Optimal Control, Trajectory Optimization, and Planning

CS 294-112: Deep Reinforcement Learning

Week 2, Lecture 2

Sergey Levine

Announcements

- 1. Assignment 1 will be out next week
- 2. Friday section
 - Review of automatic differentiation, SGD, training neural nets
 - Try the MNIST TensorFlow tutorial if you're having trouble, come to the section!
 - Fri 1/27 at 10 am
 - Sutardja Dai Hall 240
 - Chelsea Finn will teach the section

Overview

1. Last lecture: imitation learning from a human teacher



- 2. Today: can the machine make its own decisions?
 - a. How can we choose actions under perfect knowledge of the system dynamics?
 - b. Optimal control, trajectory optimization, planning
- 3. Next week: how can we learn *unknown* dynamics?
- 4. How can we then also learn policies? (e.g. by imitating optimal control)



Today's Lecture

- 1. Making decisions under known dynamics
 - Definitions & problem statement
- 2. Trajectory optimization: backpropagation through dynamical systems
- 3. Linear dynamics: linear-quadratic regulator (LQR)
- 4. Nonlinear dynamics: differential dynamic programming (DDP) & iterative LQR
- 5. Discrete systems: Monte-Carlo tree search (MCTS)
- 6. Case study: imitation learning from MCTS
- Goals:
 - Understand the terminology and formalisms of optimal control
 - Understand some standard optimal control & planning algorithms

Terminology & notation



$$\min_{\mathbf{u}_1,\ldots,\mathbf{u}_T} \frac{\sum_{t=1}^T p(\mathbf{x}_t,\mathbf{u}_t)}{\sum_{t=1}^T p(\mathbf{x}_t,\mathbf{u}_t)} \text{by.ttige}_t |\mathbf{u}_1 f(\mathbf{x}_t,\mathbf{u}_T)_{t-1})$$

Trajectory optimization

$$\min_{\mathbf{u}_1,\ldots,\mathbf{u}_T} \sum_{t=1}^T c(\mathbf{x}_t,\mathbf{u}_t) \text{ s.t. } \mathbf{x}_t = f(\mathbf{x}_{t-1},\mathbf{u}_{t-1})$$

$$\min_{\mathbf{u}_1,\ldots,\mathbf{u}_T} c(\mathbf{x}_1,\mathbf{u}_1) + c(f(\mathbf{x}_1,\mathbf{u}_1),\mathbf{u}_2) + \cdots + c(f(f(\ldots)\ldots),\mathbf{u}_T)$$

usual story: differentiate via backpropagation and optimize!

need
$$\frac{df}{d\mathbf{x}_t}, \frac{df}{d\mathbf{u}_t}, \frac{dc}{d\mathbf{x}_t}, \frac{dc}{d\mathbf{u}_t}$$

in practice, it really helps to use a 2^{nd} order method!

Shooting methods vs collocation

shooting method: optimize over actions only

 $\min_{\mathbf{u}_1,\ldots,\mathbf{u}_T} c(\mathbf{x}_1,\mathbf{u}_1) + c(f(\mathbf{x}_1,\mathbf{u}_1),\mathbf{u}_2) + \cdots + c(f(f(\ldots)\ldots),\mathbf{u}_T)$



Shooting methods vs collocation

collocation method: optimize over actions and states, with constraints

$$\min_{\mathbf{u}_1,\ldots,\mathbf{u}_T,\mathbf{x}_1,\ldots,\mathbf{x}_T} \sum_{t=1}^T c(\mathbf{x}_t,\mathbf{u}_t) \text{ s.t. } \mathbf{x}_t = f(\mathbf{x}_{t-1},\mathbf{u}_{t-1})$$





Linear case: LQR

$$\begin{array}{l} \underset{\mathbf{u}_{1},\ldots,\mathbf{u}_{T}}{\min} c(\mathbf{x}_{1},\mathbf{u}_{1}) + c(f(\mathbf{x}_{1},\mathbf{u}_{1}),\mathbf{u}_{2}) + \cdots + c(f(f(\ldots)\ldots),\mathbf{u}_{T}) \\ c(\mathbf{x}_{t},\mathbf{u}_{t}) = \frac{1}{2} \begin{bmatrix} \mathbf{x}_{t} \\ \mathbf{u}_{t} \end{bmatrix}^{T} \mathbf{C}_{t} \begin{bmatrix} \mathbf{x}_{t} \\ \mathbf{u}_{t} \end{bmatrix} + \begin{bmatrix} \mathbf{x}_{t} \\ \mathbf{u}_{t} \end{bmatrix}^{T} \mathbf{c}_{t} \quad \begin{array}{c} \text{only term that} \\ \text{depends on } \mathbf{u}_{T} \end{array}$$

$$f(\mathbf{x}_{t},\mathbf{u}_{t}) = \mathbf{F}_{t} \begin{bmatrix} \mathbf{x}_{t} \\ \mathbf{u}_{t} \end{bmatrix} + \mathbf{f}_{t}$$

