

Reinforcement Learning Based ADP



Optimal Control Problem

Problem

Design a controller u that minimizes the cost

$$J(x, u) = \int_0^{\infty} r(x(\tau), u(x(\tau))) d\tau$$

$$r(x^o, u^o) = Q(x^o) + u^{oT} R u^o$$

Subject to dynamic constraints

$$\dot{x}(t) = f(x(t)) + g(x(t)) u(t)$$

Exact Solution

Optimal value function

$$V^*(x^o) \triangleq \min_{u(\tau) \in U | \tau \in \mathbb{R}_{\geq t}} \int_t^{\infty} r(\phi^u(\tau; t, x^o), u(\tau)) d\tau$$

HJB equation

$$0 = \min_{u^o \in U} (\nabla V^*(x^o) (f(x^o) + g(x^o) u^o) + r(x^o, u^o))$$

Optimal policy

$$u^*(x^o) = -\frac{1}{2} R^{-1} g^T(x^o) (\nabla V^*(x^o))^T$$



Uncertainty

- How to make the best possible decision in the presence of uncertainty, right now
- How to solve the exploration vs. exploitation problem, while also performing system identification
 - Bellman Error extrapolation (simulation of experience) = simultaneous exploration and exploitation
 - Concurrent learning = on-line data-based system identification
 - **Sparse Neural Networks (CDC 2019)**

Expensive

- Curse of dimensionality – large computational cost
 - StaF approximation
 - Sparse NN approximation

More complex problems

- How to include constraints, embed logic-based decision making, intermittency,
 - Formal methods, hybrid/switched systems ADP, scalability
 - **Application: Orbital transfer differential game (SciTech 2020)**



Approximate BE Extrapolation

Uncertainty

Bellman error

$$\delta(x, \hat{W}_c, \hat{W}_a) = r(x, \hat{u}(x, \hat{W}_a)) + \nabla \hat{V}(x, \hat{W}_c) (f(x) + g(x) \hat{u}(x, \hat{W}_a))$$

Parametric approximation $\hat{f}(x, \hat{\theta})$

Approximate Bellman error

$$\hat{\delta}(x, \hat{W}_c, \hat{W}_a, \hat{\theta}) = r(x, \hat{u}(x, \hat{W}_a)) + \nabla \hat{V}(x, \hat{W}_c) (Y(x) \hat{\theta} + g(x) \hat{u}(x, \hat{W}_a))$$

If $\hat{\theta}(t) \rightarrow B_r(\theta)$ exponentially as $t \rightarrow \infty$,
 $\hat{\delta}_t(t) \rightarrow B_{r1}(\delta_t(t))$ exponentially as $t \rightarrow \infty$.

$$\hat{\delta}_t(t) = \hat{\delta}(x(t), \hat{W}_c(t), \hat{W}_a(t), \hat{\theta}(t)) \quad \hat{\delta}_{ti}(t) = \hat{\delta}(x_i, \hat{W}_c(t), \hat{W}_a(t), \hat{\theta}(t))$$

Use $\hat{\delta}_t(t)$ and $\hat{\delta}_{ti}(t)$ for learning

Update law

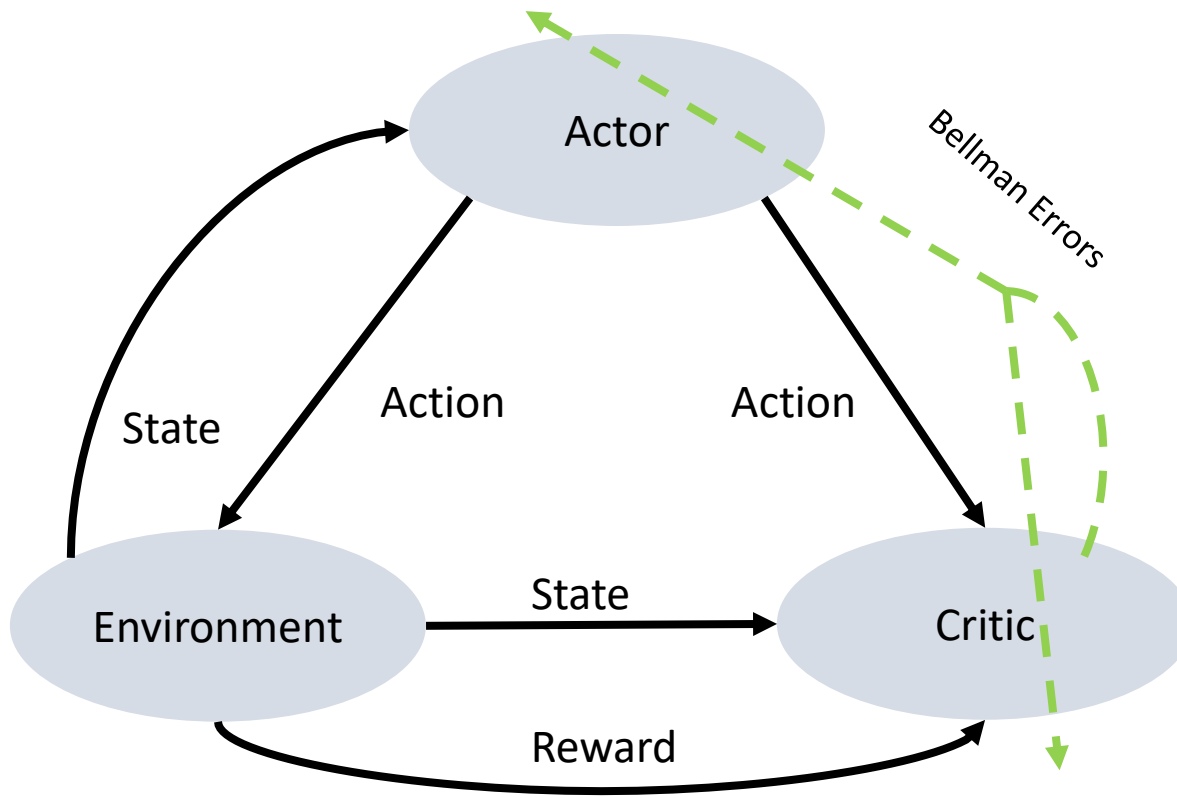
$$\dot{\hat{W}}_c = -\frac{\eta_{c1}\Gamma}{1+\nu\omega^T\Gamma\omega}\omega\hat{\delta}_t - \eta_{c2}\Gamma\sum_{i=1}^N\frac{1}{1+\nu\omega_i^T\Gamma\omega_i}\omega_i\hat{\delta}_{ti}$$

