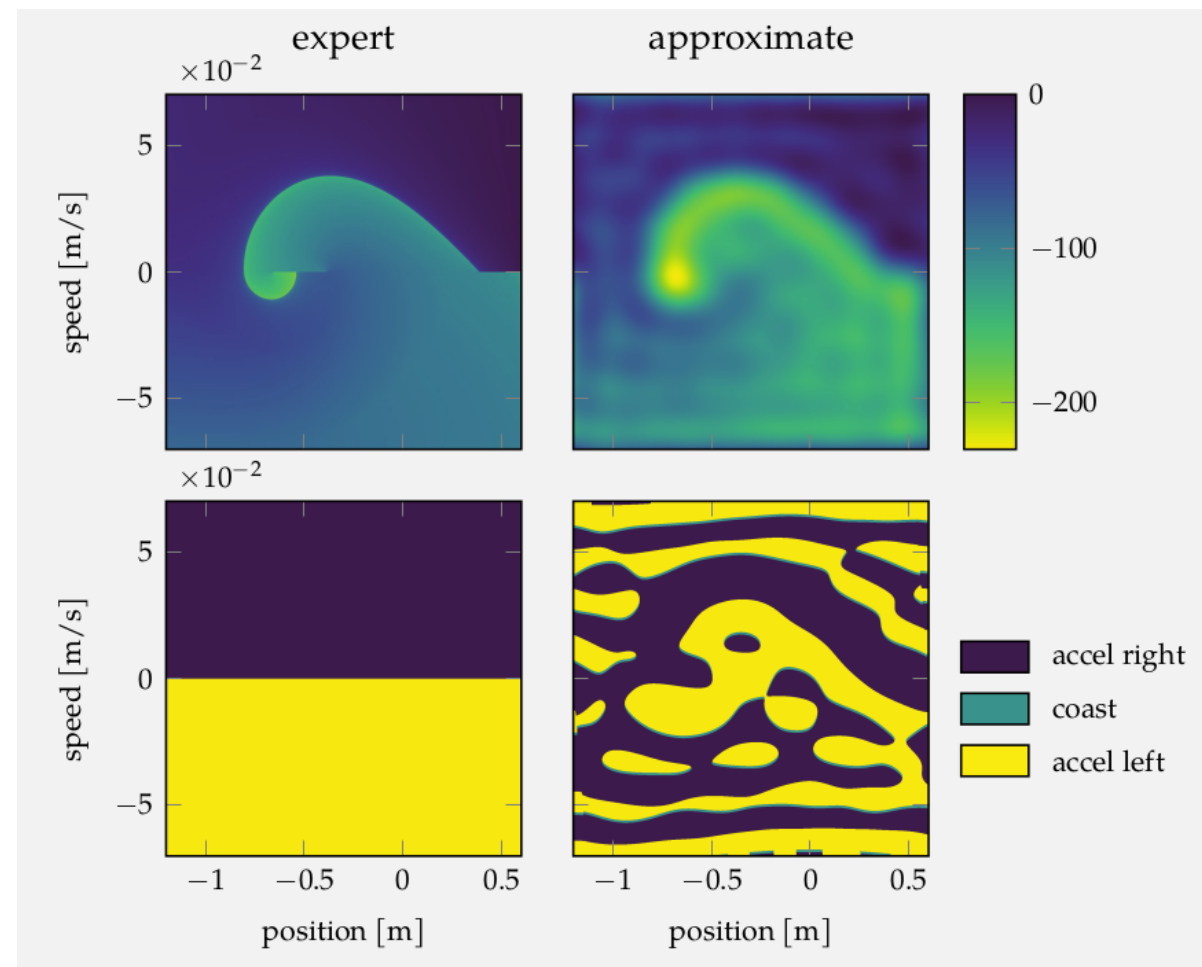
# Function Approximation: Mountain Car



(Fourier, 17 params)