Base case: solve for \mathbf{u}_T only

$$Q(\mathbf{x}_T, \mathbf{u}_T) = \text{const} + \frac{1}{2} \begin{bmatrix} \mathbf{x}_T \\ \mathbf{u}_T \end{bmatrix}^T \mathbf{C}_T \begin{bmatrix} \mathbf{x}_T \\ \mathbf{u}_T \end{bmatrix} + \begin{bmatrix} \mathbf{x}_T \\ \mathbf{u}_T \end{bmatrix}^T \mathbf{c}_T$$

$$\mathbf{C}_T = \left[egin{array}{ccc} \mathbf{C}_{\mathbf{x}_T,\mathbf{x}_T} & \mathbf{C}_{\mathbf{x}_T,\mathbf{u}_T} \ \mathbf{C}_{\mathbf{u}_T,\mathbf{x}_T} & \mathbf{C}_{\mathbf{u}_T,\mathbf{u}_T} \end{array}
ight] \ \mathbf{c}_T = \left[egin{array}{ccc} \mathbf{c}_{\mathbf{x}_T} \ \mathbf{c}_{\mathbf{u}_T} \end{array}
ight]$$

$$\nabla_{\mathbf{u}_T} Q(\mathbf{x}_T, \mathbf{u}_T) = \mathbf{C}_{\mathbf{u}_T, \mathbf{x}_T} \mathbf{x}_T + \mathbf{C}_{\mathbf{u}_T, \mathbf{u}_T} \mathbf{u}_T + \mathbf{c}_{\mathbf{u}_T}^T = 0 \qquad \mathbf{K}_T = -\mathbf{C}_{\mathbf{u}_T, \mathbf{u}_T}^{-1} \mathbf{C}_{\mathbf{u}_T, \mathbf{x}_T} \mathbf{C}_{\mathbf{u}_T, \mathbf{x}_T} \mathbf{c}_{\mathbf{u}_T, \mathbf{u}_T} \mathbf{c}_{\mathbf{u}_T, \mathbf{u}_T} \mathbf{c}_{\mathbf{u}_T, \mathbf{u}_T} \mathbf{c}_{\mathbf{u}_T, \mathbf{u}_T} \mathbf{c}_{\mathbf{u}_T, \mathbf{u}_T} \mathbf{c}_{\mathbf{u}_T, \mathbf{u}_T} \mathbf{c}_{\mathbf{u}_T} \mathbf{c}_{\mathbf{$$

$$\mathbf{u}_{T} = \mathbf{K}_{T}\mathbf{x}_{T} + \mathbf{k}_{T} \qquad \mathbf{K}_{T} = -\mathbf{C}_{\mathbf{u}_{T},\mathbf{u}_{T}}^{-1}\mathbf{C}_{\mathbf{u}_{T},\mathbf{x}_{T}} \qquad \mathbf{k}_{T} = -\mathbf{C}_{\mathbf{u}_{T},\mathbf{u}_{T}}^{-1}\mathbf{c}_{\mathbf{u}_{T}}$$
$$Q(\mathbf{x}_{T},\mathbf{u}_{T}) = \operatorname{const} + \frac{1}{2} \begin{bmatrix} \mathbf{x}_{T} \\ \mathbf{u}_{T} \end{bmatrix}^{T}\mathbf{C}_{T} \begin{bmatrix} \mathbf{x}_{T} \\ \mathbf{u}_{T} \end{bmatrix} + \begin{bmatrix} \mathbf{x}_{T} \\ \mathbf{u}_{T} \end{bmatrix}^{T}\mathbf{c}_{T}$$

Since \mathbf{u}_T is fully determined by \mathbf{x}_T , we can eliminate it via substitution!