Weight estimation error dynamics

$$\dot{\tilde{W}}_c = -\Gamma\left(\eta_{c1}\frac{\omega\omega^T}{\rho} + \eta_{c2}\sum_{i=1}^N\frac{\omega_i\omega_i^T}{\rho_i}\right)\tilde{W} + \Delta$$

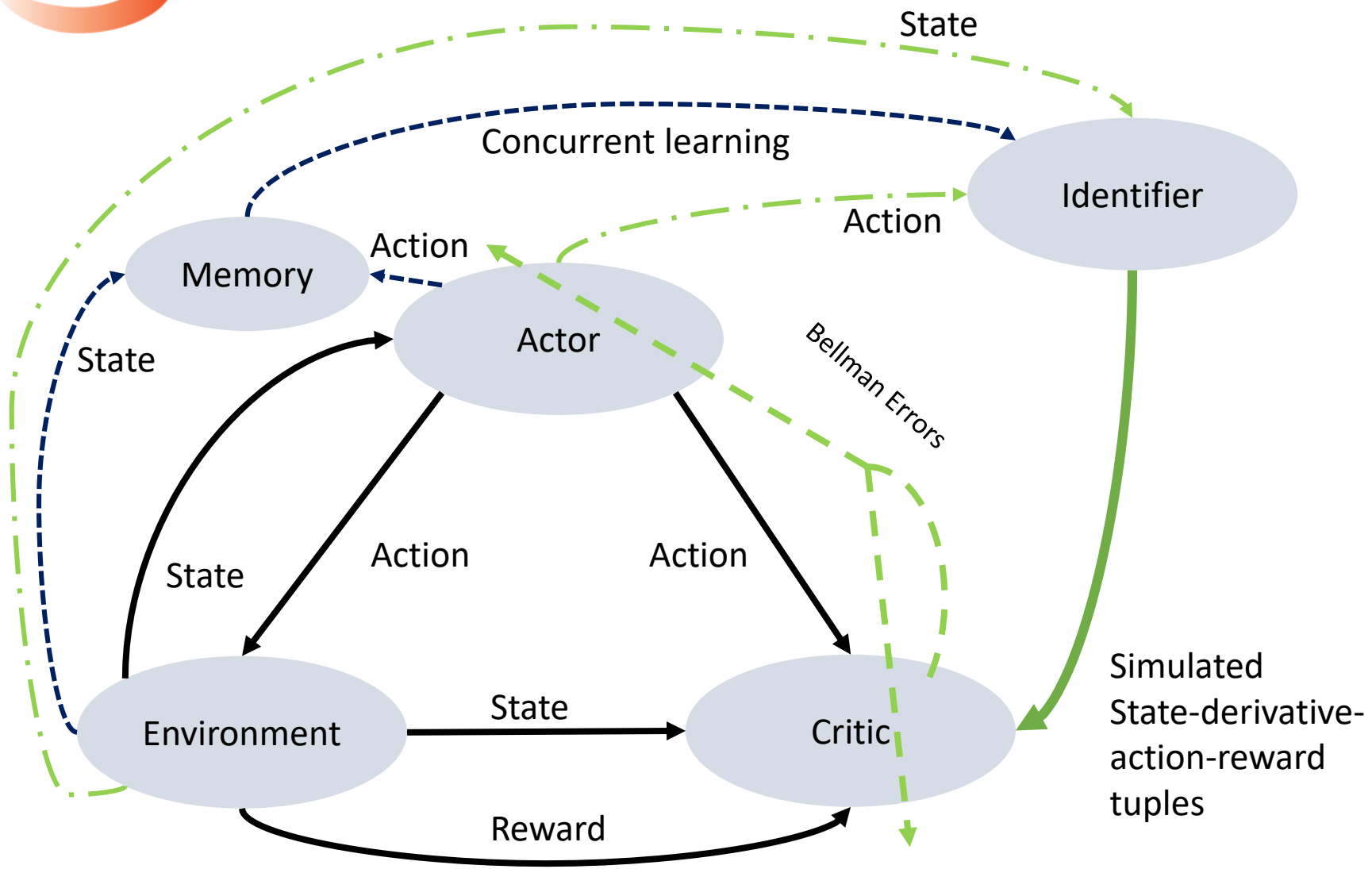


Typical RL-based ADP Approach



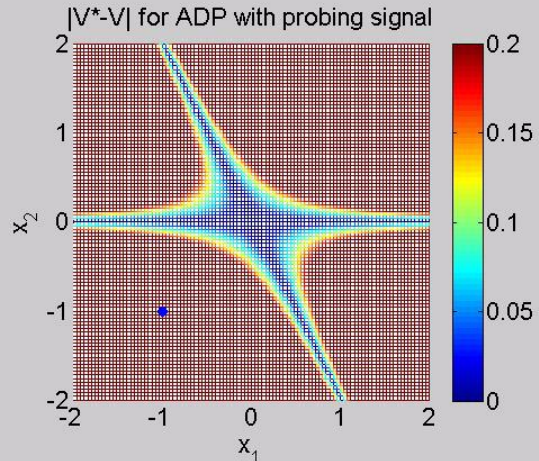
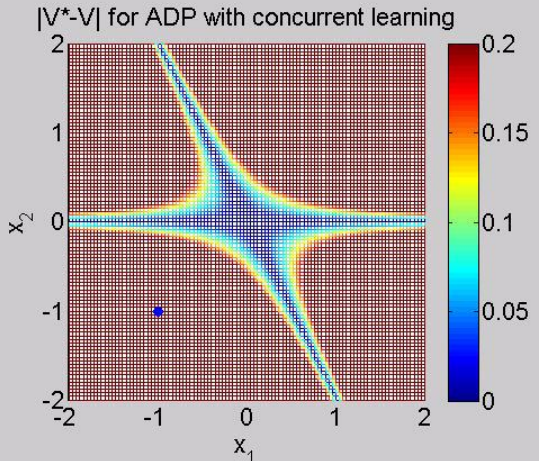


