

L11: Optimal Control

Hao Su

Spring, 2021

Contents are based on Underactuated Robotics taught at MIT by Prof. Russ Tedrake and CS287 taught at UC Berkeley by Prof. Pieter Abbeel.

Agenda

- Fully-actuated v.s. Under-actuated Systems
- Basic Idea of Optimal Control
- Linear Quadratic Control

click to jump to the section.

Fully-actuated v.s. Under-actuated Systems

Review: Control

- A *desired* trajectory to follow: $(\mathbf{q}_d, \dot{\mathbf{q}}_d, \ddot{\mathbf{q}}_d)$
- Forward dynamics $\ddot{\mathbf{q}} = \text{FD}(\mathbf{F}; \mathbf{q}, \dot{\mathbf{q}})$
- Inverse dynamics $\mathbf{F} = \text{ID}(\ddot{\mathbf{q}}; \mathbf{q}, \dot{\mathbf{q}})$
- We use **control** to deal with delay, overshoot, or steady-state error, and ensure stability.
- Steady-state error

$$e = q - q_d$$

Review: Feedforward and Feedback Control

- We need some force to match $\ddot{\mathbf{q}}_d$. This component is called the **feed-forward component**, which comes from $ID(\cdot)$:

$$\mathbf{F}_{ff} = ID(\ddot{\mathbf{q}}_d; \mathbf{q}, \dot{\mathbf{q}})$$

- We also need some additional force to correct the steady-state error, which is called the **feedback component**:

$$\mathbf{F}_{fb} = M(\mathbf{q})(-K_v \dot{e} - K_p e)$$

where $M(\mathbf{q})$ is the inertia of the system.

- The total force we exert to control the system is

$$\mathbf{F} = \mathbf{F}_{ff} + \mathbf{F}_{fb}$$

Review: Computed Torque Control

- Computed Torque Control:

$$\tau = M(\theta)(\ddot{\theta}_d - K_v\dot{e} - K_p e) + C(\theta, \dot{\theta})\dot{\theta} + g(\theta) \quad (1)$$

- By ID, the acceleration under this τ is

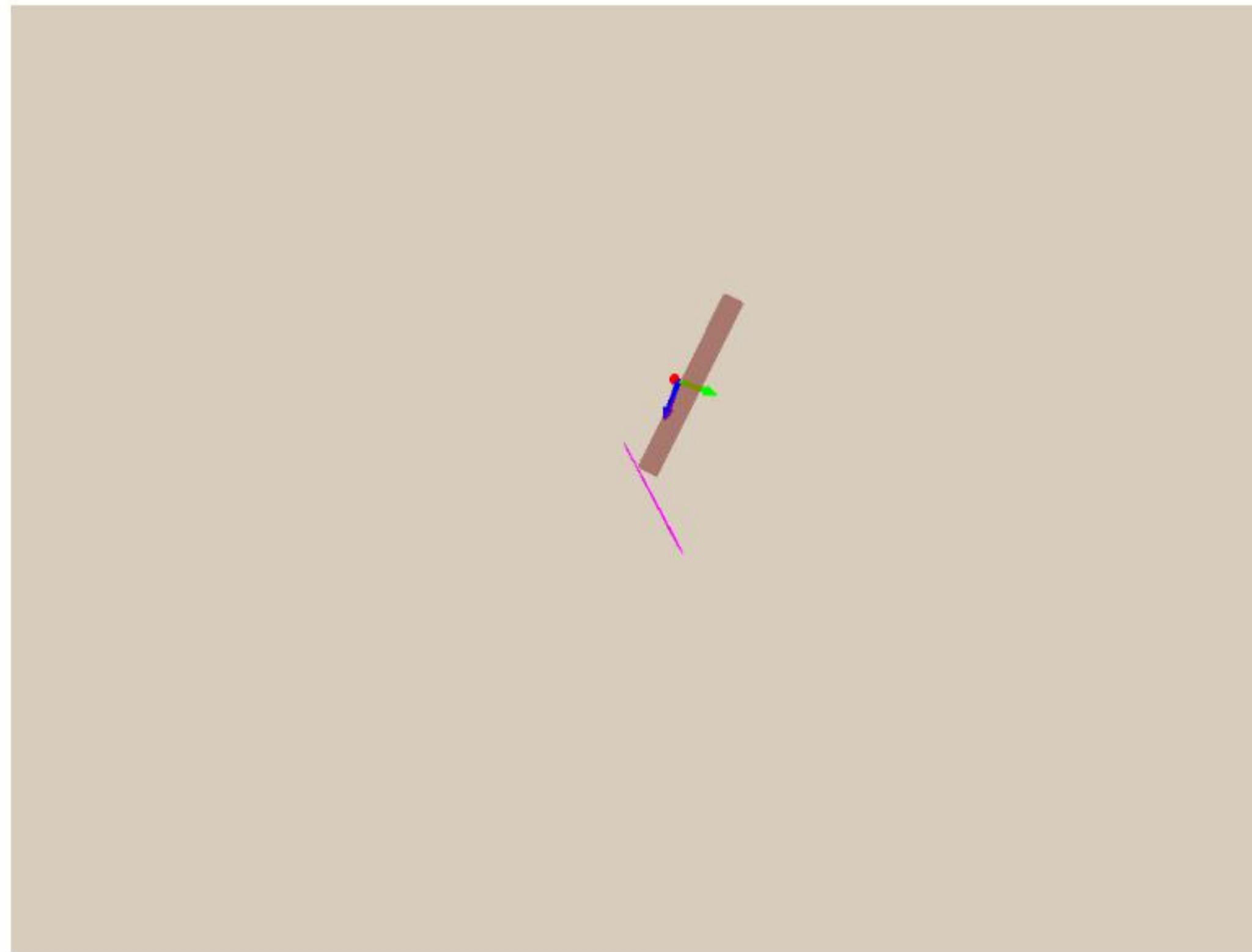
$$\tau = M(\theta)\ddot{\theta} + C(\theta, \dot{\theta})\dot{\theta} + g(\theta) \quad (2)$$

- Subtracting (2) from (1) and cancel $M(\theta)$, we get the error equation:

$$\ddot{e} + K_v\dot{e} + K_p e = 0$$

- Because $K_v, K_p \in \mathbb{S}^+$, by the theory of ODE, $e(t) = \mathcal{O}(e^{\alpha t})$, $\alpha \leq 0$.
- We say that the computed torque control law has **exponential convergence rate**.