$$V(\mathbf{x}_{T}) = \operatorname{const} + \frac{1}{2} \begin{bmatrix} \mathbf{x}_{T} \\ \mathbf{K}_{T} \mathbf{x}_{T} + \mathbf{k}_{T} \end{bmatrix}^{T} \mathbf{C}_{T} \begin{bmatrix} \mathbf{x}_{T} \\ \mathbf{K}_{T} \mathbf{x}_{T} + \mathbf{k}_{T} \end{bmatrix} + \begin{bmatrix} \mathbf{x}_{T} \\ \mathbf{K}_{T} \mathbf{x}_{T} + \mathbf{k}_{T} \end{bmatrix}^{T} \mathbf{c}_{T}$$

$$V(\mathbf{x}_{T}) = \frac{1}{2} \mathbf{x}_{T}^{T} \mathbf{C}_{\mathbf{x}_{T}, \mathbf{x}_{T}} \mathbf{x}_{T} + \frac{1}{2} \mathbf{x}_{T}^{T} \mathbf{C}_{\mathbf{x}_{T}, \mathbf{u}_{T}} \mathbf{K}_{T} \mathbf{x}_{T} + \frac{1}{2} \mathbf{x}_{T}^{T} \mathbf{K}_{T}^{T} \mathbf{C}_{\mathbf{u}_{T}, \mathbf{x}_{T}} \mathbf{x}_{T} + \frac{1}{2} \mathbf{x}_{T}^{T} \mathbf{C}_{\mathbf{u}_{T}, \mathbf{u}_{T}} \mathbf{K}_{T} \mathbf{x}_{T} + \frac{1}{2} \mathbf{x}_{T}^{T} \mathbf{C}_{\mathbf{u}_{T}, \mathbf{u}_{T}} \mathbf{k}_{T} \mathbf{x}_{T} + \frac{1}{2} \mathbf{x}_{T}^{T} \mathbf{C}_{\mathbf{u}_{T}, \mathbf{u}_{T}} \mathbf{k}_{T} + \frac{1}{2} \mathbf{x}_{T}^{T} \mathbf{C}_{\mathbf{x}_{T}, \mathbf{u}_{T}} \mathbf{k}_{T} + \mathbf{x}_{T}^{T} \mathbf{c}_{\mathbf{x}_{T}} + \mathbf{x}_{T}^{T} \mathbf{c}_{\mathbf{u}_{T}} + \mathbf{x}_{T}^{T} \mathbf{c}_{\mathbf{u}_{T}} + \mathbf{c}_{\mathbf{n}} \mathbf{x}_{T} + \mathbf{c}_{\mathbf{n}} \mathbf{x}_{T} \mathbf{c}_{\mathbf{n}} \mathbf{x}_{T} + \mathbf{c}_{\mathbf{n}} \mathbf{x}_{T} \mathbf{c}_{\mathbf{n}} \mathbf{x}_{T} + \mathbf{c}_{\mathbf{n}} \mathbf{x}_{T} \mathbf{c}_{\mathbf{n}} \mathbf{x}_{T} + \mathbf{c}_{\mathbf{n}} \mathbf{c}_{\mathbf{n}} \mathbf{x}_{T} \mathbf{c}_{\mathbf{n}} \mathbf$$

Solve for \mathbf{u}_{T-1} in terms of \mathbf{x}_{T-1}

 \mathbf{u}_{T-1} affects $\mathbf{x}_T!$

$$f(\mathbf{x}_{T-1}, \mathbf{u}_{T-1}) = \mathbf{x}_{T} = \mathbf{F}_{T-1} \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix} + \mathbf{f}_{T-1}$$

$$Q(\mathbf{x}_{T-1}, \mathbf{u}_{T-1}) = \operatorname{const} + \frac{1}{2} \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix}^{T} \mathbf{C}_{T-1} \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix} + \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix}^{T} \mathbf{c}_{T-1} + V(f(\mathbf{x}_{T-1}, \mathbf{u}_{T-1}))$$

$$V(\mathbf{x}_{T}) = \operatorname{const} + \frac{1}{2} \mathbf{x}_{T}^{T} \mathbf{V}_{T} \mathbf{x}_{T} + \mathbf{x}_{T}^{T} \mathbf{v}_{T}$$

$$V(\mathbf{x}_{T}) = \operatorname{const} + \frac{1}{2} \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix}^{T} \mathbf{F}_{T-1}^{T} \mathbf{V}_{T} \mathbf{F}_{T-1} \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix} + \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix}^{T} \mathbf{F}_{T-1}^{T} \mathbf{V}_{T} \mathbf{f}_{T-1} + \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix}^{T} \mathbf{F}_{T-1}^{T} \mathbf{v}_{T}$$