Simulated Experience

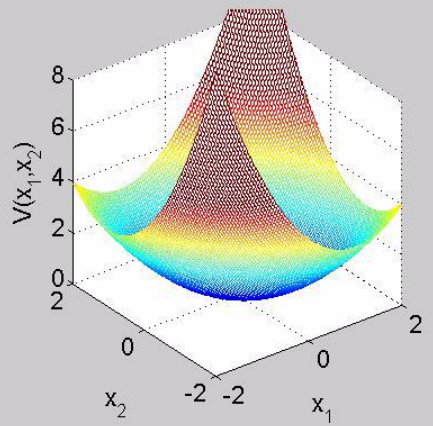




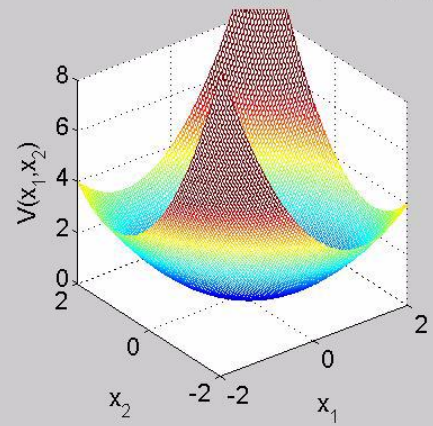
Simulation: Known Optimal Solution



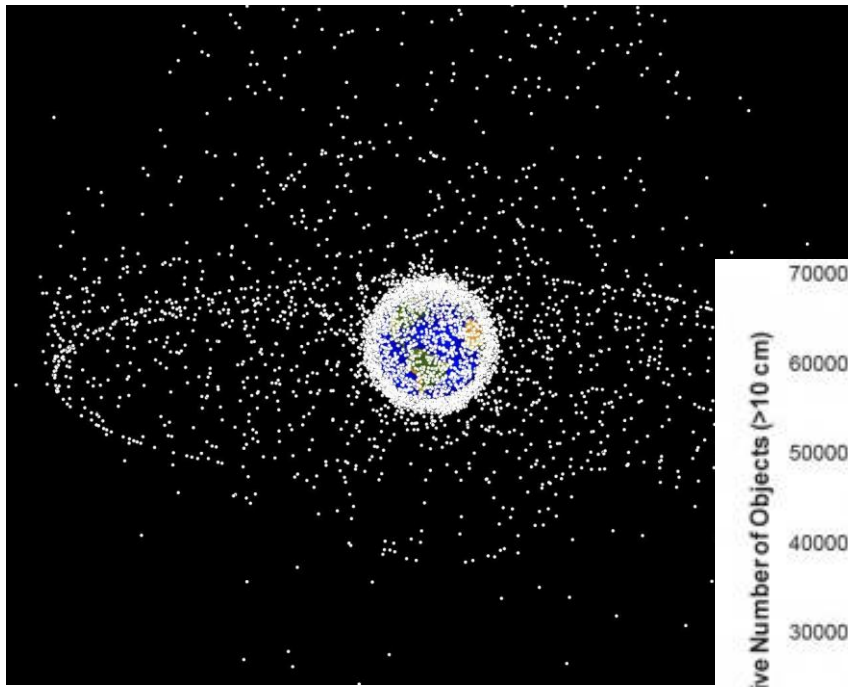
Value function for ADP with concurrent learning



Value function for ADP with probing signal

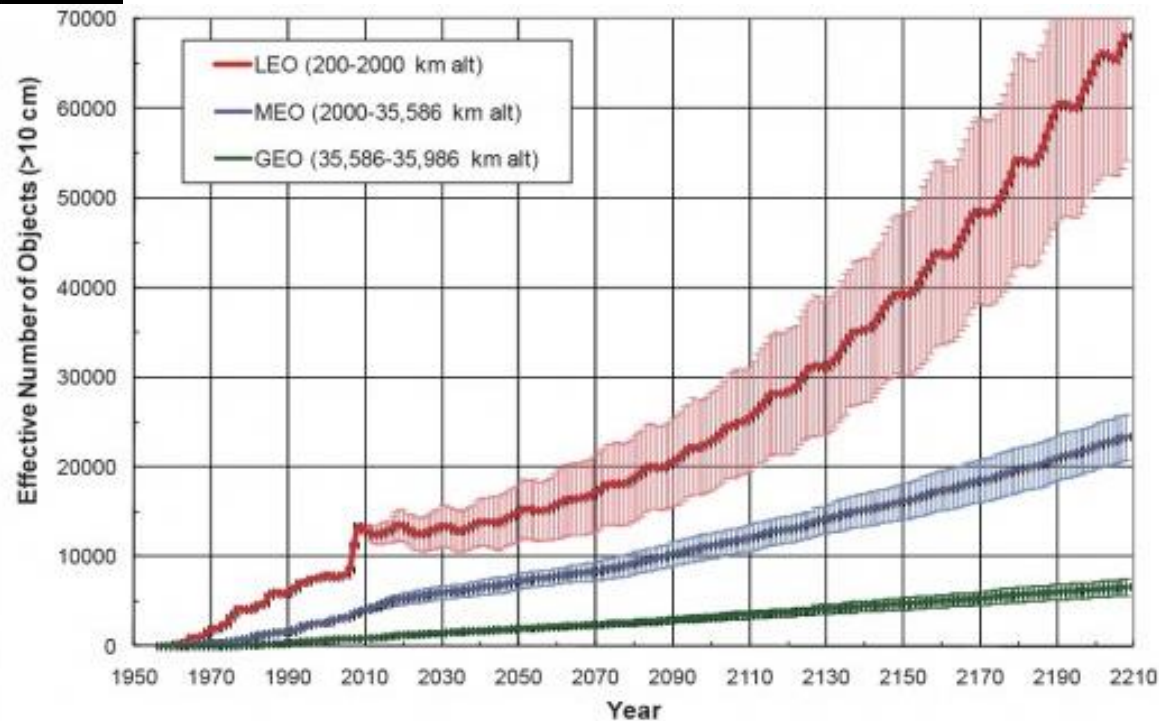


APPROXIMATE OPTIMAL ORBIT TRANSFER OF NON-COOPERATIVE DEBRIS



[NASA]