# Function Approximation: Mountain Car
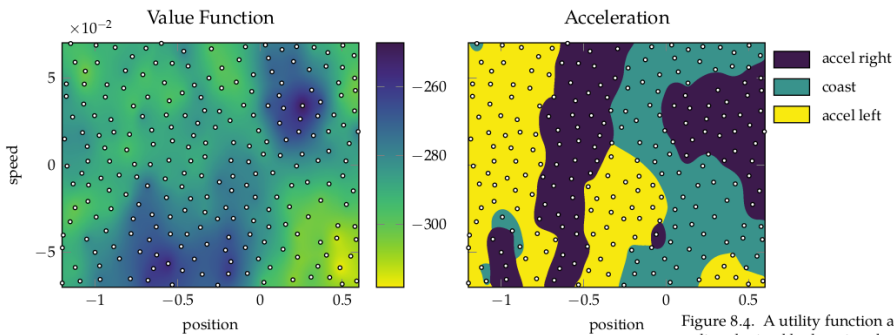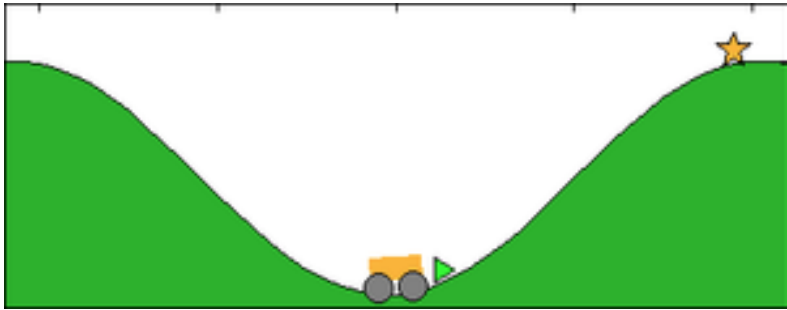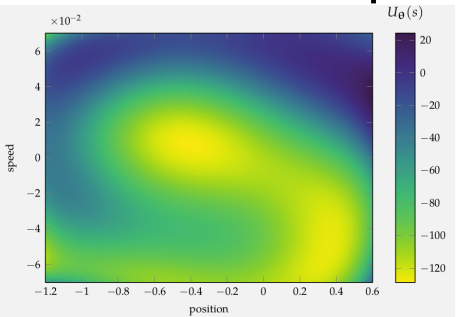


(Kernel, > 100 params)

Figure 8.4. A utility function and policy obtained by learning the action values for a finite set of states (white) in the mountain car problem using the distance function $\|\mathbf{s} - \mathbf{s}'\|_2 + 0.1$.

(Fourier, 17 params)