$$Iinear$$

$$\begin{aligned} Q(\mathbf{x}_{T-1}, \mathbf{u}_{T-1}) &= \mathrm{const} + \frac{1}{2} \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix}^T \mathbf{C}_{T-1} \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix} + \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix}^T \mathbf{c}_{T-1} + V(f(\mathbf{x}_{T-1}, \mathbf{u}_{T-1})) \\ \\ V(\mathbf{x}_T) &= \mathrm{const} + \frac{1}{2} \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix}^T \frac{\mathbf{F}_{T-1}^T \mathbf{V}_T \mathbf{F}_{T-1}}{\mathbf{q} \mathsf{u} \mathsf{d} \mathsf{r} \mathsf{a} \mathsf{t} \mathsf{c}} \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix} + \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix}^T \frac{\mathbf{F}_{T-1}^T \mathbf{V}_T \mathbf{f}_{T-1}}{\mathbf{l} \mathsf{mear}} + \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix}^T \frac{\mathbf{F}_{T-1}^T \mathbf{v}_T}{\mathsf{l} \mathsf{mear}} \\ \\ Q(\mathbf{x}_{T-1}, \mathbf{u}_{T-1}) &= \mathrm{const} + \frac{1}{2} \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix}^T \mathbf{Q}_{T-1} \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix} + \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix}^T \mathbf{q}_{T-1} \\ \\ \mathbf{Q}_{T-1} &= \mathbf{C}_{T-1} + \mathbf{F}_{T-1}^T \mathbf{V}_T \mathbf{F}_{T-1} \\ \\ \mathbf{q}_{T-1} &= \mathbf{c}_{T-1} + \mathbf{F}_{T-1}^T \mathbf{V}_T \mathbf{f}_{T-1} + \mathbf{F}_{T-1}^T \mathbf{v}_T \\ \\ \nabla_{\mathbf{u}_{T-1}} Q(\mathbf{x}_{T-1}, \mathbf{u}_{T-1}) &= \mathbf{Q}_{\mathbf{u}_{T-1}, \mathbf{x}_{T-1}} + \mathbf{Q}_{\mathbf{u}_{T-1}, \mathbf{u}_{T-1}} \mathbf{u}_{T-1} + \mathbf{q}_{\mathbf{u}_{T-1}}^T = 0 \\ \\ \mathbf{u}_{T-1} &= \mathbf{K}_{T-1} \mathbf{x}_{T-1} + \mathbf{k}_{T-1} \\ \end{aligned}$$

Backward recursion

for
$$t = T$$
 to 1:
 $\mathbf{Q}_t = \mathbf{C}_t + \mathbf{F}_t^T \mathbf{V}_{t+1} \mathbf{F}_t$
 $\mathbf{q}_t = \mathbf{c}_t + \mathbf{F}_t^T \mathbf{V}_{t+1} \mathbf{f}_t + \mathbf{F}_t^T \mathbf{v}_{t+1}$
 $Q(\mathbf{x}_t, \mathbf{u}_t) = \operatorname{const} + \frac{1}{2} \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix}^T \mathbf{Q}_t \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix} + \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix}^T \mathbf{u}_t$
 $\mathbf{u}_t \leftarrow \arg\min_{\mathbf{u}_t} Q(\mathbf{x}_t, \mathbf{u}_t) = \mathbf{K}_t \mathbf{x}_t + \mathbf{k}_t$
 $\mathbf{K}_t = -\mathbf{Q}_{\mathbf{u}_t, \mathbf{u}_t}^{-1} \mathbf{Q}_{\mathbf{u}_t, \mathbf{x}_t}$
 $\mathbf{k}_t = -\mathbf{Q}_{\mathbf{u}_t, \mathbf{u}_t}^{-1} \mathbf{Q}_{\mathbf{u}_t}$
 $\mathbf{V}_t = \mathbf{Q}_{\mathbf{x}_t, \mathbf{x}_t} + \mathbf{Q}_{\mathbf{x}_t, \mathbf{u}_t} \mathbf{K}_t + \mathbf{K}_t^T \mathbf{Q}_{\mathbf{u}_t, \mathbf{x}_t} + \mathbf{K}_t^T \mathbf{Q}_{\mathbf{u}_t, \mathbf{u}_t} \mathbf{K}_t$
 $\mathbf{v}_t = \mathbf{q}_{\mathbf{x}_t} + \mathbf{Q}_{\mathbf{x}_t, \mathbf{u}_t} \mathbf{k}_t + \mathbf{K}_t^T \mathbf{Q}_{\mathbf{u}_t} + \mathbf{K}_t^T \mathbf{Q}_{\mathbf{u}_t, \mathbf{u}_t} \mathbf{k}_t$