Approximate Optimal Orbit Transfer of Non-cooperative Debris



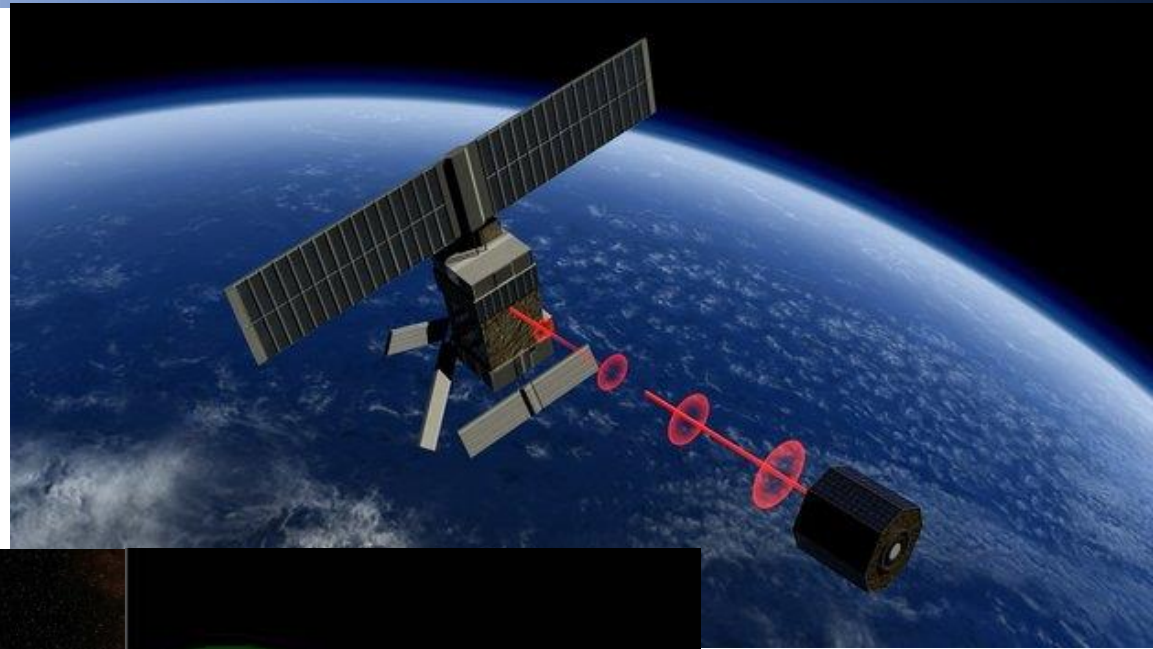
[NASA]

Debris Removal

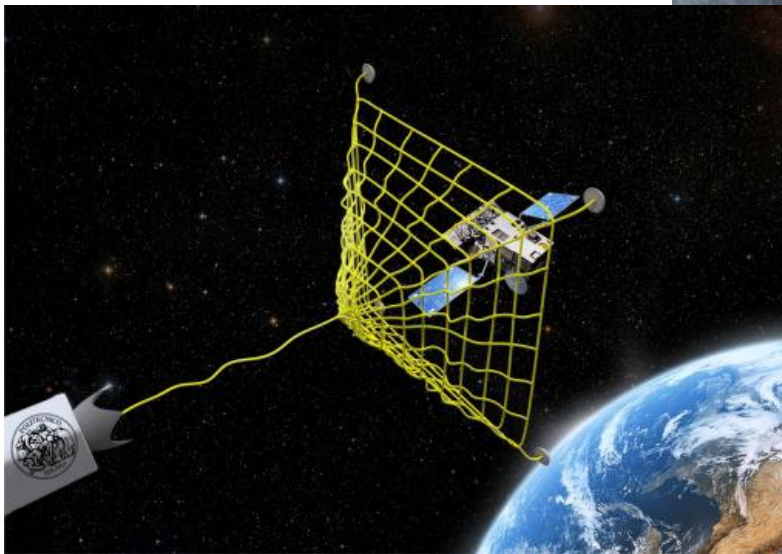


Removal Methods

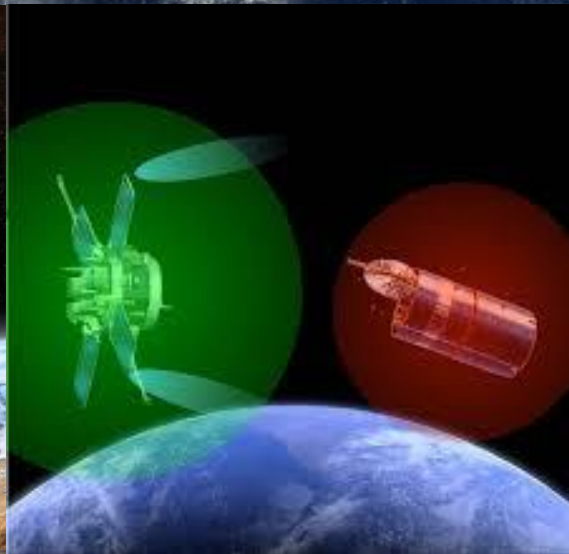
- Laser-Based
- Ion-Beam
- Tether-Based
- Sail-Based
- Satellite-Based
- Unconventional



[ESA]



[Benvenuto, Salvi, Lavagna]



[Hughes, Schaub]



Problem Formulation

Agent Dynamics

$$\dot{\eta}(t) = h(\eta(t)) + g(z(t), \eta(t), t)$$

$$\dot{z}(t) = d(t) \triangleq f(z(t), \eta(t), t)$$

Error System

Servicing Satellite (SS), $\eta(t)$, intercept Disabled Satellite (DS), $z(t)$, at a distance, r_d

$$e_1 \triangleq z(t) - (\eta(t) - r_d)$$

Servicing Satellite (SS), $\eta(t)$, regulate the Disabled Satellite (DS), $z(t)$, to z_g

$$e_2 \triangleq (\eta(t) - r_d) - z_g - k_1(z(t) - z_g)$$

Compact Dynamics

$$\dot{x}(t) = F(x(t)) + G(x(t))u(t) + Kd(t)$$



Problem Formulation

Control Objective (2 player zero sum game)

Design a controller, u , which minimizes a cost function:

$$J(x, u, d, t_0) = \int_{t_0}^{\infty} (x(\tau)^T Q x(\tau) + u(\tau)^T R u(\tau) - \gamma^2 d^T(\tau) d(\tau)) d\tau$$

Cost-to-Go

Optimal value function:

$$V^*(x(t)) = \min_{u(\tau)} \max_{d(\tau)} \int_t^{\infty} (x(\tau)^T Q x(\tau) + u(\tau)^T R u(\tau) - \gamma^2 d^T(\tau) d(\tau)) d\tau$$

HJI Equation

Design a controller, u , which minimizes a cost function:

$$0 = x(t)^T Q x(t) + u(t)^T R u(t) - \gamma^2 d^T(t) d(t) \\ + \nabla V^*(x(t)) \left(F(x(t)) + G(x(t)) u^*(x(t)) + K d^*(x(t)) \right)$$



Problem Formulation

Universal Function Approximation Property

$$V^*(x) = W(x)^T \sigma(c(x)) + \epsilon_W(x)$$

$$u^*(x) = -\frac{R^{-1}G(x)^T}{2} \left(\nabla \sigma(c(x))^T W(x) + \epsilon_W(x)^T \right)$$

$$d^*(x) = \frac{K^T}{2\gamma^2} \left(\nabla \sigma(c(x))^T W(x) + \epsilon_W(x)^T \right)$$

