

Problem Setup

The system considered is described as a stochastic differential equation (SDE)

$$dx(t) = f(x(t), t)dt + G(x(t), t)u(t)dt + \Sigma(x(t), t)dw(t).$$

The objective is to find a control sequence that minimizes the following cost

$$J^u(x, t) = \mathbb{E} \left[g(x(T)) + \int_t^T \left(q(x(s)) + \frac{1}{2} u(s)^T R u(s) \right) ds \mid x(t) = x \right]$$

with state constraints

$$c_{\min} \leq c_s(x) \leq c_{\max}$$

and control saturation

$$u \in \mathcal{U} = \{u \mid |u_i| \leq U_{i, \max}\}.$$

State & Control Constraint

The control is saturated as

$$u^*(x, t) = U_{\max} * \text{sig}(-R^{-1}G^T(t, x)V_x).$$

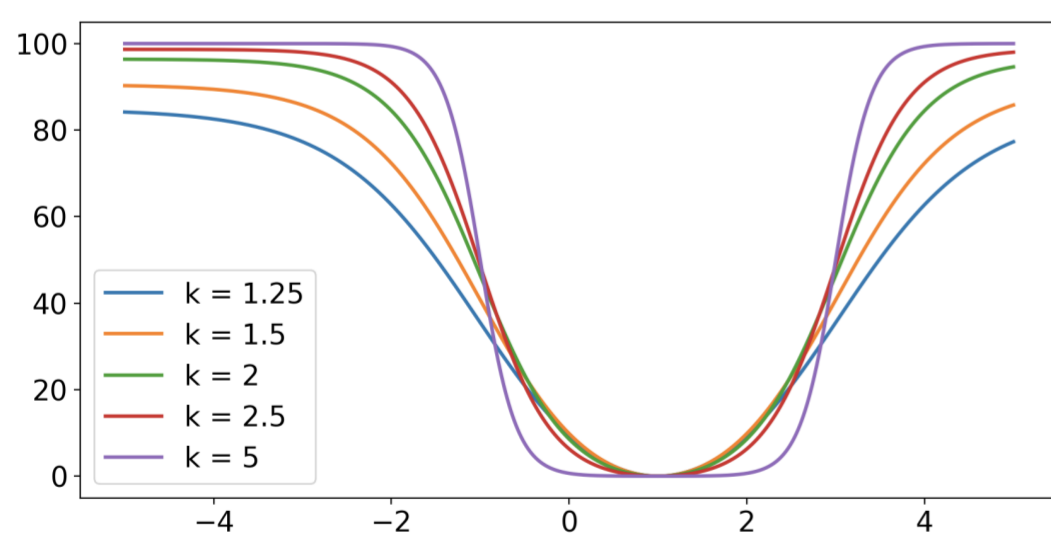
The corresponding control cost is

$$S_i(u_i) = c_i \int_0^{u_i} \text{sig}^{-1}\left(\frac{v}{U_{i, \max}}\right) dv.$$

The state constraint is applied via a penalty function

$$p(x) = \frac{L}{1 + e^{-k(c_s(x) - c_{\max})}} - \frac{L}{1 + e^{-k(c_s(x) - c_{\min})}} + L - \frac{2L}{1 + e^{-k(\mu - c_{\max})}}.$$

For a state constraint of $[-1, 3]$, the penalty function with different parameterizations are shown in the figure below.



Taking both state constraints and control saturation into consideration, the overall cost function has the form

$$\mathbb{E} \left[g(x(T)) + \int_t^T \left(q(x(s)) + p(x(s)) + \sum_{i=1}^m S_i(u_i(s)) \right) ds \mid x(t) = x \right].$$

Adaptive Update Scheme

To ensure numerical stability, we use the square root of state cost variance over a fixed number of iterations as the update threshold, and gradually harden the penalty function $p(x)$. Since the state cost variance would never decrease to zero, we also set a minimum value for the threshold.

$$\begin{aligned} k &\leftarrow k + \delta \\ \delta &\leftarrow \delta - \Delta\delta \\ \beta &\leftarrow \gamma\beta \\ \gamma &\leftarrow \gamma + \Delta \end{aligned}$$

Deep FBSDE

The problem under the updated cost function can be recasted as a forward-backward stochastic differential equation (FBSDE) as shown on the right, where V_x is the partial derivative of the value function w.r.t. the state, Φ represents learned weights, and the Hamiltonian is defined as