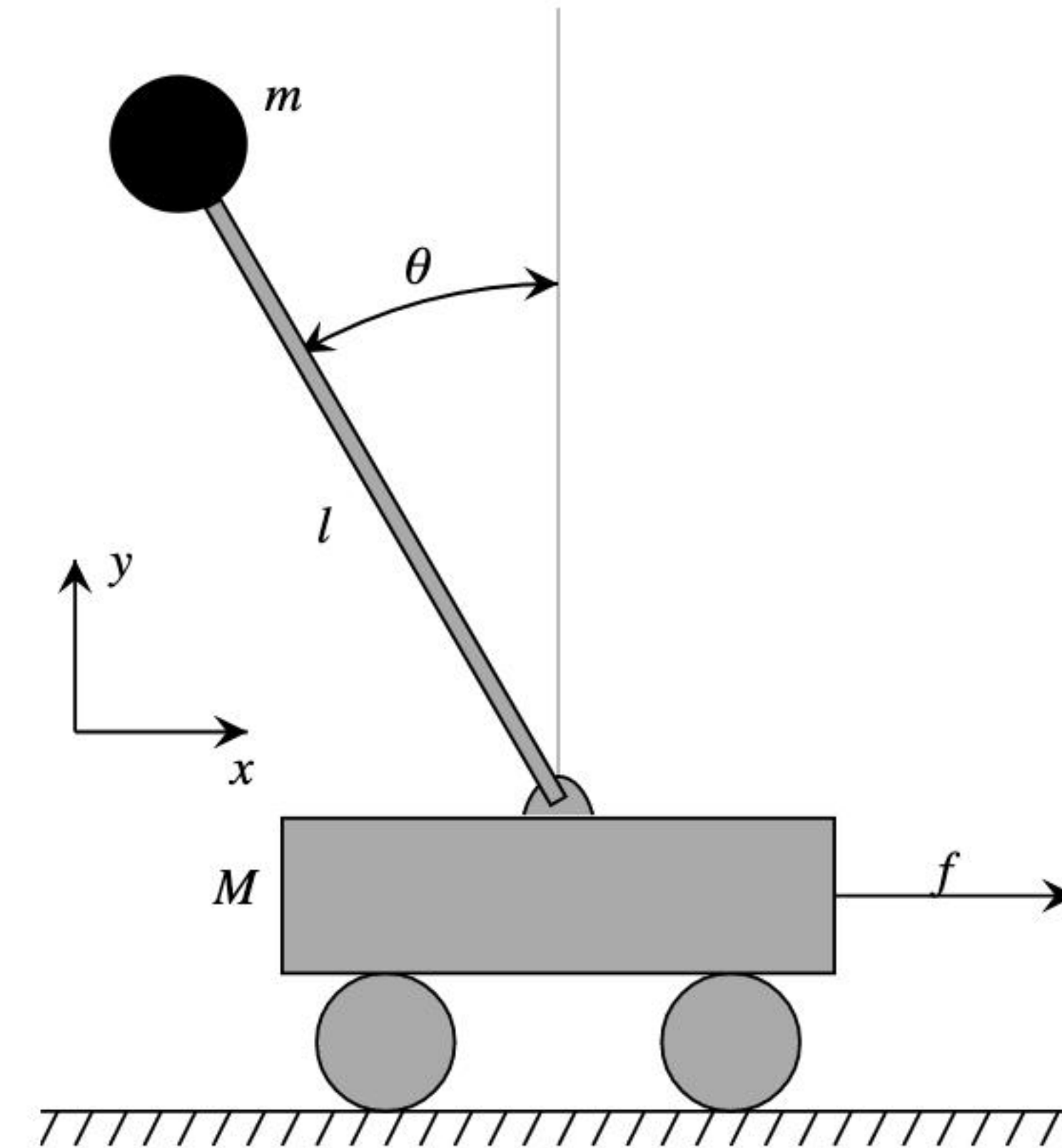
Example of Single-Joint Manipulator (Inverted Pendulum)



Using feedback control (PID), it is very easy to control this swing-up inverted pendulum.

Limitation of PID Control

- Consider the cart-pole example.
- Our desired position is that the rod is upright straight.
- Can we control it by the computed torque law?



A schematic drawing of the inverted pendulum on a cart. The rod is considered massless.

