



Computer Engineering Department

Reinforcement Learning: Model Based RL

Mohammad Hossein Rohban, Ph.D.

Spring 2024

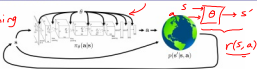
Courtesy: Most of slides are adopted from CS 235 Berkeley.

Overview

- Introduction to model-based reinforcement learning
- What if we know the dynamics? How can we make decisions?
- Stochastic optimization methods
- Monte Carlo tree search (MCTS)
- Trajectory optimization
- Goal: Understand how we can perform planning with known dynamics models in discrete and continuous spaces

Recap: Model-Free RL

Planning

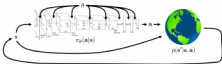


Learning

$$\underbrace{p(s|s_1, a_1, \dots, s_T, a_T)}_{\pi_\theta(r)} = p(s_1) \prod_{t=1}^T \pi_\theta(a_t|s_t) p(s_{t+1}|s_t, a_t)$$

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{s \sim p(s)} \left[\sum_t r(s_t, a_t) \right]$$

Recap: Model-Free RL



$$p_{\theta}(s_1, a_1, \dots, s_T, a_T) = p(s_1) \prod_{t=2}^T \pi_{\theta}(a_t | s_t) \cancel{p(s_t | s_{t-1}, a_{t-1})}$$

assume this is unknown

$$\theta^* = \arg \max_{\theta} E_{r \sim p_{\theta}(r)} \left[\sum_i r(a_i, a_i) \right]$$

What if we knew the transition dynamics?

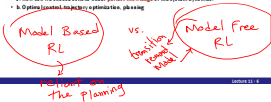
- Often we do know the dynamics
 - Games (e.g., Atari games, chess, Go)
 - Easily modeled systems (e.g., navigating a car)
 - Simulated environments (e.g., simulated robots, video games)
- Often we can learn the dynamics
 - System identification – fit unknown parameters of a known model
 - Learning – fit a general-purpose model to observed transition data

Does knowing the dynamics make things easier?

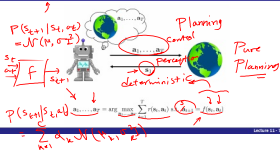
Often, yes!

Model-based RL

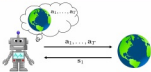
- Model-based reinforcement learning: learn the transition dynamics, then figure out how to choose actions.
- Today: how can we make decisions if we know the dynamics?
 - a. How can we choose actions under **perfect knowledge** of the system dynamics?
 - b. Optimal control, trajectory optimization, planning



The deterministic case



The stochastic open-loop case



$$p(a_1, \dots, a_T | s_1, \dots, s_T) = p(a_1) \prod_{t=1}^T p(a_{t+1} | s_t, a_t)$$

$$a_1, \dots, a_T = \arg \max_{a_1, \dots, a_T} E \left[\sum_{t=1}^T r(s_t, a_t) | a_1, \dots, a_T \right]$$

why is this suboptimal?

The stochastic open-loop case

[illegible]

Journal of Management Inquiry 16(4)

چونکہ انہی پر ہم حالتِ بے اختیار ہے، لہذا انہی کو "گنہگار" سے تعبیر کیا جاتا ہے۔

© 2007 The Authors
Journal compilation © 2007 Blackwell Publishing Ltd

[illegible]

© 2006 The Authors
Journal compilation © 2006 Blackwell Publishing Ltd

www.elsevier.com/locate/jmb

1000

Source: *Author's calculations*.

W. J. G. B. van den Broek

المجلس الأعلى للدراسات والبحوث

Downloaded At: 11:53 11 September 2009

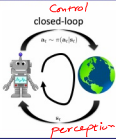
بعد الزمان لم يبق لنا "أولادنا" المصطفى المصطفى

Small family is the primary capital resource for the majority of poor people.