Stone-Weierstrauss Approximation Theorem

$$\hat{V}(x, \hat{W}_c) = \hat{W}_c^T \sigma(c(x))$$

$$\hat{u}(x, \hat{W}_a) = -\frac{R^{-1}G(x)^T}{2} \nabla \sigma(c(x))^T \hat{W}_a$$

$$\hat{d}(x, \hat{W}_d) = \frac{K^T}{2\gamma^2} \nabla \sigma(c(x))^T \hat{W}_d$$



Problem Formulation

Bellman Error

$$\delta_t(x, \hat{W}_c, \hat{W}_a, \hat{W}_d) = x(t)^T Q x(t) + \hat{u}(x, \hat{W}_a)^T R \hat{u}(x, \hat{W}_a) - \gamma^2 \hat{d}(x, \hat{W}_d)^T \hat{d}(x, \hat{W}_d) + \nabla \hat{V}(x, \hat{W}_c) \left(F(x) + G(x) \hat{u}(x, \hat{W}_a) + K \hat{d}(x, \hat{W}_d) \right)$$

Update Laws

$$\dot{\hat{W}}_c(t) = -\Gamma_c(t) \frac{k_{c1}}{N+1} \frac{\omega(t)}{\rho(t)^2} \delta_t(t) - \Gamma_c(t) \frac{k_{c2}}{N+1} \sum_{i=1}^N \frac{\omega_i(t)}{\rho_i(t)^2} \delta_{ti}(t)$$

$$\dot{\hat{W}}_a(t) = \text{proj} \left(-K_a k_{a1} \left(\hat{W}_a(t) - \hat{W}_c(t) \right) \right)$$

$$\dot{\hat{W}}_d(t) = \text{proj} \left(-K_d k_{d1} \left(\hat{W}_d(t) - \hat{W}_c(t) \right) \right)$$

$$\dot{\Gamma}_c(t) = \beta_c \Gamma_c(t) - \Gamma_c(t) \frac{k_{c1}}{N+1} \frac{\omega(t) \omega^T(t)}{\rho(t)^2} \Gamma_c(t) - \Gamma_c(t) \frac{k_{c2}}{N+1} \sum_{i=1}^N \frac{\omega_i(t) \omega_i^T(t)}{\rho_i(t)^2} \Gamma_c(t)$$



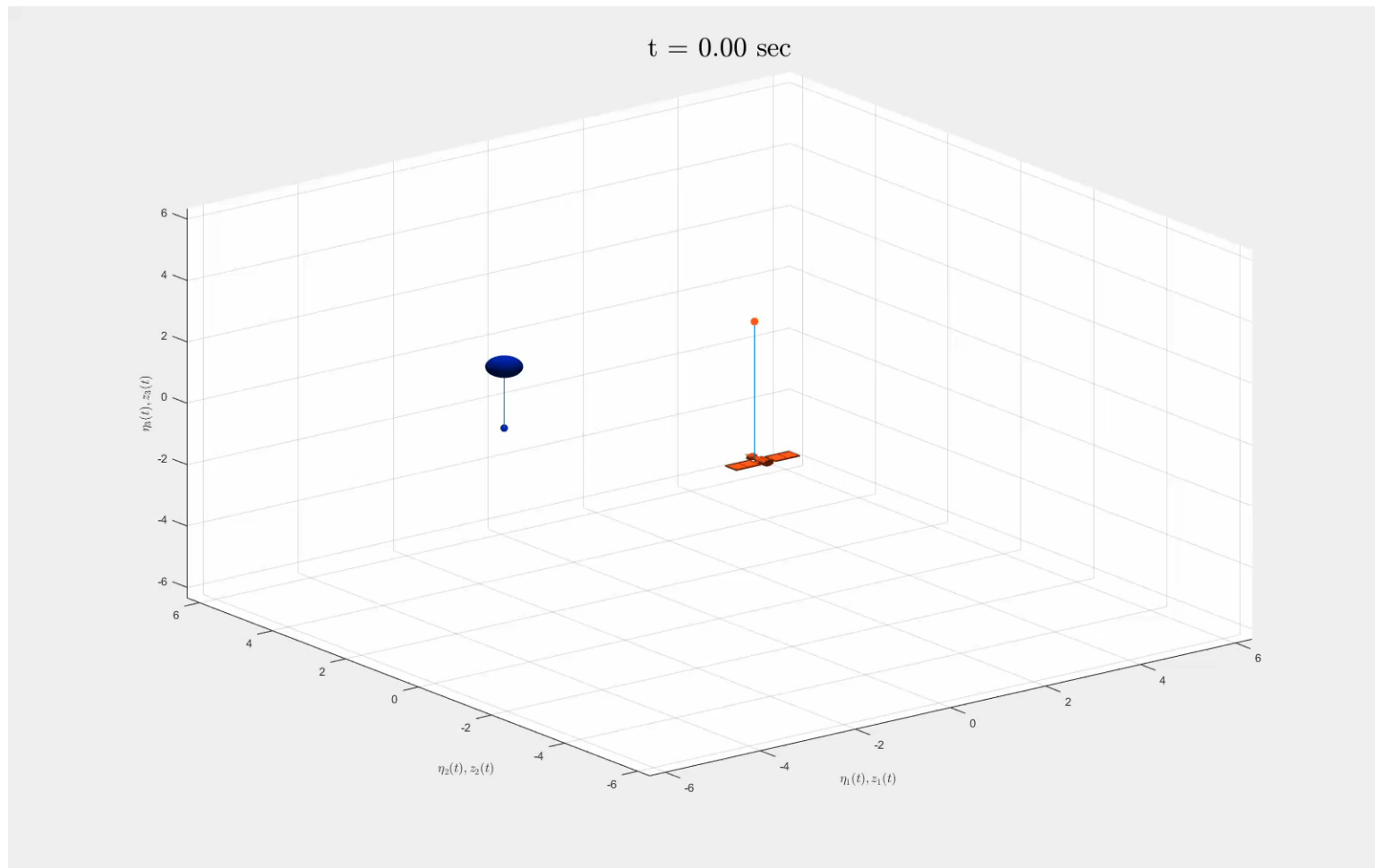
Lyapunov Stability Analysis

$$V_L(Z_L, t) = V^*(x) + \frac{1}{2} \tilde{W}_c^T \Gamma_c^{-1}(t) \tilde{W}_c + \frac{1}{2} \tilde{W}_a^T K_a^{-1} \tilde{W}_a + \frac{1}{2} \tilde{W}_d^T K_d^{-1} \tilde{W}_d$$

$$\begin{aligned} \dot{V}_L(Z_L, t) = & \nabla V^*(x, \hat{W}_c)(F + G\hat{u} + K\hat{d}) + \tilde{W}_c^T \Gamma_c^{-1} \dot{\tilde{W}}_c + \tilde{W}_a^T K_a^{-1} \dot{\tilde{W}}_a \\ & + \tilde{W}_d^T K_d^{-1} \dot{\tilde{W}}_d - \frac{1}{2} \tilde{W}_c^T (\Gamma_c^{-1} \dot{\Gamma}_c \Gamma_c^{-1}) \end{aligned}$$

