

4/21: model-based RL

Announcements:

- HW4 due Thurs, 4/23
- Guest lecture, Prof. Abhishek Kawrigi, 4/30

Last time:

- sysID

Today:

- model-based RL
- CBFs

Last time:



$$H(z)\ddot{q} + C(z,\dot{q})\dot{q} + g(z) = \tau$$

"system identification"

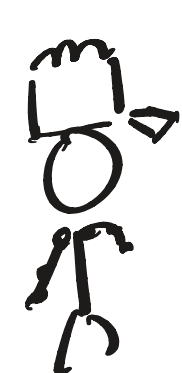
$$D = \sum (x_t, u_t, x_{t+1})$$

$$\min_{\theta} \sum_{t=0}^{T-1} \frac{1}{2} \|x_{t+1} - f(x_t, u_t; \theta)\|^2 = \sum_{t=0}^{T-1} r_t(\theta)^2$$

↑ mujoco → CEM

Q: when does $\theta^* \rightarrow \theta_{true}$

e.g.1 (robo-cheff)



goal: use a fork to pick and twirl spaghetti

u: joint torques

y: encoders, image

x: ? $q, \dot{q}, T_{base}, v_{base}$

T_{fork}, v_{fork}

1. parameters?
2. state space
3. state estimation

things that are hard to simulate:

- fluids
- flexibles
- people

another approach: what if we let our simulator be a neural net?

$$\min_{\theta} \sum_{t=0}^{T-1} \|x_{t+1} - f_{\theta}(x_t, u_t)\|^2$$

↑ deep network

typical name for this: "model-based RL"

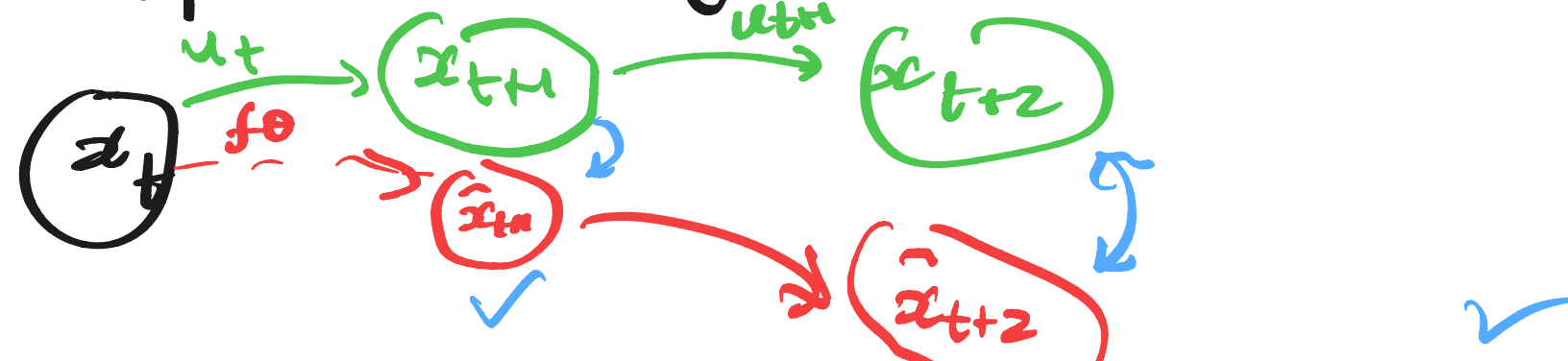
- Alg:
1. do MPC with current model \hat{f}_{θ}
 2. apply those actions on the real system
 3. Add new trajectories $\sum (x_t, u_t, x_{t+1})$ to dataset; update \hat{f}

[Nagabandi '18] [Dreamer '20]
Berkeley Deepmind

Important design decisions:

1. Horizon

suppose f_{θ} is good for single-step



solution: $x_{t+1:T} = f_{\theta}(x_t, u_{t:T-1})$

aside: works for physics-based

$$\hat{x}_{t+1} = f_{\theta}(x_t, u_t)$$

$$\hat{x}_{t+2} = f_{\theta}(\hat{x}_{t+1}, u_{t+1})$$

$$\frac{\partial x_{t+T}}{\partial \theta} = \frac{\partial f_{t+T}}{\partial \theta} + \frac{\partial f_{t+T-1}}{\partial \theta}$$

2. modeling uncertainty

[Nagabandi] $\hat{x}_{t+1} \sim \mathcal{W}_i(\mu_{\theta}(x_t, u_t), \Sigma)$

$$\max_{\theta} \sum_i \log p(x_{t+1} | x_t, u_t, \theta)$$

↑ possibly learned

- why?
1. often our system is stochastic
 2. including ensembles is good for measuring confidence

3. Handling observations

state estimator

optimal control

$$y_{1:T}, u_{0:T-1} \rightarrow \hat{x}_t \Rightarrow \hat{x}_t \rightarrow u_{t:T-1}, x_{t:T}^*$$

"world model"

$$\sum (y_0, u_0, y_1, \dots, y_t)$$

* observations, not actions

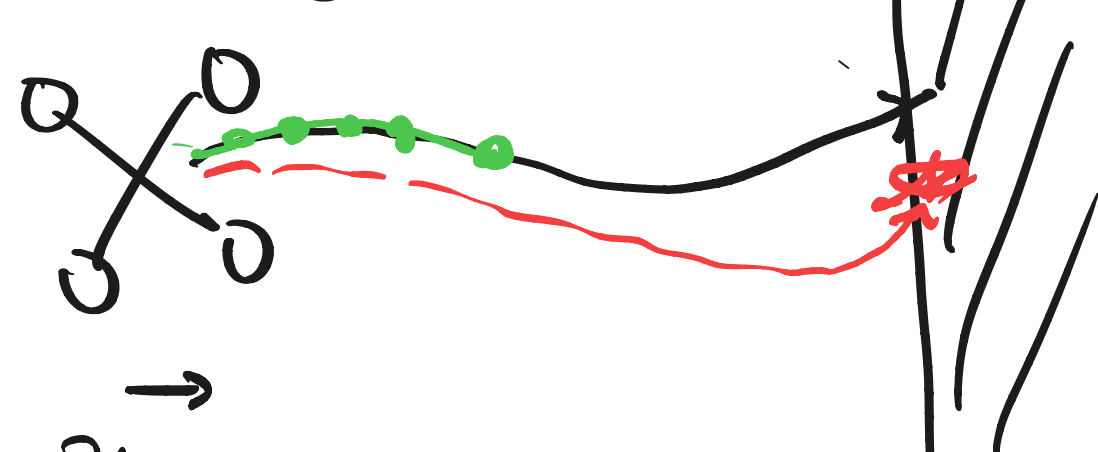
$$z_t = \phi(y_t) \text{ "image encoder"}$$

$$z_{t+1} = f_{\theta}(z_t, u_t) \quad \min_{\phi, \theta} \|y_{t+1} - \phi^{-1}(f_{\theta}(\phi(y_t), u_t))\|$$

$$\hat{y}_{t+1} = \phi^{-1}(z_{t+1})$$

- physics vs. world models
- + right structures (conservation of mass/energy)
 - + very sample efficient
 - observations
 - modeling everything
- don't exhibit object permanence
 - sample inefficient
 - + visual observations
 - + generalize to physics that is hard to model

"safety"



wall $h(x) \geq 0$

$$px \leq 0 - \epsilon$$

ways things go wrong: $x=0$

1. model error
2. MPC can be greedy

answer: add safety constraints @ multiple levels of stack