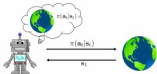
1000

These results indicate that the model is able to capture the main features of the data, and that the model is able to capture the main features of the data.

open-loop vs. closed-loop case



The stochastic open-loop case



$$p(s_1, a_1, \dots, s_T, a_T) = p(s_1) \prod_{t=1}^T \pi(a_t | s_t) p(s_{t+1} | s_t, a_t)$$

$$\pi = \arg \max_{\pi} E_{s \sim p(s)} \left[\sum_i r(s_i, a_i) \right]$$

form of π ?

neural net

global

time-varying linear

$$\mathbf{K}_t \mathbf{s}_t + \mathbf{k}_t$$

local

Stochastic optimization

abstract away optimal control/planning:

$$\mathbf{a}_1, \dots, \mathbf{a}_T = \arg \max_{\mathbf{a}_1, \dots, \mathbf{a}_T} \underbrace{J(\mathbf{a}_1, \dots, \mathbf{a}_T)}_{\text{don't care what this is}} \qquad \mathbf{A} = \arg \max_{\mathbf{A}} J(\mathbf{A})$$

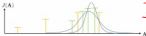
simplest method: guess & check “random shooting method”

1. pick $\mathbf{A}_1, \dots, \mathbf{A}_N$ from some distribution (e.g., uniform)
2. choose \mathbf{A}_i based on $\arg \max_i J(\mathbf{A}_i)$

Cross-entropy Method (CEM)

1. pick $\mathbf{A}_1, \dots, \mathbf{A}_N$ from some distribution (e.g., uniform)

2. choose \mathbf{A}_c based on $\arg \max_i J(\mathbf{A}_i)$ can we do better?



cross-entropy method with continuous-valued inputs:

1. sample $\mathbf{A}_1, \dots, \mathbf{A}_N$ from $p(\mathbf{A})$

2. evaluate $J(\mathbf{A}_1), \dots, J(\mathbf{A}_N)$

3. pick the elites $\mathbf{A}_1, \dots, \mathbf{A}_M$ with the highest values, where $M < N$

4. refit $p(\mathbf{A})$ to the elites $\mathbf{A}_1, \dots, \mathbf{A}_M$

$$\frac{1}{Z} e^{+J(\mathbf{A})}$$

$$(a_1 \dots a_T)$$

$$P(\mathbf{A})$$

$$P(\mathbf{A})$$

$$f(s_{i-1}, a_i)$$

$$\rightarrow J = \sum r(s_i, a_i)$$

$$\begin{bmatrix} J(A_1) \\ \vdots \\ J(A_N) \end{bmatrix} \xrightarrow{CE} \left(\sum \begin{bmatrix} e^{J(A_1)} \\ \vdots \\ e^{J(A_N)} \end{bmatrix} \right),$$

$$P(A) = \sum P(A_i) \quad (e^{J(A_i)})$$

$$= \sum e^{J(A_i)} \log P(A_i)$$

$$\log P(A_i) \propto e^{J(A_i)}$$

$$A_1 \dots A_N \longrightarrow J(A_1) \dots J(A_N)$$

$$\longrightarrow \left[\frac{e^{J(A_1)}}{\sum_i e^{J(A_i)}} \dots \frac{e^{J(A_N)}}{\sum_i e^{J(A_i)}} \right]$$

$$\begin{array}{c} \longrightarrow \\ \text{sample} \\ N \end{array} \quad \longrightarrow \dots$$

Pros and Cons

- Pros

- Could be very fast (Parallelizable)
- Extremely simple

- Cons

- Very harsh dimensionality limit
 - Only open-loop planning
- 

Discrete Case: Monte Carlo Tree Search



Discrete Case: Monte Carlo Tree Search

$$V(s) \leftarrow \left(\frac{N}{N+1} V(s) + \frac{\hat{V}(s)}{N+1} \right)$$



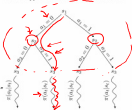
s_1



s_2

e.g., random policy

how to approximate value without full tree?



Discrete Case: Monte Carlo Tree Search

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



s_t



s_t



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

Discrete Case: Monte Carlo Tree Search

generic MCTS sketch

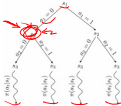
1. find a leaf s_t using $\text{TreePolicy}(s_1)$
 2. evaluate the leaf using $\text{DefaultPolicy}(s_t)$
 3. update all values in tree between s_1 and s_t
- take best action from s_1

UCT $\text{TreePolicy}(s_1)$

- if s_1 not fully expanded, choose new s_2
- else choose child with best $\text{Score}(s_{i+1})$

$$\text{Score}(s_i) = \frac{Q(s_i)}{N(s_i)} + 2C \sqrt{\frac{2 \ln N(s_{i-1})}{N(s_i)}}$$

Handwritten red annotations: A red arrow points to the $Q(s_i)$ term, another points to the $N(s_i)$ term in the denominator, and a third points to the $N(s_i)$ term in the denominator of the square root. A red bracket is under the $N(s_i)$ term in the denominator of the square root. Below the equation, the text $N(s_2)$ is written in red.