- System state (x) , weight estimation errors $(\tilde{W}_c, \tilde{W}_a, \tilde{W}_d)$, and control policy $u(t)$ are **Uniformly Ultimately Bounded**
- By using a State Following (StaF) neural network architecture a local result is achieved

Simulation Results

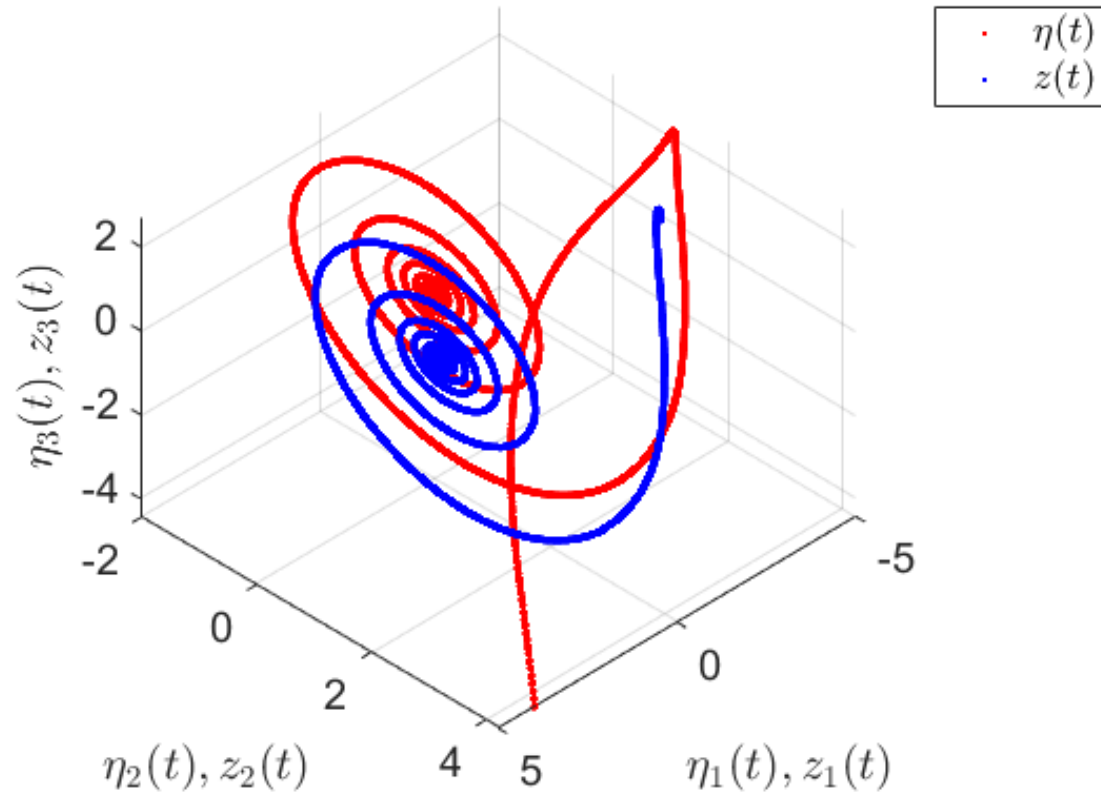


Simulation Results



Future Work

- Consider Orbital Mechanics
- Higher Fidelity Model of Electrostatic Dynamics



Reinforcement Learning With Sparse Bellman Error Extrapolation For Infinite-horizon Approximate Optimal Regulation



Needs?

- Computational Efficiency
- Fewer active NN basis functions
- BE extrapolation in operational subset
- Global Approximation

State Following (StaF) Approximation

- Local Value Function Approximation
- Explores Around the Current State
- Reduced Number of Basis Functions
- Trades Global Optimality for Computational Efficiency



StaF and MBRL Approximation

- Combination of Techniques
- StaF Outside of Boxed Region
- MBRL Inside of Boxed Region

Sparse NN

- Decrease active neurons
- Switch between smaller MBRL segments
- StaF/MBRL Method



BE Extrapolation

Instantaneous BE: Residual from HJB

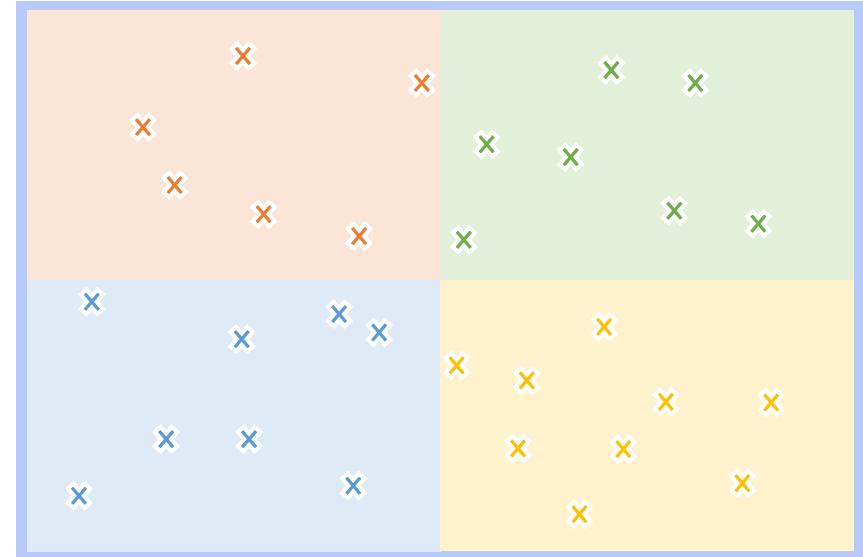
$$\hat{\delta}_i(x, t) \triangleq \hat{\delta} \left(x_i, \hat{W}_c(t), \hat{W}_a(t) \right)$$