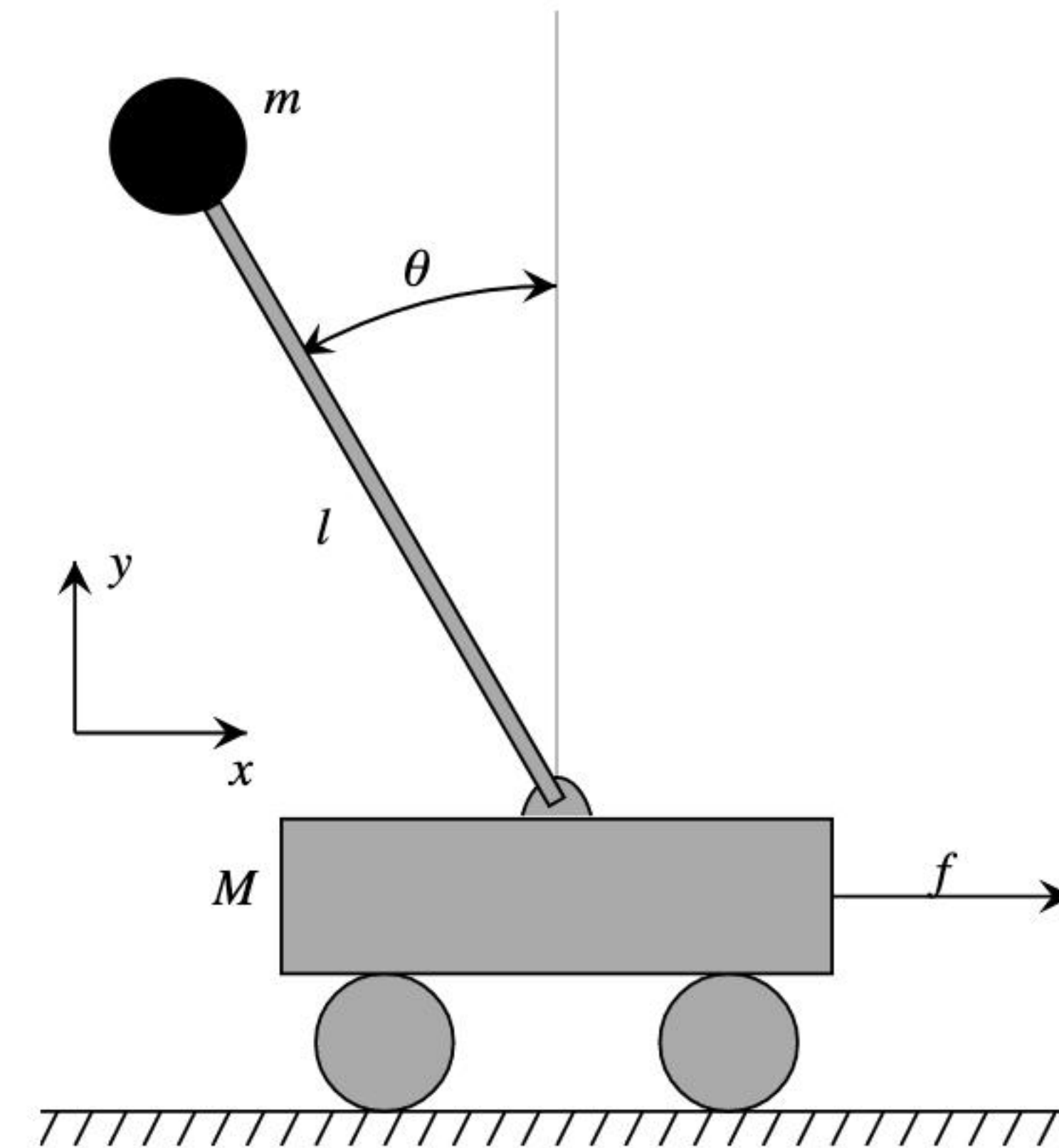
https://en.wikipedia.org/wiki/Inverted_pendulum

Limitation of PID Control

- Consider the cart-pole example.
- Our desired position is that the rod is upright straight.
- Can we control it by the computed torque law?
- Let us try the feedforward-feedback control law.
- First of all, we can only control the 1D force f . It is also obvious that $f_{feedforward} = 0$.
- The feedback force should be

$$f_{feedback} = M(\mathbf{q})(-K_v\dot{e} - K_p e)$$

- Note that $e = \mathbf{q} - \mathbf{q}_d = \begin{bmatrix} x \\ \theta \end{bmatrix} - \begin{bmatrix} x_d \\ \theta_d \end{bmatrix}$, thus the feedback force should be 2D. But we can only control the 1D force!
- Our convergence analysis does not apply here.
- Turns out that tuning the PID for cart pole is hard.



A schematic drawing of the inverted pendulum on a cart. The rod is considered massless.

https://en.wikipedia.org/wiki/Inverted_pendulum

Underactuated Control Differential Equations

- Notions:
 - **Second-order** control dynamical system: $\ddot{\mathbf{q}} = f(\mathbf{q}, \dot{\mathbf{q}}, \mathbf{u}, t)$
 - Control vector: $\mathbf{u} \in \mathcal{U}$
- **Underactuated Control Differential Equations:**

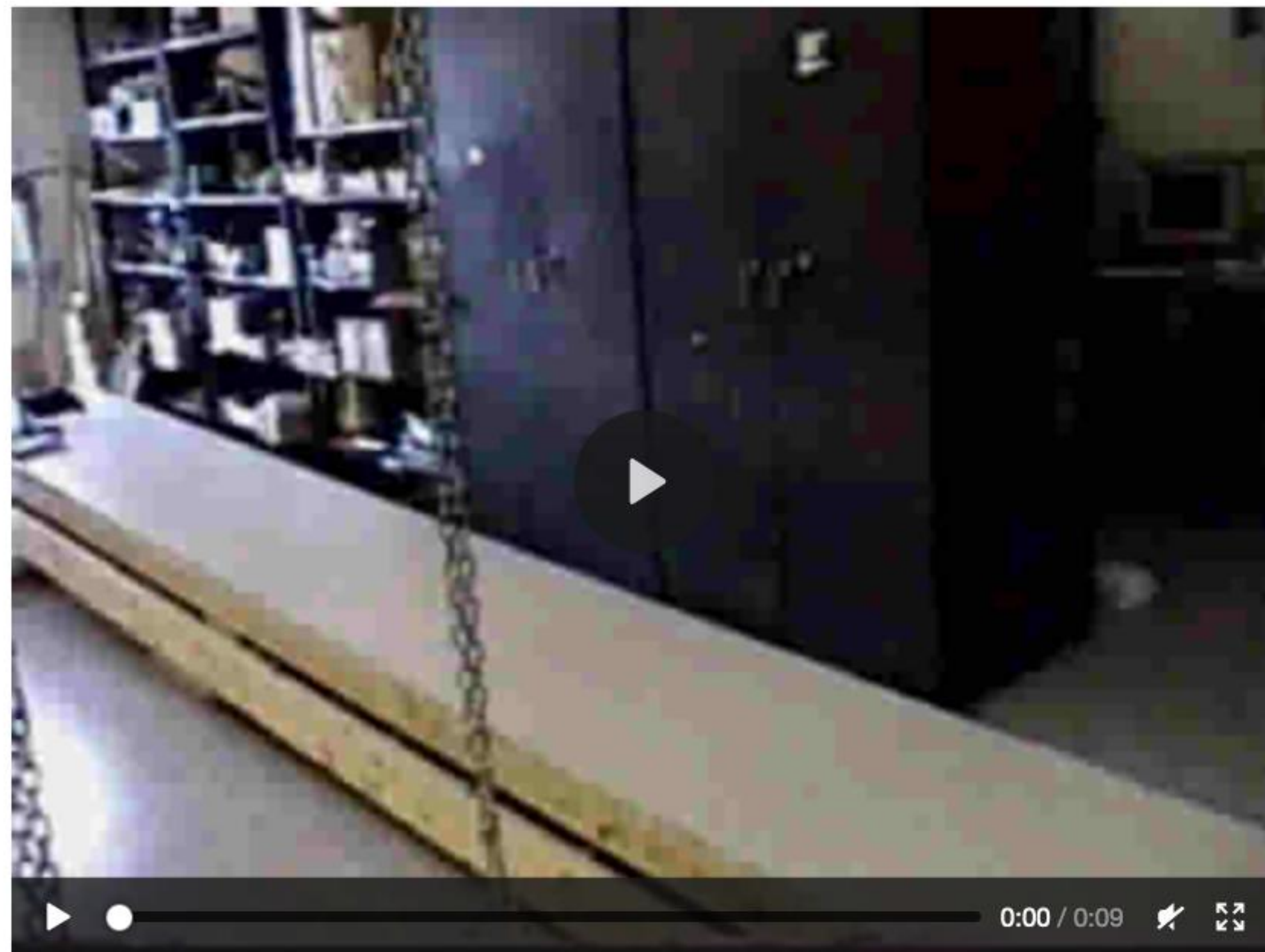
A second-order control differential equation described by the equations