# Function Approximation: Mountain Car
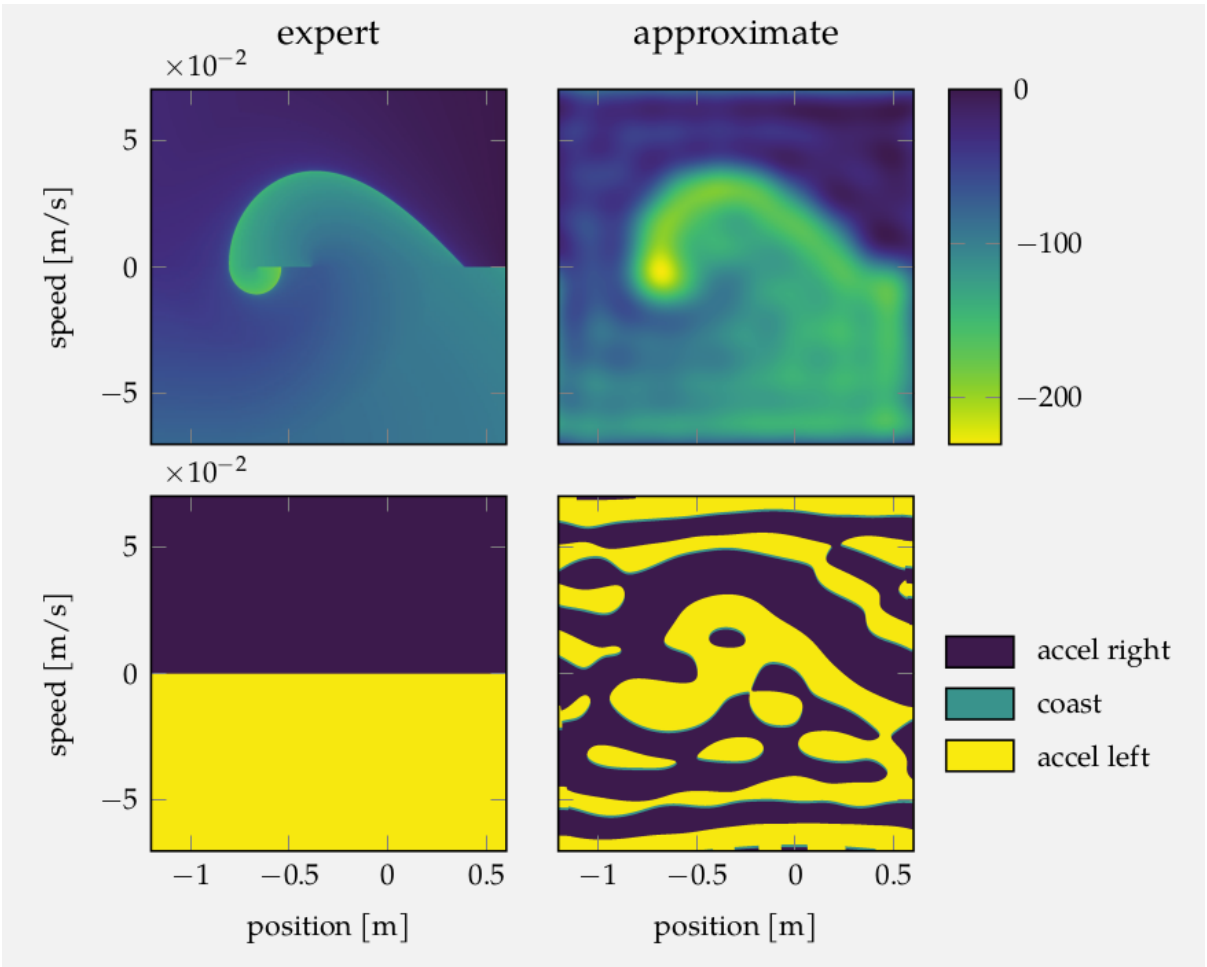


(Kernel, > 100 params)