Weight Update Laws using MBRL

$$\begin{aligned} \dot{\hat{W}}_c(t) &= -\eta_{c1} \Gamma \frac{\omega(t)}{\rho(t)} \hat{\delta} + \eta_{c2} \frac{1}{N_j} \sum_{i=1}^{N_j} \frac{\omega_i(t)}{\rho_i(t)} \hat{\delta}_i(t) && \text{Actual system} && \text{Extrapolated system} \\ \dot{\Gamma}(t) &= \left(\lambda \Gamma(t) - \frac{\eta_{c1} \Gamma(t) \omega(t) \omega(t)^T \Gamma(t)}{\rho(t)} - \Gamma(t) \eta_{c2} \left(\frac{1}{N_j} \sum_{i=1}^{N_j} \frac{\omega_i(t) \omega_i^T(t)}{\rho_i(t)} \hat{\delta}_i(t) \right) \Gamma(t) \right) \mathbf{1}_{\{\underline{\Gamma} \leq \|\Gamma\| \leq \bar{\Gamma}\}} \\ \dot{\hat{W}}_a(t) &= -\eta_{c1} \left(\hat{W}_a(t) - \hat{W}_c(t) \right) - \eta_{a2} \hat{W}_a(t) + \frac{\eta_{c1} G_\sigma^T(t) \hat{W}_a(t) \omega(t)^T}{4\rho(t)} \hat{W}_c(t) \\ &\quad + \left(\frac{\eta_{c2}}{4N_j} \sum_{i=1}^{N_j} \frac{G_{i\sigma}^T \hat{W}_a(t) \omega_i(t)}{\rho_i(t)} \hat{\delta}_i(t) \right) \hat{W}_c(t) \end{aligned}$$

Segmentation

- Separate Operating Domain
- BE Extrapolation Contained to Segment
- Smaller History Stack (one for each segment)
- Switches Depending on Region
- Less Active History
- A further step towards including memory (cognition)
- Lyapunov analysis shows SG-UUB



Update Laws

$$\dot{\hat{W}}_c(t) = -\eta_{c1} \Gamma \frac{\omega(t)}{\rho(t)} \hat{\delta}(t) - \eta_{c2} \sum_c^j(t)$$

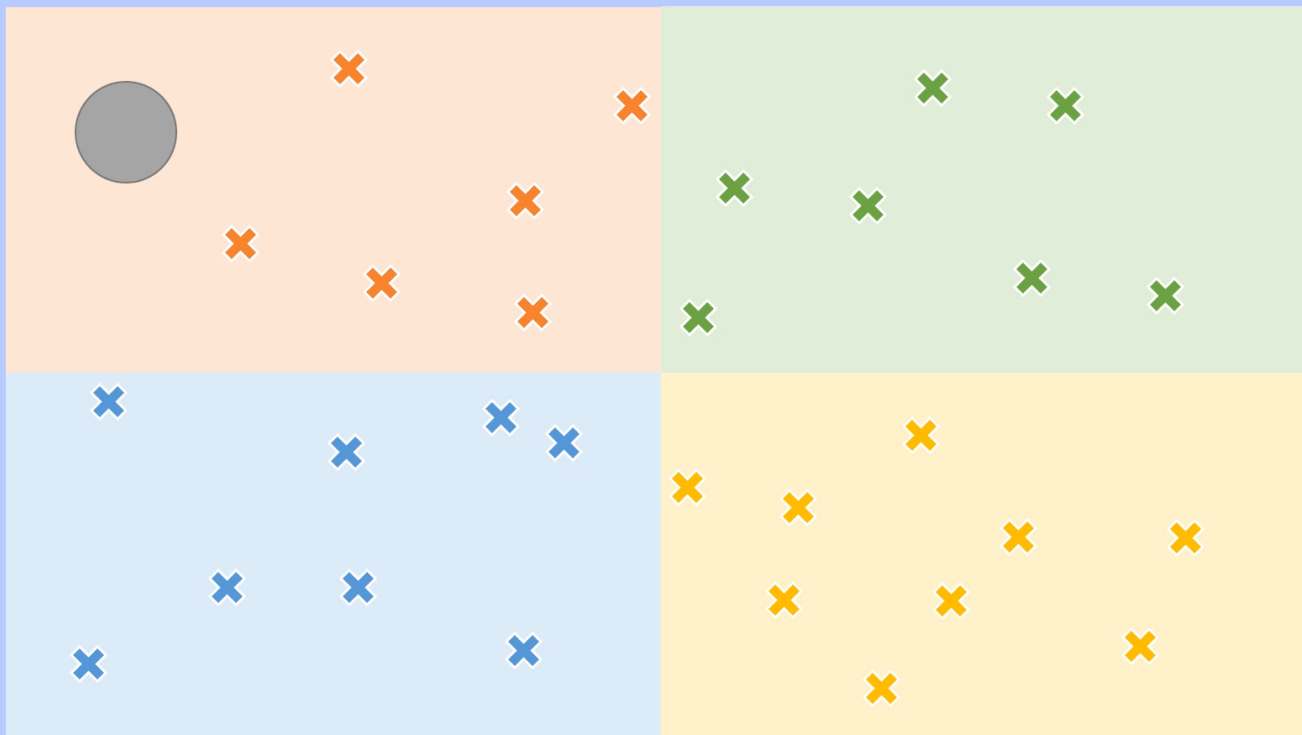
Switching
History Stack
Term

$$\dot{\Gamma}(t) = \left(\lambda \Gamma(t) - \eta_{c1} \frac{\Gamma(t) \omega(t) \omega(t)^T \Gamma(t)}{\rho(t)} - \Gamma(t) \eta_{c2} \left(\sum_{\Gamma}^j(t) \right) \Gamma(t) \right) \mathbf{1}_{\{\underline{\Gamma} \leq \|\Gamma\| \leq \bar{\Gamma}\}}$$

$$\dot{\hat{W}}_a(t) = -\eta_{a1} (\hat{W}_a(t) - \hat{W}_c(t)) - \eta_{a2} \hat{W}_a(t) + \frac{\eta_{c1} G_{\sigma}(t)^T \hat{W}_a(t) \omega(t)^T}{4\rho(t)} \hat{W}_c(t) + \left(\eta_{c2} \sum_a^j(t) \right) \hat{W}_c(t)$$



Switching Segments



$$\dot{\hat{W}}_c(t) = -\eta_{c1} \Gamma \frac{\omega(t)}{\rho(t)} \hat{\delta} + \eta_{c2} \sum_c^j(t)$$

Simulation Results

- Dynamics

- $\dot{x} = f(x) + g(x)u(t)$

- $f(x) = \begin{bmatrix} -x_1 + x_2 \\ -\frac{1}{2}x_1 - \frac{1}{2}x_2(1 - (\cos(2x_1) + 2)^2) \end{bmatrix}$

- $g(x) = \begin{bmatrix} 0 \\ \cos(2x_1) + 2 \end{bmatrix}$

- Initial Conditions

- $x(0) = [-10, 10]^T$

- $\widehat{W}_a = \widehat{W}_c = \left[\frac{1}{2}, \frac{1}{2}, \frac{1}{2}\right]^T$