Some useful definitions

Backward recursion

for t = T to 1: total cost from now until end if we take \mathbf{u}_t from state \mathbf{x}_t $\mathbf{Q}_t = \mathbf{C}_t + \mathbf{F}_t^T \mathbf{V}_{t+1} \mathbf{F}_t$ $\mathbf{q}_t = \mathbf{c}_t + \mathbf{F}_t^T \mathbf{V}_{t+1} \mathbf{f}_t + \mathbf{F}_t^T \mathbf{v}_{t+1}$ $Q(\mathbf{x}_t, \mathbf{u}_t) = \text{const} + \frac{1}{2} \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix}^T \mathbf{Q}_t \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix} + \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix}^T \mathbf{q}_t$ $\mathbf{u}_t \leftarrow \arg\min_{\mathbf{u}_t} Q(\mathbf{x}_t, \mathbf{u}_t) = \mathbf{K}_t \mathbf{x}_t + \mathbf{k}_t$ $\mathbf{K}_t = -\mathbf{Q}_{\mathbf{u}_t,\mathbf{u}_t}^{-1}\mathbf{Q}_{\mathbf{u}_t,\mathbf{x}_t}$ $\mathbf{k}_t = -\mathbf{Q}_{\mathbf{u}_t,\mathbf{u}_t}^{-1}\mathbf{q}_{\mathbf{u}_t}$ $\mathbf{V}_t = \mathbf{Q}_{\mathbf{x}_t,\mathbf{x}_t} + \mathbf{Q}_{\mathbf{x}_t,\mathbf{u}_t}\mathbf{K}_t + \mathbf{K}_t^T\mathbf{Q}_{\mathbf{u}_t,\mathbf{x}_t} + \mathbf{K}_t^T\mathbf{Q}_{\mathbf{u}_t,\mathbf{u}_t}\mathbf{K}_t$ $\mathbf{v}_t = \mathbf{q}_{\mathbf{x}_t} + \mathbf{Q}_{\mathbf{x}_t,\mathbf{u}_t}\mathbf{k}_t + \mathbf{K}_t^T\mathbf{Q}_{\mathbf{u}_t} + \mathbf{K}_t^T\mathbf{Q}_{\mathbf{u}_t,\mathbf{u}_t}\mathbf{k}_t$ $V(\mathbf{x}_t) = \text{const} + \frac{1}{2}\mathbf{x}_t^T \mathbf{V}_t \mathbf{x}_t + \mathbf{x}_t^T \mathbf{v}_t \quad \longleftarrow \begin{array}{l} \text{total cost from now until end from state } \mathbf{x}_t \\ V(\mathbf{x}_t) = \min Q(\mathbf{x}_t, \mathbf{u}_t) \end{array}$

Stochastic dynamics

$$f(\mathbf{x}_t, \mathbf{u}_t) = \mathbf{F}_t \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix} + \mathbf{f}_t$$
$$\mathbf{x}_{t+1} \sim p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t)$$
$$p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t) = \mathcal{N} \left(\mathbf{F}_t \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix} + \mathbf{f}_t, \Sigma_t \right)$$

Solution: choose actions according to $\mathbf{u}_t = \mathbf{K}_t \mathbf{x}_t + \mathbf{k}_t$

 $\mathbf{x}_t \sim p(\mathbf{x}_t)$, no longer deterministic, but $p(\mathbf{x}_t)$ is Gaussian

no change to algorithm! can ignore Σ_t due to symmetry of Gaussians (checking this is left as an exercise; hint: the expectation of a quadratic under a Gaussian has an analytic solution)

Linear-quadratic assumptions:

$$f(\mathbf{x}_t, \mathbf{u}_t) = \mathbf{F}_t \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix} + \mathbf{f}_t$$
$$c(\mathbf{x}_t, \mathbf{u}_t) = \frac{1}{2} \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix}^T \mathbf{C}_t \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix} + \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix}^T \mathbf{c}_t$$

Can we *approximate* a nonlinear system as a linear-quadratic system?