$$\ddot{\mathbf{q}} = f(\mathbf{q}, \dot{\mathbf{q}}, \mathbf{u}, t)$$

*is **fully actuated** in state $\mathbf{x} = (\mathbf{q}, \dot{\mathbf{q}})$ and time t if the resulting map f is surjective: for every $\ddot{\mathbf{q}}$ there exists a \mathbf{u} which produces the desired response. Otherwise, it is **underactuated** (in \mathbf{x} at time t).*

- For example, our cart pole system is underactuated (1-D f has to control \ddot{x} and $\ddot{\theta}$).

The Power of Underactuated System: A Passive Dynamic Walker Example



A 3D passive dynamic walker by Steve Collins and Andy Ruina at Cornell.

Many *Interesting* Problems in Robotics are Underactuated

- **Legged robots are underactuated.** Consider a legged machine with N internal joints and N actuators. If the robot is not bolted to the ground, then the degrees of freedom of the system include both the internal joints and the six degrees of freedom which define the position and orientation of the robot in space. Since $\mathbf{u} \in \mathbb{R}^N$ and $\mathbf{q} \in \mathbb{R}^{N+6}$, the system is underactuated.
- **(Most) Swimming and flying robots are underactuated.** The story is the same here as for legged machines. Each control surface adds one actuator and one DOF. And this is already a simplification, as the true state of the system should really include the (infinite-dimensional) state of the flow.
- **Robot manipulation is (often) underactuated.** Consider a fully-actuated robotic arm. When this arm is manipulating an object with degrees of freedom (even a brick has six), it can become underactuated.

Read by Yourself

Concepts and Main Theoretical Results of Optimal Control

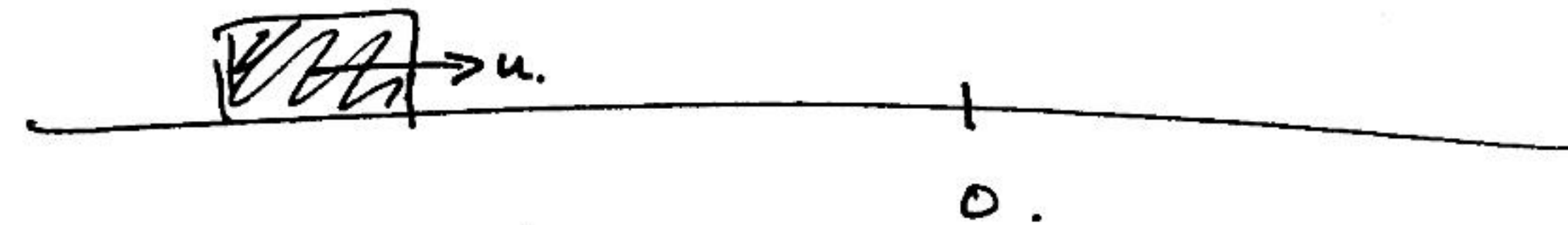
Control as Optimization Problem

- Key idea: **Design an optimization problem whose solution is the control signal.**

- Example

- Consider the simple second-order control dynamics system:

$$\ddot{q} = u, \quad |u| \leq 1$$



credit: Sec 11 of Underactuated Robotics.

- The task is to design a control system, $u = \pi(\mathbf{x}, t)$, $\mathbf{x} = [q, \dot{q}]^T$ to regulate this brick to $\mathbf{x} = [0, 0]^T$.
- Optimization problem:
 - Minimum time: $\min_{\pi} t_f$ subject to $\mathbf{x}(t_0) = \mathbf{x}_0, \mathbf{x}(t_f) = 0$.
 - Quadratic cost: $\min_{\pi} \int_{t_0}^{\infty} \mathbf{x}^T(t) \mathbf{Q} \mathbf{x}(t) dt, \mathbf{Q} \succ 0$.

Transition Function

$$\dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{u}) \quad (\text{transition function})$$

- Example 1 – Brick:

$$\begin{bmatrix} \dot{q} \\ \ddot{q} \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} q \\ \dot{q} \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u,$$

which is of $\dot{\mathbf{x}} = A\mathbf{x} + B\mathbf{u}$ form.

- Example 2 – Manipulator :

$$\begin{bmatrix} \dot{\theta} \\ \ddot{\theta} \end{bmatrix} = \begin{bmatrix} 0 & I \\ 0 & M_{\theta}^{-1} \end{bmatrix} \begin{bmatrix} \theta \\ \dot{\theta} \end{bmatrix} + \begin{bmatrix} 0 \\ M_{\theta}^{-1} C_{\theta, \dot{\theta}} \end{bmatrix} \tau + \begin{bmatrix} 0 \\ -M_{\theta}^{-1} g_{\theta} \end{bmatrix},$$

which is of $\dot{\mathbf{x}} = A_x \mathbf{x} + B_x \mathbf{u} + C_x$ form. Let $\mathbf{x}'_x = \mathbf{x} + A_x^{-1} C_x$, then $\dot{\mathbf{x}}' = A_x \mathbf{x}' + B_x \mathbf{u}$.