Discrete Case: Monte Carlo Tree Search

Algorithm 7 (Monte-Carlo Tree Search)

Input : MDP $M = (S, a_0, A, P_s(s' | s), r(s, a, s'))$, base Q-function Q , time limit T

Output : updated Q-function \hat{Q}

while $\text{currentTime} < T$ do

$\text{selected_node} \leftarrow \text{Select}(a_0)$

$\text{child} \leftarrow \text{Expand}(\text{selected_node})$ — expand and choose a child to simulate

$G \leftarrow \text{Simulate}(\text{child})$ — simulate from child

$\text{Backpropagate}(\text{selected_node}, \text{child}, \hat{Q}, G)$

return \hat{Q}

Discrete Case: Monte Carlo Tree Search

Algorithm 8 (Function – *Select*($s : S$))

Input : state s

Output : unexpanded states

while s is fully expanded do

 Select action a to apply in s using a multi-armed bandit algorithm

 Choose one outcome s' according to $P_s(s' | s)$

$s \leftarrow s'$

return s

Discrete Case: Monte Carlo Tree Search

✱ Algorithm 9 (Function – $\text{Expand}(s \in S)$)

Input : state s

Output : expanded state s'

if s is fully expanded then

 Randomly select action a to apply in s

 Expand one outcome s' according to $P_a(s' \mid s)$ and observe reward r

return s'

Discrete Case: Monte Carlo Tree Search

Algorithm 10 (Function – Backpropagation($s : \mathcal{S}; a : A; Q : \mathcal{S} \times A \rightarrow \mathbb{R}; G : \mathbb{R}$))

Input : state-action pair (s, a) , Q -function Q , rewards G

Output : none

do

$N(s, a) \leftarrow N(s, a) + 1$

$G \leftarrow r + \gamma G$

$Q(s, a) \leftarrow Q(s, a) + \frac{1}{N(s, a)} (G - Q(s, a))$

$s \leftarrow \text{parent of } s$

$a \leftarrow \text{parent action of } s$

while $s \neq s_0$

Discrete Case: Monte Carlo Tree Search



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.

Trajectory Optimization: Can we use derivatives?

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

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

usual story: differentiate via backpropagation and optimize!

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

\mathbf{x}_t - state

\mathbf{u}_t - action



\mathbf{x}_t - state

\mathbf{u}_t - action



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

Shooting methods vs collocation

shooting method: optimize over actions only

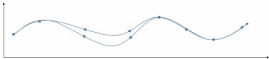
$$\min_{\mathbf{u}_1, \dots, \mathbf{u}_T} c(\mathbf{x}_1, \mathbf{u}_1) + c(f(\mathbf{x}_1, \mathbf{u}_1), \mathbf{u}_2) + \dots + c(f(f(\dots)), \mathbf{u}_T)$$



Shooting methods vs collocation

collocation method: optimize over actions and states, with constraints

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



Linear case: LQR

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

$$\underbrace{f(\mathbf{x}_t, \mathbf{u}_t) = \mathbf{F}_t \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix} + \mathbf{f}_t}_{\text{linear}}$$

$$\underbrace{c(\mathbf{x}_t, \mathbf{u}_t) = \frac{1}{2} \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix}^T \mathbf{C}_t \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix} + \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix}^T \mathbf{c}_t}_{\text{quadratic}}$$

Linear case: LQR

$$\min_{u_1, \dots, u_T} c(x_1, u_1) + c(f(x_1, u_1), u_2) + \dots + \overbrace{c(f(f(\dots), \dots), u_T)}^{x_T \text{ (unknown)}}$$

$$c(x_t, u_t) = \frac{1}{2} \begin{bmatrix} x_t \\ u_t \end{bmatrix}^T \underbrace{C_t}_{\text{only term that depends on } u_T} \begin{bmatrix} x_t \\ u_t \end{bmatrix} + \begin{bmatrix} x_t \\ u_t \end{bmatrix}^T \underbrace{e_t}_{\text{only term that depends on } u_T}$$

$$\underbrace{f(x_t, u_t) = \underbrace{F_t}_{\text{only term that depends on } u_T} \begin{bmatrix} x_t \\ u_t \end{bmatrix} + \underbrace{f_t}_{\text{only term that depends on } u_T}}$$