$$f(\mathbf{x}_{t}, \mathbf{u}_{t}) \approx f(\hat{\mathbf{x}}_{t}, \hat{\mathbf{u}}_{t}) + \nabla_{\mathbf{x}_{t}, \mathbf{u}_{t}} f(\hat{\mathbf{x}}_{t}, \hat{\mathbf{u}}_{t}) \begin{bmatrix} \mathbf{x}_{t} - \hat{\mathbf{x}}_{t} \\ \mathbf{u}_{t} - \hat{\mathbf{u}}_{t} \end{bmatrix}$$

$$c(\mathbf{x}_{t}, \mathbf{u}_{t}) \approx c(\hat{\mathbf{x}}_{t}, \hat{\mathbf{u}}_{t}) + \nabla_{\mathbf{x}_{t}, \mathbf{u}_{t}} c(\hat{\mathbf{x}}_{t}, \hat{\mathbf{u}}_{t}) \begin{bmatrix} \mathbf{x}_{t} - \hat{\mathbf{x}}_{t} \\ \mathbf{u}_{t} - \hat{\mathbf{u}}_{t} \end{bmatrix} + \frac{1}{2} \begin{bmatrix} \mathbf{x}_{t} - \hat{\mathbf{x}}_{t} \\ \mathbf{u}_{t} - \hat{\mathbf{u}}_{t} \end{bmatrix}^{T} \nabla_{\mathbf{x}_{t}, \mathbf{u}_{t}}^{2} c(\hat{\mathbf{x}}_{t}, \hat{\mathbf{u}}_{t}) \begin{bmatrix} \mathbf{x}_{t} - \hat{\mathbf{x}}_{t} \\ \mathbf{u}_{t} - \hat{\mathbf{u}}_{t} \end{bmatrix}^{T}$$

$$f(\mathbf{x}_{t}, \mathbf{u}_{t}) \approx f(\hat{\mathbf{x}}_{t}, \hat{\mathbf{u}}_{t}) + \nabla_{\mathbf{x}_{t}, \mathbf{u}_{t}} f(\hat{\mathbf{x}}_{t}, \hat{\mathbf{u}}_{t}) \begin{bmatrix} \mathbf{x}_{t} - \hat{\mathbf{x}}_{t} \\ \mathbf{u}_{t} - \hat{\mathbf{u}}_{t} \end{bmatrix}$$
$$c(\mathbf{x}_{t}, \mathbf{u}_{t}) \approx c(\hat{\mathbf{x}}_{t}, \hat{\mathbf{u}}_{t}) + \nabla_{\mathbf{x}_{t}, \mathbf{u}_{t}} c(\hat{\mathbf{x}}_{t}, \hat{\mathbf{u}}_{t}) \begin{bmatrix} \mathbf{x}_{t} - \hat{\mathbf{x}}_{t} \\ \mathbf{u}_{t} - \hat{\mathbf{u}}_{t} \end{bmatrix} + \frac{1}{2} \begin{bmatrix} \mathbf{x}_{t} - \hat{\mathbf{x}}_{t} \\ \mathbf{u}_{t} - \hat{\mathbf{u}}_{t} \end{bmatrix}^{T} \nabla_{\mathbf{x}_{t}, \mathbf{u}_{t}}^{2} c(\hat{\mathbf{x}}_{t}, \hat{\mathbf{u}}_{t}) \begin{bmatrix} \mathbf{x}_{t} - \hat{\mathbf{x}}_{t} \\ \mathbf{u}_{t} - \hat{\mathbf{u}}_{t} \end{bmatrix}^{T}$$

$$\bar{f}(\delta \mathbf{x}_{t}, \delta \mathbf{u}_{t}) = \mathbf{F}_{t} \begin{bmatrix} \delta \mathbf{x}_{t} \\ \delta \mathbf{u}_{t} \end{bmatrix} \qquad \bar{c}(\delta \mathbf{x}_{t}, \delta \mathbf{u}_{t}) = \frac{1}{2} \begin{bmatrix} \delta \mathbf{x}_{t} \\ \delta \mathbf{u}_{t} \end{bmatrix}^{T} \mathbf{C}_{t} \begin{bmatrix} \delta \mathbf{x}_{t} \\ \delta \mathbf{u}_{t} \end{bmatrix} + \begin{bmatrix} \delta \mathbf{x}_{t} \\ \delta \mathbf{u}_{t} \end{bmatrix}^{T} \mathbf{C}_{t} \\ \nabla_{\mathbf{x}_{t}, \mathbf{u}_{t}} f(\hat{\mathbf{x}}_{t}, \hat{\mathbf{u}}_{t}) \qquad \nabla_{\mathbf{x}_{t}, \mathbf{u}_{t}} c(\hat{\mathbf{x}}_{t}, \hat{\mathbf{u}}_{t})$$