Figure 8.4. A utility function and policy obtained by learning the action values for a finite set of states (white) in the mountain car problem using the distance function $\|s - s'\|_2 + 0.1$.
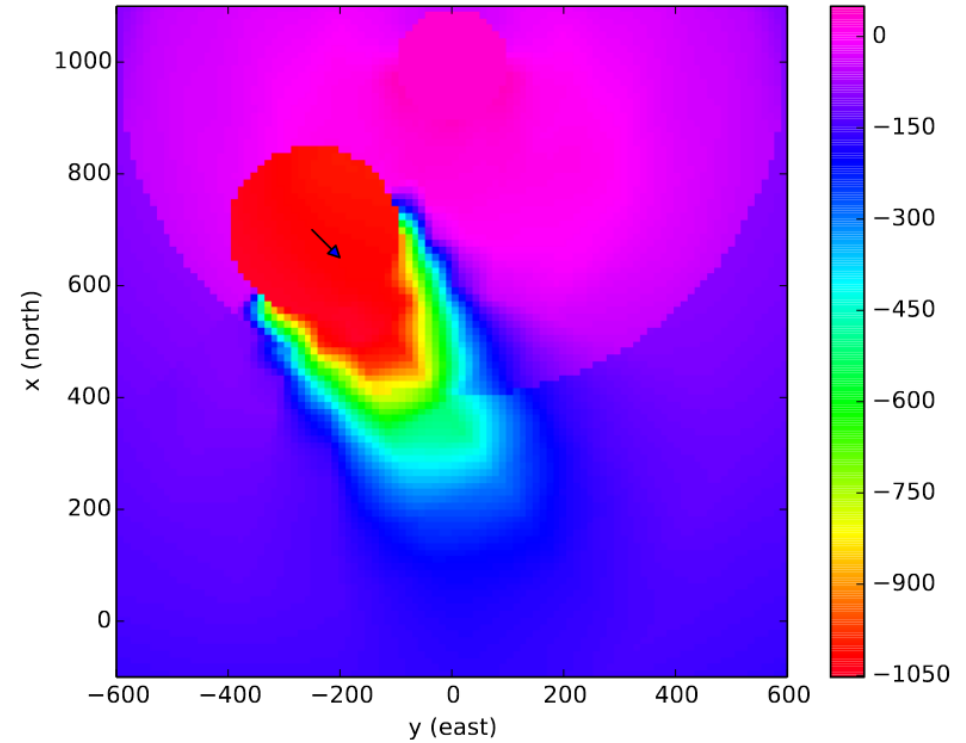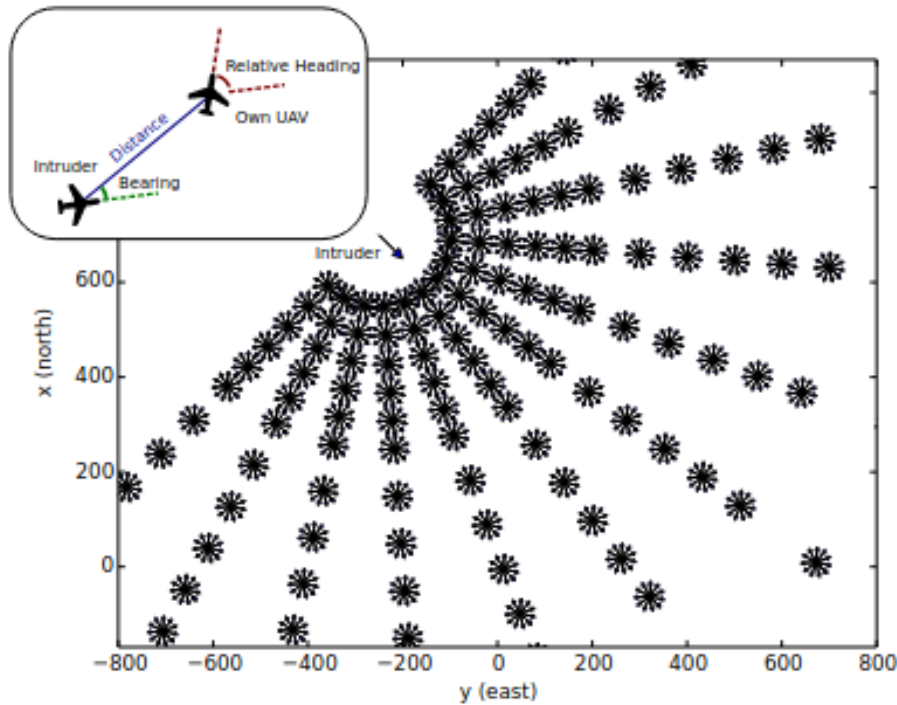
(Polynomial, 28 params)

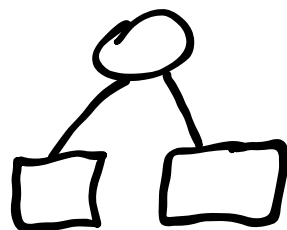(Fourier, 17 params)

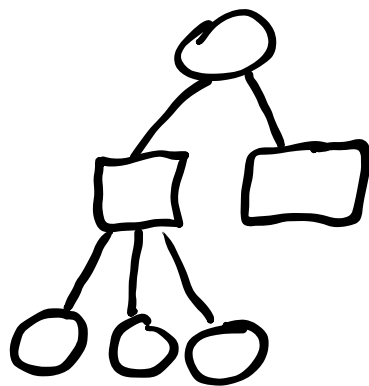# Function Approximation

# Break

What will a Monte Carlo Tree Search tree look like if run on a problem with continuous spaces?
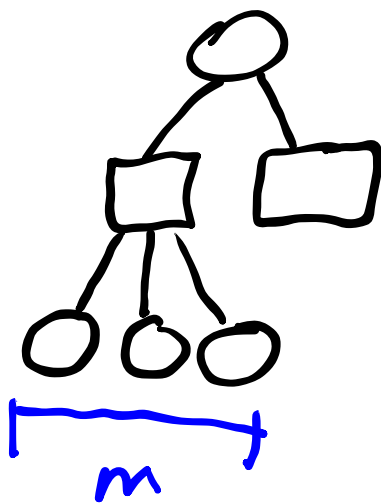
# 3. Sparse Tree Search/Progressive Widening