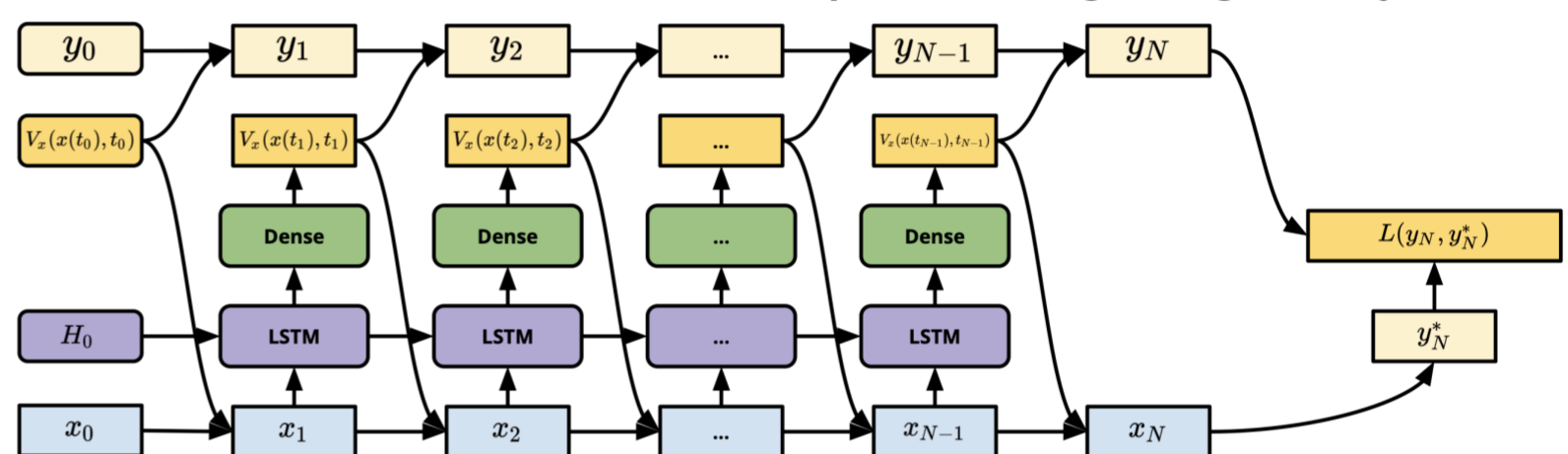
$$h(x, V_x, t, u^*) = q(x) + V_x^T G(x, t)u^*(x, t) + \sum_{i=1}^m S_i(u_i^*).$$

$$\begin{aligned} dy(t) &= \left(-h(x(t), V_x(x(t), t; \theta), t, u(t)) \right. \\ &\quad \left. + V_x^T(x(t), t; \theta)G(x(t), t)u(x(t), t) \right) dt \\ &\quad + V_x^T(x(t), t; \theta)\Sigma(x(t), t)dw(t) \end{aligned}$$

$$\begin{aligned} dx(t) &= \left(f(x(t), t) + G(x(t), t)u(x(t), t) \right) dt \\ &\quad + \Sigma(x(t), t)dw(t) \end{aligned}$$

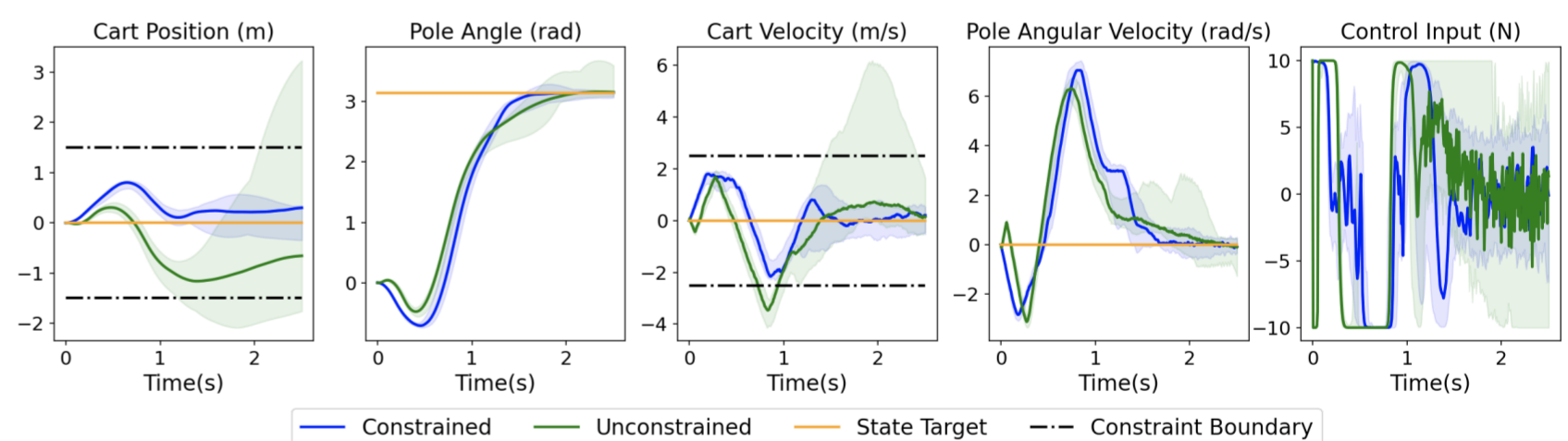
$$\begin{aligned} u(t) &= U_{\max} * \text{sig}(-R^{-1}G^T(x(t), t)V_x(x(t), t; \theta)) \\ y(0) &= V(\phi) \\ dy(0) &= V_x(\phi) \\ x(0) &= x_0. \end{aligned}$$