Base case: solve for u_T only

$$C_T = \begin{bmatrix} C_{x_T, u_T} & C_{u_T, u_T} \end{bmatrix}$$

$$e_T = \begin{bmatrix} e_{x_T} \\ e_{u_T} \end{bmatrix}$$

$$Q(x_T, u_T) = \text{const} + \frac{1}{2} \begin{bmatrix} x_T \\ u_T \end{bmatrix}^T C_T \begin{bmatrix} x_T \\ u_T \end{bmatrix} + \begin{bmatrix} x_T \\ u_T \end{bmatrix}^T e_T$$

$$\nabla_{u_T} Q(x_T, u_T) = C_{u_T, u_T} u_T + e_{u_T}^T = 0$$

$$K_T = -C_{u_T, u_T}^{-1} C_{x_T, u_T}$$

$$u_T = -C_{u_T, u_T}^{-1} (C_{x_T, u_T} x_T + e_{u_T}) \quad u_T = K_T x_T + k_T$$

$$k_T = -C_{u_T, u_T}^{-1} e_{u_T}$$

Linear Case: LQR

$$\mathbf{u}_T = \mathbf{K}_T \mathbf{x}_T + \mathbf{k}_T$$

$$\mathbf{K}_T = -\mathbf{C}_{\mathbf{u}_T, \mathbf{u}_T}^{-1} \mathbf{C}_{\mathbf{u}_T, \mathbf{x}_T}$$

$$\mathbf{k}_T = -\mathbf{C}_{\mathbf{u}_T, \mathbf{u}_T}^{-1} \mathbf{e}_{\mathbf{u}_T}$$

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

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

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

$$V(\mathbf{x}_T) = \frac{1}{2} \mathbf{x}_T^T \mathbf{C}_{\mathbf{x}_T, \mathbf{x}_T} \mathbf{x}_T + \frac{1}{2} \mathbf{x}_T^T \mathbf{C}_{\mathbf{x}_T, \mathbf{u}_T} \mathbf{K}_T \mathbf{x}_T + \frac{1}{2} \mathbf{x}_T^T \mathbf{K}_T^T \mathbf{C}_{\mathbf{u}_T, \mathbf{x}_T} \mathbf{x}_T + \frac{1}{2} \mathbf{x}_T^T \mathbf{K}_T^T \mathbf{C}_{\mathbf{u}_T, \mathbf{u}_T} \mathbf{K}_T \mathbf{x}_T + \\ \mathbf{x}_T^T \mathbf{K}_T^T \mathbf{C}_{\mathbf{u}_T, \mathbf{u}_T} \mathbf{k}_T + \frac{1}{2} \mathbf{x}_T^T \mathbf{C}_{\mathbf{x}_T, \mathbf{u}_T} \mathbf{k}_T + \mathbf{x}_T^T \mathbf{e}_{\mathbf{x}_T} + \mathbf{x}_T^T \mathbf{K}_T^T \mathbf{e}_{\mathbf{u}_T} + \text{const}$$

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

$$\mathbf{V}_T = \mathbf{C}_{\mathbf{x}_T, \mathbf{x}_T} + \mathbf{C}_{\mathbf{x}_T, \mathbf{u}_T} \mathbf{K}_T + \mathbf{K}_T^T \mathbf{C}_{\mathbf{u}_T, \mathbf{x}_T} + \mathbf{K}_T^T \mathbf{C}_{\mathbf{u}_T, \mathbf{u}_T} \mathbf{K}_T$$

$$\mathbf{v}_T = \mathbf{e}_{\mathbf{x}_T} + \mathbf{C}_{\mathbf{x}_T, \mathbf{u}_T} \mathbf{k}_T + \mathbf{K}_T^T \mathbf{e}_{\mathbf{u}_T} + \mathbf{K}_T^T \mathbf{C}_{\mathbf{u}_T, \mathbf{u}_T} \mathbf{k}_T$$

Linear Case: LQR

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

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

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

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

$$\overbrace{V(\mathbf{x}_T) = \text{const} + \frac{1}{2} \mathbf{x}_T^T \mathbf{V}_T \mathbf{x}_T + \mathbf{x}_T^T \mathbf{v}_T}^{\wedge}$$