Transition Function

$$\dot{\boldsymbol{x}} = f(\boldsymbol{x}, \boldsymbol{u}) \quad (\text{transition function})$$

- Understanding the linear transition function $\dot{\boldsymbol{x}} = A\boldsymbol{x} + B\boldsymbol{u}$ is fundamentally important, because *locally* it approximates any differentiable transitions.
 - For example, we do Taylor's expansion around stationary state $\boldsymbol{x}^*, 0 = f(\boldsymbol{x}^*, \boldsymbol{u}^*)$.
 - Stationary state: both velocity and acceleration are 0.
 - Take $\boldsymbol{x}' = \boldsymbol{x} - \boldsymbol{x}^*, \boldsymbol{u}' = \boldsymbol{u} - \boldsymbol{u}^*$, plug in the transition function, and use Taylor's expansion to approximate f around \boldsymbol{x}^* and \boldsymbol{u}^* , then

$$\dot{\boldsymbol{x}}' \approx \nabla_{\boldsymbol{x}} f(\boldsymbol{x}^*, \boldsymbol{u}^*) \boldsymbol{x}' + \nabla_{\boldsymbol{u}} f(\boldsymbol{x}^*, \boldsymbol{u}^*) \boldsymbol{u}'$$

Additive Cost

$$\int_0^T \ell(\mathbf{x}(t), \mathbf{u}(t)) dt$$

(additive cost)

- For example, the quadratic cost for the brick example is an additive cost:

$$\min_{\pi} \int_{t_0}^{\infty} \mathbf{x}^T(t) \mathbf{Q} \mathbf{x}(t) dt, \quad \mathbf{Q} \succ 0.$$

- Additive cost is a favorable choice in optimal control, because
 - it admits an optimality condition of elegant form.
 - in discrete case, it also implies a "dynamic programming" solution (we will see in the next lecture of reinforcement learning).

Cost-to-go Function

- Consider a *time-invariant* dynamic system $\dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{u})$ with an *infinite-horizon* additive cost $\int_0^\infty \ell(\mathbf{x}, \mathbf{u}) dt$
 - Time-invariant: f and ℓ do not directly depend on t .
- Suppose that we will control the system using a **policy** $\mathbf{u} = \pi(\mathbf{x})$
- The **cost-to-go function** at a certain starting state \mathbf{x}_0 is:

$$J^\pi(\mathbf{x}_0) = \int_0^\infty \ell(\mathbf{x}(t), \mathbf{u}(t)) dt \quad (\text{cost-to-go function})$$

where $\mathbf{x}(0) = \mathbf{x}_0$ and $\mathbf{u}(t) = \pi(\mathbf{x}(t))$.

- You will see that the cost-to-go-function and the value function in reinforcement learning are essentially the same thing, except that one needs to be minimized and the other maximized.

The Hamilton-Jacobi-Bellman Equation (Optimality Condition)

- Consider a time-invariant system $\dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{u})$ with an additive cost $\int_0^\infty \ell(\mathbf{x}, \mathbf{u}) dt$, in which $f, \ell \in C^\infty(\mathcal{X} \times \mathcal{U})$, $\text{Hess}(\ell) \succ 0$, $\ell(\mathbf{x}) = 0$ iff $\mathbf{x} = \mathbf{x}^*$.
- Under some technical conditions on the existence and boundedness of solutions (see *Thm 7.1 in Underactuated Robotics*), a sufficient condition for the existence of optimal policy is:

$\exists J : \mathcal{X} \rightarrow \mathbb{R} \in C^\infty(\mathcal{X})$, $\text{Hess}(J) \succ 0$ such that $\forall \mathbf{x} \in \mathcal{X}$,

$$\min_{\mathbf{u} \in \mathcal{U}} \left[\ell(\mathbf{x}, \mathbf{u}) + \left(\frac{\partial J(\mathbf{x})}{\partial \mathbf{x}} \right)^T f(\mathbf{x}, \mathbf{u}) \right] = 0 \quad (\text{HJB equation})$$

- $J(\mathbf{x})$ is the **cost-to-go** function.
- The **optimal policy** is $\pi^*(\mathbf{x}) = \arg \min_{\mathbf{u} \in \mathcal{U}} \left[\ell(\mathbf{x}, \mathbf{u}) + \left(\frac{\partial J(\mathbf{x})}{\partial \mathbf{x}} \right)^T f(\mathbf{x}, \mathbf{u}) \right]$

The Hamilton-Jacobi-Bellman Equation (Optimality Condition)

- The HJB equation is the continuous version of the Bellman equation in discrete dynamic programming/reinforcement learning.
- A (sloppy) justification:
 - The Bellman equation (for optimal policy π^*) of the discrete problem:

$$J(\mathbf{x}_t) = \ell(\mathbf{x}_t, \mathbf{u}_t)\Delta t + J(\mathbf{x}_{t+1}), \text{ where } \frac{\mathbf{x}_{t+1} - \mathbf{x}_t}{\Delta t} = f(\mathbf{x}_t, \mathbf{u}_t) \text{ and } \mathbf{u}_t = \pi^*(\mathbf{x}_t)$$

- Therefore,

$$\begin{aligned} 0 = \ell(\mathbf{x}_t, \mathbf{u}_t) + \frac{J(\mathbf{x}_{t+1}) - J(\mathbf{x}_t)}{\Delta t} &\iff 0 = \ell(\mathbf{x}_t, \mathbf{u}_t) + \frac{J(\mathbf{x}_{t+1}) - J(\mathbf{x}_t)}{\mathbf{x}_{t+1} - \mathbf{x}_t} \frac{\mathbf{x}_{t+1} - \mathbf{x}_t}{\Delta t} \\ &\iff 0 = \ell(\mathbf{x}_t, \mathbf{u}_t) + \frac{J(\mathbf{x}_{t+1}) - J(\mathbf{x}_t)}{\mathbf{x}_{t+1} - \mathbf{x}_t} f(\mathbf{x}_t, \mathbf{u}_t) \end{aligned}$$

- Let $\Delta t \rightarrow 0$, $0 = \ell(\mathbf{x}_t, \mathbf{u}_t) + \left(\frac{\partial J(\mathbf{x})}{\partial \mathbf{x}}\right)^T f(\mathbf{x}, \mathbf{u})$

Linear-Quadratic Regulator

Slides are based on CH8 of Underactuated Robotics

Canonical Form

- Consider a linear time-invariant system in state-space form:

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$$

with the objective function to minimize:

$$\int_0^{\infty} \ell(t) dt = \int_0^{\infty} [\mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u}] dt$$

in which $\mathbf{Q} = \mathbf{Q}^T \succeq 0$, $\mathbf{R} = \mathbf{R}^T \succ 0$ are two designed matrices.