The Deep FBSDE neural network architecture is shown below, at each time step the neural network estimates V_x . The neural network is optimized using a weight decayed L2 loss.

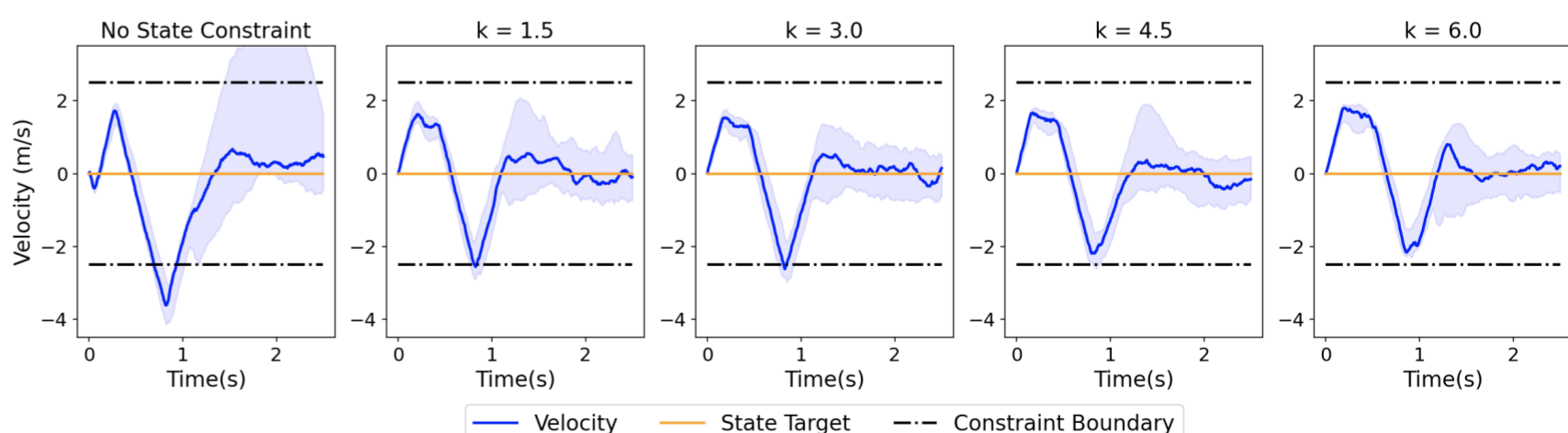


Experiments

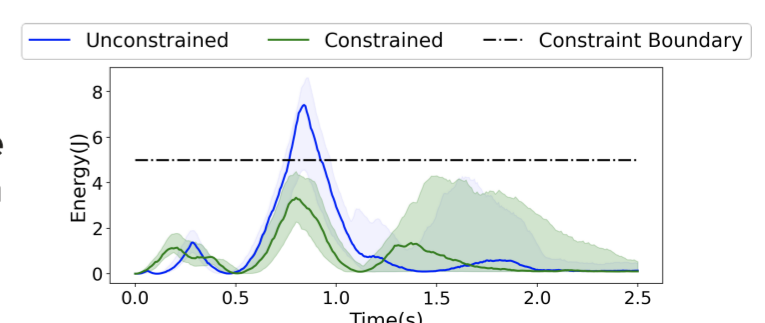
Comparison between constrained and unconstrained controller



Effectiveness of adaptive update scheme



Energy constraint comparison



The presented experimental results are conducted on the cart-pole swing-up task. Two state constraint settings were tested: (i) constraining cart position and cart velocity; (ii) constraining the sum of kinetic and potential energy. We see that in both settings, the learned controller is able to respect the constraint boundaries.