- Define Segmentation

- $\Omega_1 \triangleq \{x \in \mathbb{R}^2: 3.5 \leq |x_1|, |x_2| \leq 10\}$
 - Outer Segment

- $\Omega_2 \triangleq \{x \in \mathbb{R}^2: |x_1|, |x_2| < 3.5\}$
 - Inner Segment

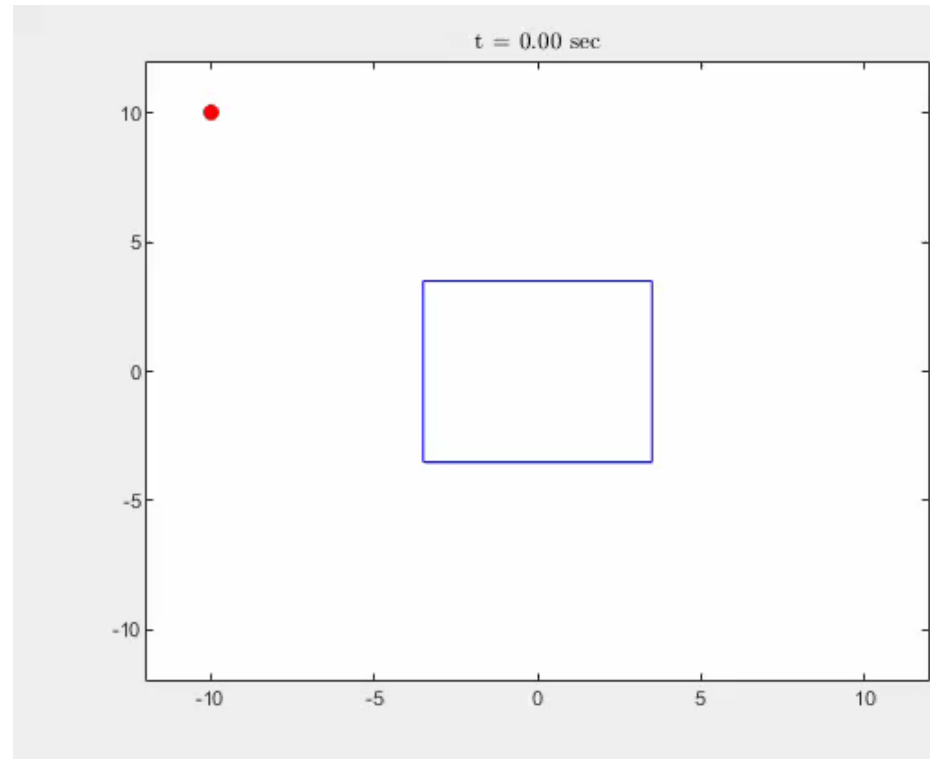
- Define BE Bases

- For $x_i \in \Omega_1$ (Outer Segment)

- $\sigma_i(x_i) = [x_{1,i}^2 \quad x_{1,i}x_{2,i} \quad x_{2,i}^2]^T$

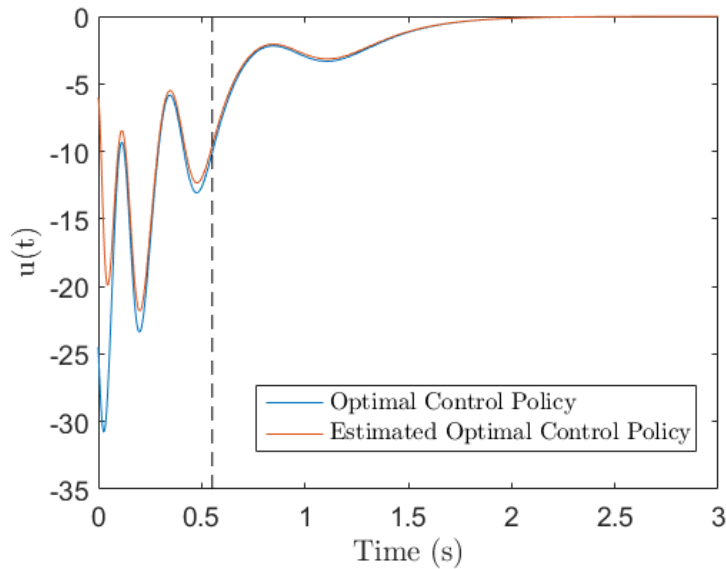
- For $x_i \in \Omega_2$ (Inner Segment)

- $\sigma_i(x_i) = [x_{1,i}^2 \quad 0 \quad x_{2,i}^2]^T$

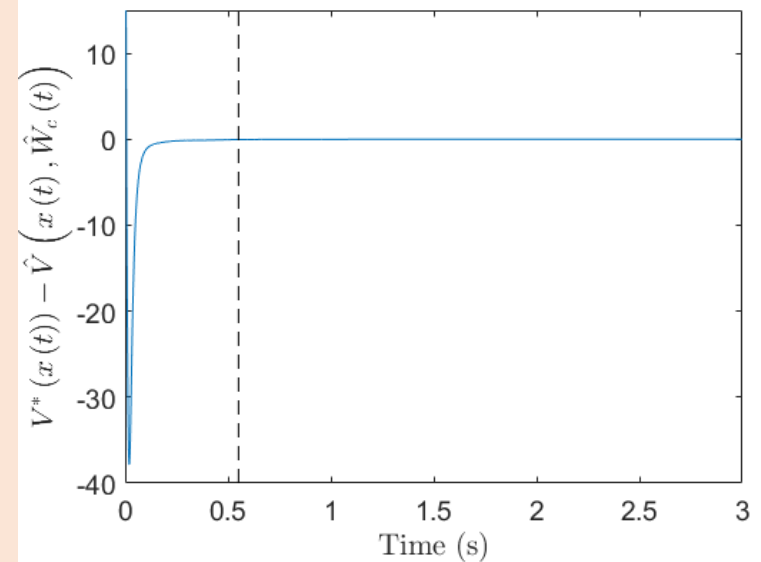


Simulation Results

Control Policy



Value Function





Benefits of Sparsity and Segmentation

- Decrease Active Memory Used
- Decrease Number of Computations in NN
- Establish Sparse, Switching Framework

Future Work (on-going collaboration with Scott Nivison (RW))

- Quantify Sparsity Benefit
 - Number of Computations
 - Compare to Existing ADP Methods
 - Computation Time
- Utilize Barrier Functions
 - Guarantee System Safety
- Application to Multi-Agent Systems