# 3. Sparse Tree Search/Progressive Widening
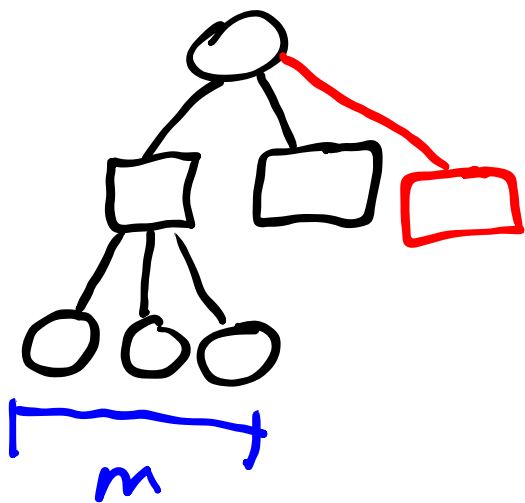
# 3. Sparse Tree Search/Progressive Widening
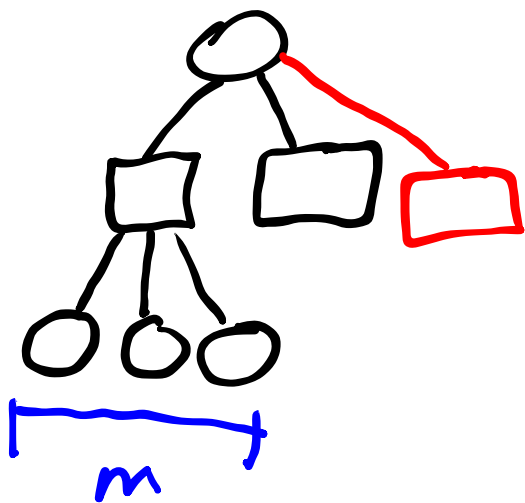
# 3. Sparse Tree Search/Progressive Widening

# 3. Sparse Tree Search/Progressive Widening
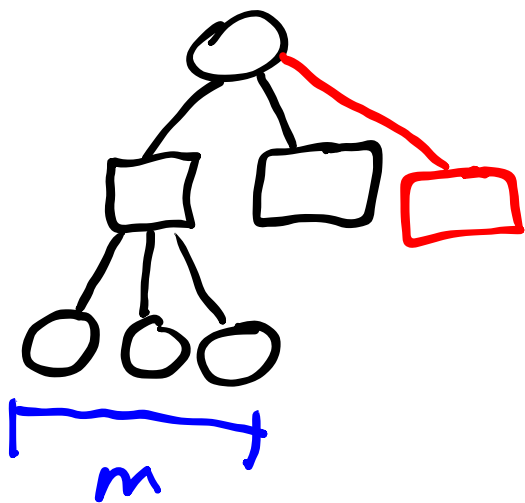


$m$

# 3. Sparse Tree Search/Progressive Widening



add new branch if $C < kN^\alpha$  $(\alpha < 1)$

# 3. Sparse Tree Search/Progressive Widening

add new branch if $C < kN^{\alpha}$ $\quad (\alpha < 1)$



$m$

Online Tree Search Planner     Voronoi Progressive Widening

Action Space

# 4. Model Predictive Control

(Use off-the-shelf optimization software, e.g. Ipopt)

# 4. Model Predictive Control

## (Use off-the-shelf optimization software, e.g. Ipopt)

Certainty-
Equivalent

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

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

# 4. Model Predictive Control

## (Use off-the-shelf optimization software, e.g. Ipopt)

Certainty-
Equivalent

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

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

Open-Loop

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

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

# 4. Model Predictive Control

## (Use off-the-shelf optimization software, e.g. Ipopt)

Certainty-Equivalent

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

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

Open-Loop

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

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

Hindsight Optimization

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

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

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

# Guiding Questions

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