 $\delta \mathbf{x}_t = \mathbf{x}_t - \hat{\mathbf{x}}_t$ $\delta \mathbf{u}_t = \mathbf{u}_t - \hat{\mathbf{u}}_t$ Now we can run LQR with dynamics \bar{f} , cost \bar{c} , state $\delta \mathbf{x}_t$, and action $\delta \mathbf{u}_t$

Iterative LQR (simplified pseudocode)

♦ until convergence:

- $\mathbf{F}_t = \nabla_{\mathbf{x}_t, \mathbf{u}_t} f(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t)$
- $\mathbf{c}_t = \nabla_{\mathbf{x}_t, \mathbf{u}_t} c(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t)$
- $\mathbf{C}_t = \nabla_{\mathbf{x}_t, \mathbf{u}_t}^2 c(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t)$

Run LQR backward pass on state $\delta \mathbf{x}_t = \mathbf{x}_t - \hat{\mathbf{x}}_t$ and action $\delta \mathbf{u}_t = \mathbf{u}_t - \hat{\mathbf{u}}_t$ Run forward pass with real nonlinear dynamics and $\mathbf{u}_t = \mathbf{K}_t(\mathbf{x}_t - \hat{\mathbf{x}}_t) + \mathbf{k}_t$ Update $\hat{\mathbf{x}}_t$ and $\hat{\mathbf{u}}_t$ based on states and actions in forward pass

Why does this work?

Compare to Newton's method for computing $\min_{\mathbf{x}} g(\mathbf{x})$:

until convergence: $\mathbf{g} = \nabla_{\mathbf{x}} g(\hat{\mathbf{x}})$ $\mathbf{H} = \nabla_{\mathbf{x}}^{2} g(\hat{\mathbf{x}})$ $\hat{\mathbf{x}} \leftarrow \arg\min_{\mathbf{x}} \frac{1}{2} (\mathbf{x} - \hat{\mathbf{x}})^{T} \mathbf{H} (\mathbf{x} - \hat{\mathbf{x}}) + \mathbf{g}^{T} (\mathbf{x} - \hat{\mathbf{x}})$

Iterative LQR (iLQR) is the same idea: locally approximate a complex nonlinear function via Taylor expansion

In fact, iLQR is an approximation of Newton's method for solving

 $\min_{\mathbf{u}_1,\ldots,\mathbf{u}_T} c(\mathbf{x}_1,\mathbf{u}_1) + c(f(\mathbf{x}_1,\mathbf{u}_1),\mathbf{u}_2) + \cdots + c(f(f(\ldots)\ldots),\mathbf{u}_T)$

In fact, iLQR is an approximation of Newton's method for solving

 $\min_{\mathbf{u}_1,\ldots,\mathbf{u}_T} c(\mathbf{x}_1,\mathbf{u}_1) + c(f(\mathbf{x}_1,\mathbf{u}_1),\mathbf{u}_2) + \cdots + c(f(f(\ldots)\ldots),\mathbf{u}_T)$

To get Newton's method, need to use *second order* dynamics approximation:

$$f(\mathbf{x}_t, \mathbf{u}_t) \approx f(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) + \nabla_{\mathbf{x}_t, \mathbf{u}_t} f(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) \begin{bmatrix} \delta \mathbf{x}_t \\ \delta \mathbf{u}_t \end{bmatrix} + \frac{1}{2} \left(\nabla_{\mathbf{x}_t, \mathbf{u}_t}^2 f(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) \cdot \begin{bmatrix} \delta \mathbf{x}_t \\ \delta \mathbf{u}_t \end{bmatrix} \right) \begin{bmatrix} \delta \mathbf{x}_t \\ \delta \mathbf{u}_t \end{bmatrix}$$

differential dynamic programming (DDP)

$$\hat{\mathbf{x}} \leftarrow \arg\min_{\mathbf{x}} \frac{1}{2} (\mathbf{x} - \hat{\mathbf{x}})^T \mathbf{H} (\mathbf{x} - \hat{\mathbf{x}}) + \mathbf{g}^T (\mathbf{x} - \hat{\mathbf{x}})$$

why is this a bad idea?