- Remarks about the objective function:
 - The cost $\ell(\mathbf{x}, \mathbf{u})$ is quadratic: brings in advantage in doing analysis on the HJB equation.
 - The cost involves penalty for large control signals (the $\mathbf{u}^T \mathbf{R} \mathbf{u}$ term), which favors *low energy* to achieve a certain goal.

State Tracking Example

- Sometimes, we need a bit of derivation to arrive at the canonical form.
- Consider the task of tracking a desired state trajectory $\mathbf{x}_d(t)$. It is natural to write down the following optimization problem:

$$\underset{\mathbf{u} \in \mathcal{U}}{\text{minimize}} \quad \int_0^\infty (\mathbf{x} - \mathbf{x}_d)^T \mathbf{Q} (\mathbf{x} - \mathbf{x}_d) + \mathbf{u}^T \mathbf{R} \mathbf{u}$$

in which $\mathbf{Q} \succeq 0$, $\mathbf{R} \succ 0$, and $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$.

- Introduce $\tilde{\mathbf{x}} = \begin{bmatrix} \mathbf{x} - \mathbf{x}_d \\ \mathbf{x}_d \end{bmatrix}$, $\tilde{\mathbf{A}} = \begin{bmatrix} \mathbf{A} & \mathbf{A} \\ 0 & 0 \end{bmatrix}$, $\tilde{\mathbf{B}} = \begin{bmatrix} \mathbf{B} \\ 0 \end{bmatrix}$, $\tilde{\mathbf{Q}} = \begin{bmatrix} \mathbf{Q} & 0 \\ 0 & 0 \end{bmatrix}$, and the above optimization problem becomes the canonical form:

$$\underset{\mathbf{u} \in \mathcal{U}}{\text{minimize}} \quad \int_0^\infty \tilde{\mathbf{x}}^T \tilde{\mathbf{Q}} \tilde{\mathbf{x}} + \mathbf{u}^T \mathbf{R} \mathbf{u}$$

in which $\mathbf{Q} \succeq 0$, $\mathbf{R} \succ 0$, and $\dot{\tilde{\mathbf{x}}} = \tilde{\mathbf{A}}\tilde{\mathbf{x}} + \tilde{\mathbf{B}}\mathbf{u}$.

LQR Controller (Infinite-horizon, Closed-form)

- Recall the HJB optimality condition:

$$\exists J : \mathcal{X} \rightarrow \mathbb{R} \in \mathcal{C}^\infty(\mathcal{X}), \text{Hess}(J) \succ 0 \text{ such that } \forall \mathbf{x} \in \mathcal{X},$$

$$\min_{\mathbf{u} \in \mathcal{U}} L(\mathbf{x}, \mathbf{u}) = 0, \quad \text{where } L(\mathbf{x}, \mathbf{u}) = \ell(\mathbf{x}, \mathbf{u}) + \left(\frac{\partial J(\mathbf{x})}{\partial \mathbf{x}} \right)^T f(\mathbf{x}, \mathbf{u})$$

- For our linear-quadratic system setup, it is known that $J(\mathbf{x}) = \mathbf{x}^T \mathbf{S} \mathbf{x}$.
- Let us plug this $J(\mathbf{x})$ with unknown \mathbf{S} into $L(\mathbf{x}, \mathbf{u})$ to solve \mathbf{S} .
- Sketch:
 1. For some given \mathbf{x} , we first find $\mathbf{u}^* = \arg \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u})$;
 2. Then we plug \mathbf{u}^* back to $L(\mathbf{x}, \mathbf{u})$ and obtain an equation system that includes \mathbf{S}
 3. Solve \mathbf{S} from the equation system.

LQR Controller (Infinite-horizon, Closed-form)

Details:

- By $\ell(\mathbf{x}, \mathbf{u}) = \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u}$, $f(\mathbf{x}, \mathbf{u}) = \mathbf{A} \mathbf{x} + \mathbf{B} \mathbf{u}$, and $J(\mathbf{x}) = \mathbf{x}^T \mathbf{S} \mathbf{x}$,

$$L(\mathbf{x}, \mathbf{u}) = \ell(\mathbf{x}, \mathbf{u}) + \left(\frac{\partial J(\mathbf{x})}{\partial \mathbf{x}} \right)^T f(\mathbf{x}, \mathbf{u}) = \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u} + 2\mathbf{x}^T \mathbf{S} \mathbf{A} \mathbf{x} + 2\mathbf{x}^T \mathbf{S} \mathbf{B} \mathbf{u}$$

- To solve \mathbf{u}^* given \mathbf{x} , we have $\frac{\partial L}{\partial \mathbf{u}^*} = 2\mathbf{R} \mathbf{u}^* + 2\mathbf{B}^T \mathbf{S} \mathbf{x} = 0$

$$\therefore \mathbf{u}^* = -\mathbf{R}^{-1} \mathbf{B}^T \mathbf{S} \mathbf{x} = -\mathbf{K} \mathbf{x}$$

- Plug \mathbf{u}^* back into L and by $\mathbf{x}^T \mathbf{S} \mathbf{A} \mathbf{x} = \mathbf{x}^T \mathbf{A}^T \mathbf{S} \mathbf{x}$,

$$\min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) = \mathbf{x}^T (\mathbf{Q} - \mathbf{S} \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^T \mathbf{S} + \mathbf{S} \mathbf{A} + \mathbf{A}^T \mathbf{S}) \mathbf{x}$$

- Since $\min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) \equiv 0, \forall \mathbf{x} \in \mathcal{X}$, it must be true that

$$\mathbf{Q} - \mathbf{S} \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^T \mathbf{S} + \mathbf{S} \mathbf{A} + \mathbf{A}^T \mathbf{S} = 0 \quad (\text{algebraic Riccati equation})$$

Algebraic Riccati Equation

Our goal is to solve \mathbf{S} from the algebraic Riccati equation:

$$\mathbf{Q} - \mathbf{SBR}^{-1}\mathbf{B}^T\mathbf{S} + \mathbf{SA} + \mathbf{A}^T\mathbf{S} = 0$$

The existence of solution depends on a so-called "controllable" condition.

- Controllable:
 - A system is said to be **controllable** if we can reach any target state \mathbf{x}^* from any start state \mathbf{x}_0 .
 - A system is said to be **t -time controllable** if we can reach any target state \mathbf{x}^* from any start state \mathbf{x}_0 within t time period.
- The equation has a single positive-definite solution if and only if the system is controllable.
- There are good numerical methods for finding that solution, even in high-dimensional problems.

A Bit of Retrospection

- Our $\mathbf{u}^* = -\mathbf{K}\mathbf{x}$. Plug in the transition function $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$, and

$$\dot{\mathbf{x}} = (\mathbf{A} - \mathbf{B}\mathbf{K})\mathbf{x}$$

- The solution takes the form

$$\mathbf{x}(t) = e^{(\mathbf{A} - \mathbf{B}\mathbf{K})t} \mathbf{x}(0),$$

- Plug $\mathbf{x}(t)$ in the definition of the cost function, and we see that the cost takes the form

$$J = \mathbf{x}^T(0) \mathbf{S} \mathbf{x}(0),$$

which is a quadratic form, consistent with our assumption that J is quadratic.

iterative LQR (iLQR) for Non-linear System

- Consider the general optimal control problem:

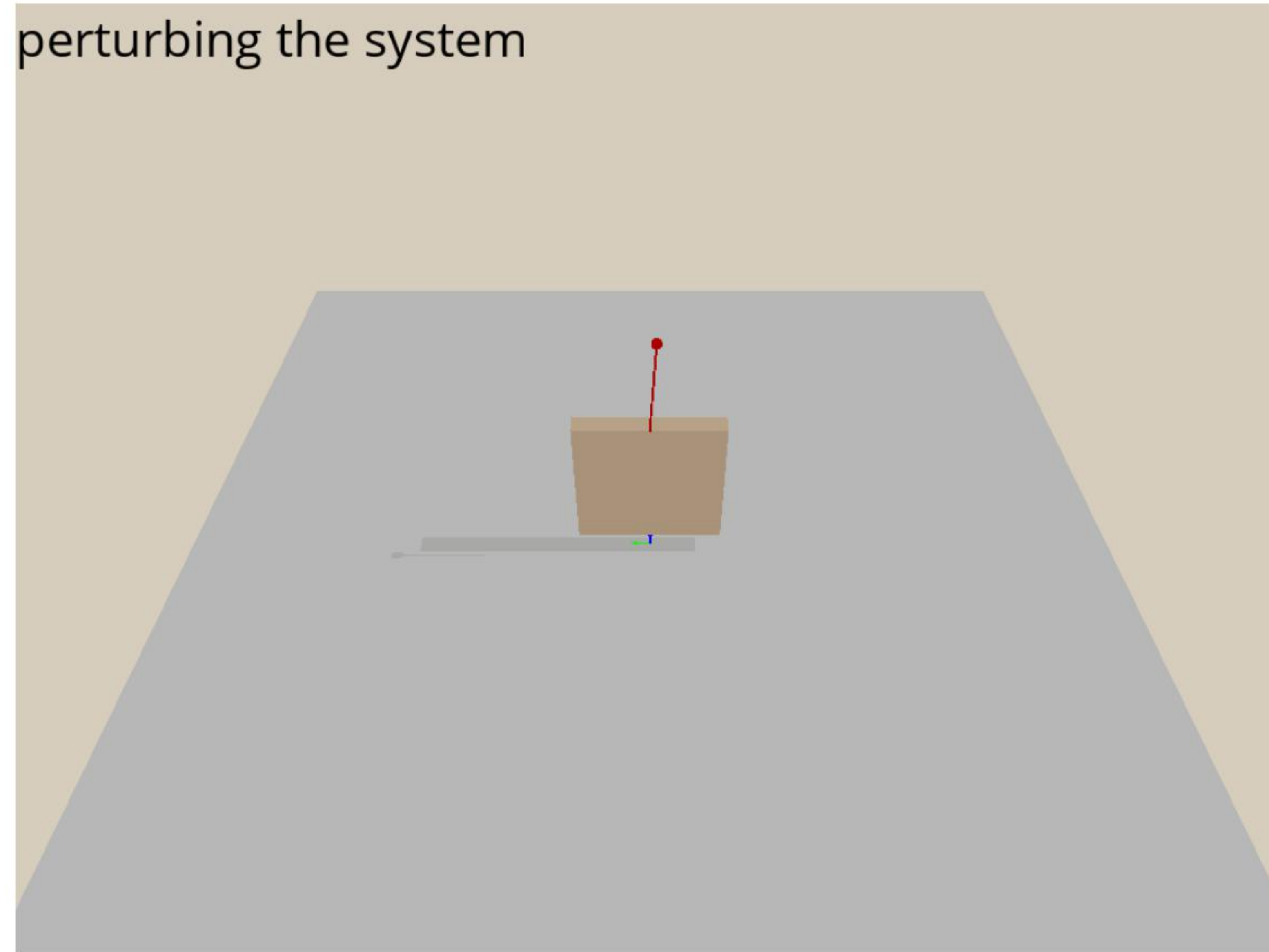
$$\dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{u})$$

with the general objective function to minimize:

$$\int_0^{\infty} \ell(t) dt$$

- We can build a local LQR problem at some key time steps and iteratively apply the LQR controller:
 1. Compute the first-order Taylor expansion of the dynamics model $\dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{u})$ and the second-order Taylor expansion of the cost function $\ell(\mathbf{x}, \mathbf{u})$;
 2. Use the LQR to solve the optimal control policy and execute the policy for Δt ;
 3. Go to (1) and recompute the approximation;

Example: Cart Pole



However, it does not always converge, if the initial perturbation stage gets longer.

Discrete-time LQR

Contents are based on CS287-FA19 by Prof. Pieter Abbeel at UC Berkeley

Problem Setup

- Transition function: $\mathbf{x}_{t+1} = A\mathbf{x}_t + B\mathbf{u}_t$
 \mathbf{x}_t : state at time t
 \mathbf{u}_t : input at time t
- 1-step cost function: $\ell(\mathbf{x}_t, \mathbf{u}_t) = \mathbf{x}_t^T Q \mathbf{x}_t + \mathbf{u}_t^T R \mathbf{u}_t$ with $Q \succeq 0, R \succ 0$
- We consider the **finite-horizon** setup, which needs more delicate treatment compared with the infinite-horizon setup.
- Suppose that our total cost only includes the cost for the first n steps:

$$\sum_{t=0}^{n-1} \ell(\mathbf{x}_t, \mathbf{u}_t),$$

Finite-horizon Cost-to-go Function

- We introduce the k -step cost-to-go function as:

$$J_k^\pi(\mathbf{x}) = \sum_{t=0}^{k-1} \ell(\mathbf{x}_t, \mathbf{u}_t),$$

where $\mathbf{x}_0 = \mathbf{x}$ and $\mathbf{u}_t = \pi(\mathbf{x}_t)$.