• until convergence:

 $\mathbf{F}_{t} = \nabla_{\mathbf{x}_{t},\mathbf{u}_{t}} f(\hat{\mathbf{x}}_{t},\hat{\mathbf{u}}_{t})$ $\mathbf{c}_{t} = \nabla_{\mathbf{x}_{t},\mathbf{u}_{t}} c(\hat{\mathbf{x}}_{t},\hat{\mathbf{u}}_{t})$ $\mathbf{C}_{t} = \nabla_{\mathbf{x}_{t},\mathbf{u}_{t}}^{2} c(\hat{\mathbf{x}}_{t},\hat{\mathbf{u}}_{t})$



search over α until improvement achieved

Run LQR backward pass on state $\delta \mathbf{x}_t = \mathbf{x}_t - \hat{\mathbf{x}}_t$ and action $\delta \mathbf{u}_t = \mathbf{u}_t - \hat{\mathbf{u}}_t$ Run forward pass with real nonlinear dynamics and $\mathbf{u}_t = \mathbf{K}_t((\mathbf{x}_t - \hat{\mathbf{x}}_t)) + d\mathbf{k}_t$ Update $\hat{\mathbf{x}}_t$ and $\hat{\mathbf{u}}_t$ based on states and actions in forward pass

Additional reading

- 1. Mayne, Jacobson. (1970). Differential dynamic programming.
 - Original differential dynamic programming algorithm.
- 2. Tassa, Erez, Todorov. (2012). Synthesis and Stabilization of Complex Behaviors through Online Trajectory Optimization.
 - Practical guide for implementing non-linear iterative LQR.
- 3. Levine, Abbeel. (2014). Learning Neural Network Policies with Guided Policy Search under Unknown Dynamics.
 - Probabilistic formulation and trust region alternative to deterministic line search.

Case study: nonlinear model-predictive control

Synthesis and Stabilization of Complex Behaviors through Online Trajectory Optimization

Yuval Tassa, Tom Erez and Emanuel Todorov University of Washington

every time step: observe the state \mathbf{x}_t use iLQR to plan $\mathbf{u}_t, \ldots, \mathbf{u}_T$ to minimize $\sum_{t'=t}^{t+T} c(\mathbf{x}_{t'}, \mathbf{u}_{t'})$ execute action \mathbf{u}_t , discard $\mathbf{u}_{t+1}, \ldots, \mathbf{u}_{t+T}$ Synthesis of Complex Behaviors with Online Trajectory Optimization

Yuval Tassa, Tom Erez & Emo Todorov

IEEE International Conference on Intelligent Robots and Systems 2012



discrete planning as a search problem



how to approximate value without full tree?



can't search all paths – where to search first?



intuition: choose nodes with best reward, but also prefer rarely visited nodes

generic MCTS sketch

- ▶ 1. find a leaf s_l using TreePolicy (s_1)
 - 2. evaluate the leaf using $DefaultPolicy(s_l)$
- **3**. update all values in tree between s_1 and s_l take best action from s_1

UCT TreePolicy (s_t)

if s_t not fully expanded, choose new a_t else choose child with best $Score(s_{t+1})$

$$Score(s_t) = \frac{Q(s_t)}{N(s_t)} + 2C\sqrt{\frac{2\ln N(s_{t-1})}{N(s_t)}}$$



Additional reading

- Browne, Powley, Whitehouse, Lucas, Cowling, Rohlfshagen, Tavener, Perez, Samothrakis, Colton. (2012). A Survey of Monte Carlo Tree Search Methods.
 - Survey of MCTS methods and basic summary.

Case study: imitation learning from MCTS

Deep Learning for Real-Time Atari Game Play Using Offline Monte-Carlo Tree Search Planning

Xiaoxiao Guo Computer Science and Eng. University of Michigan guoxiao@umich.edu

Honglak Lee Computer Science and Eng. University of Michigan honglak@umich.edu Satinder Singh Computer Science and Eng. University of Michigan baveja@umich.edu

Richard Lewis Department of Psychology University of Michigan rickl@umich.edu Xiaoshi Wang Computer Science and Eng. University of Michigan xiaoshiw@umich.edu



Case study: imitation learning from MCTS

DAgger

- 1. train $\pi_{\theta}(\mathbf{u}_t | \mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{u}_1, \dots, \mathbf{o}_N, \mathbf{u}_N\}$ 2. run $\pi_{\theta}(\mathbf{u}_t | \mathbf{o}_t)$ to get dataset $\mathcal{D}_{\pi} = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
 - 3. Abbdsernations labest Ateswith Dactising UACTS

4. Aggregate: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_{\pi}$

Why train a policy?

- In this case, MCTS is too slow for real-time play
- Other reasons perception, generalization, etc.: more on this later

What's wrong with known dynamics?





Next time: learning the dynamics model