- The optimal policy must take the optimal action at each step. Therefore, the $(k + 1)$ -step cost-to-go function for the optimal policy is:

$$J_{k+1}^{\pi^*}(\mathbf{x}) = \min_{\mathbf{u}} [\mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u} + J_k^{\pi^*}(\mathbf{A} \mathbf{x} + \mathbf{B} \mathbf{u})] \quad (\text{Bellman equation})$$

Dynamic Programming

- Final state: When no step can be taken, there is a constant cost: $J_0(\mathbf{x}) = \mathbf{x}^T \mathbf{Q} \mathbf{x}$.

$$\boxed{J_0(\mathbf{x}) = \mathbf{x}^T \mathbf{P}_0 \mathbf{x}, \quad \text{where } \mathbf{P}_0 = \mathbf{Q}}$$

- Last step (n -th step):

$$\begin{aligned} J_1(\mathbf{x}) &= \min_{\mathbf{u}} [\mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u} + J_0(\mathbf{A} \mathbf{x} + \mathbf{B} \mathbf{u})] \\ &= \min_{\mathbf{u}} [\mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u} + (\mathbf{A} \mathbf{x} + \mathbf{B} \mathbf{u})^T \mathbf{P}_0 (\mathbf{A} \mathbf{x} + \mathbf{B} \mathbf{u})] \end{aligned} \quad (1)$$

Setting the gradient w.r.t. to \mathbf{u} to zero and solve \mathbf{u}_1 :

$$\mathbf{u}_1 = -(\mathbf{R} + \mathbf{B}^T \mathbf{Q} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{P}_0 \mathbf{A} \mathbf{x} = -\mathbf{K}_1 \mathbf{x} \quad (2)$$

(2) into (1):

$$\boxed{J_1(\mathbf{x}) = \mathbf{x}^T \mathbf{P}_1 \mathbf{x}, \quad \text{where } \mathbf{P}_1 = \mathbf{Q} + \mathbf{K}_1^T \mathbf{R} \mathbf{K}_1 + (\mathbf{A} - \mathbf{B} \mathbf{K}_1)^T \mathbf{P}_0 (\mathbf{A} - \mathbf{B} \mathbf{K}_1)}$$

Dynamic Programming (cont')

- $(n - 1)$ -th step:

$$J_2(\mathbf{x}) = \mathbf{x}^T \mathbf{P}_2 \mathbf{x}, \quad \text{where } \mathbf{P}_2 = \mathbf{Q} + \mathbf{K}_2^T \mathbf{R} \mathbf{K}_2 + (\mathbf{A} - \mathbf{B} \mathbf{K}_2)^T \mathbf{P}_1 (\mathbf{A} - \mathbf{B} \mathbf{K}_2)$$

Control signal: $\mathbf{u}_2 = -\mathbf{K}_2 \mathbf{x}$, where $\mathbf{K}_2 = (\mathbf{R} + \mathbf{B}^T \mathbf{P}_1 \mathbf{B})^{-1} \mathbf{B}^T \mathbf{P}_1 \mathbf{A}$.

- Notice the similarity between \mathbf{P}_1 and \mathbf{P}_2 .

Summary of Dynamic Programming for LQR

1. Set $\mathbf{P}_0 = \mathbf{Q}$
2. for $i = 1, 2, 3, \dots$

$$\mathbf{K}_i = (\mathbf{R} + \mathbf{B}^T \mathbf{P}_{i-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{P}_{i-1} \mathbf{A}$$
$$\mathbf{P}_i = \mathbf{Q} + \mathbf{K}_i^T \mathbf{R} \mathbf{K}_i + (\mathbf{A} - \mathbf{B} \mathbf{K}_i)^T \mathbf{P}_{i-1} (\mathbf{A} - \mathbf{B} \mathbf{K}_i)$$

The optimal policy for a i -step horizon is given by:

$$\pi(\mathbf{x}) = \mathbf{K}_i \mathbf{x}$$

The cost-to-go function for a i -step horizon is given by: $J_i(\mathbf{x}) = \mathbf{x}^T \mathbf{P}_i \mathbf{x}$

Some Last Words about discrete LQR

- Solving the infinite-horizon optimal policy in closed-form for discrete LQR is not easy.
- However, if we decide to back-up for infinite steps, DP converges to the infinite-horizon optimal policy *if and only if* the dynamics (A, B) is such that there exists a policy that can drive the state to 0 (*source of information*).
- If converged, it is often most convenient to use the steady-state feedback K for all times.
- Similar to the infinite-horizon LQR example, we can extend our method to more sophisticated scenarios using iterative LQR (iLQR).

What I would like to but do not have the bandwidth to cover

- Stability theory
- Lyapunov analysis with convex optimization
- Trajectory optimization-based control
- Robust control
- Discrete/Continuous Hybrid control
- ...

Recommend to check out *Underactuated Robots* by Prof. Russ Tedrake if you are interested in these topics.

L11: Optimal Control

Blue Sky

Spring 2021

Agenda

- Fully-actuated v.s. Under-actuated Systems
- Basic Idea of Optimal Control
- Linear Quadratic Control

Fully-actuated v.s. Under-actuated Systems

Review: Control

- Control systems: transfer function, block diagram
- Block diagram reduction
- State-space representation
- Control design: PID, root locus, frequency response

Review: Feedforward and Feedback Control

- Feedforward control: disturbance rejection
- Feedback control: stability, performance
- Combined feedforward and feedback control

Review: Computed Torque Control

- Computed torque control: model-based control
- Control law: inverse dynamics
- Control design: stability, performance

Example of Single-joint Manipulator (Inverted Pendulum)



Limitation of PID Control

- Steady-state error
- Integral windup
- Derivative kick
- Control signal saturation