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

Linear Case: LQR

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

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

$$Q(\mathbf{x}_{T-1}, \mathbf{u}_{T-1}) = \text{const} + \frac{1}{2} \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix}^T \mathbf{Q}_{T-1} \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix} + \begin{bmatrix} \mathbf{x}_{T-1} \\ \mathbf{u}_{T-1} \end{bmatrix}^T \mathbf{q}_{T-1}$$

$$\mathbf{Q}_{T-1} = \mathbf{C}_{T-1} + \mathbf{F}_{T-1}^T \mathbf{V}_T \mathbf{F}_{T-1}$$

$$\mathbf{q}_{T-1} = \mathbf{e}_{T-1} + \mathbf{F}_{T-1}^T \mathbf{V}_T \mathbf{f}_{T-1} + \mathbf{F}_{T-1}^T \mathbf{v}_T$$

$$\nabla_{\mathbf{u}_{T-1}} Q(\mathbf{x}_{T-1}, \mathbf{u}_{T-1}) = \mathbf{Q}_{\mathbf{u}_{T-1}, \mathbf{x}_{T-1}} \mathbf{x}_{T-1} + \mathbf{Q}_{\mathbf{u}_{T-1}, \mathbf{u}_{T-1}} \mathbf{u}_{T-1} + \mathbf{q}_{\mathbf{u}_{T-1}}^T = 0$$

$$\begin{aligned} \mathbf{u}_{T-1} &= \mathbf{K}_{T-1} \mathbf{x}_{T-1} + \mathbf{k}_{T-1} & \mathbf{K}_{T-1} &= -\mathbf{Q}_{\mathbf{u}_{T-1}, \mathbf{u}_{T-1}}^{-1} \mathbf{Q}_{\mathbf{u}_{T-1}, \mathbf{x}_{T-1}} \\ & & \mathbf{k}_{T-1} &= -\mathbf{Q}_{\mathbf{u}_{T-1}, \mathbf{u}_{T-1}}^{-1} \mathbf{q}_{\mathbf{u}_{T-1}} \end{aligned}$$

Linear Case: LQR

Backward recursion

for $t = T$ to 1:

$$Q_t = C_t^T V_{t+1} C_t$$

$$q_t = a_t + F_t^T V_{t+1} b_t + F_t^T v_{t+1}$$

$$Q(x_t, u_t) = \text{const} + \frac{1}{2} \begin{bmatrix} x_t \\ u_t \end{bmatrix}^T Q_t \begin{bmatrix} x_t \\ u_t \end{bmatrix} + \begin{bmatrix} x_t \\ u_t \end{bmatrix}^T q_t$$

$$u_t = \arg \min_{u_t} Q(x_t, u_t) = K_t x_t + k_t$$

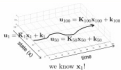
$$K_t = -Q_{u_t, u_t}^{-1} Q_{u_t, x_t}$$

$$k_t = -Q_{u_t, u_t}^{-1} q_{u_t}$$

$$V_t = Q_{x_t, x_t} + Q_{x_t, u_t} K_t + K_t^T Q_{u_t, x_t} + K_t^T Q_{u_t, u_t} K_t$$

$$v_t = q_{x_t} + Q_{x_t, u_t} k_t + K_t^T q_{u_t} + K_t^T Q_{u_t, u_t} k_t$$

$$V(x_t) = \text{const} + \frac{1}{2} x_t^T V_t x_t + x_t^T v_t$$



Forward recursion

for $t = 1$ to T :

$$u_t = K_t x_t + k_t$$

$$x_{t+1} = f(x_t, u_t)$$

LQR for Stochastic Dynamics

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

$$\mathbf{x}_{t+1} \sim p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t)$$

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

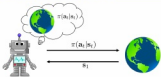
Solution: choose actions according to $\mathbf{u}_t = \mathbf{K}_t \mathbf{x}_t + \mathbf{k}_t$

$\mathbf{x}_t \sim p(\mathbf{x}_t)$, 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)

The stochastic closed-loop case



form of π ?

$$p(s_1, a_1, \dots, s_T, a_T) = p(s_1) \prod_{t=1}^T \pi(a_t | s_t) p(s_{t+1} | s_t, a_t)$$

time-varying linear

$$K_t s_t + k_t$$

$$\pi = \arg \max_{\pi} E_{s \sim p(s)} \left[\sum_t r(s_t, a_t) \right]$$

Nonlinear case: DDP/iterative LQR

Linear-quadratic assumptions:

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

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

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

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

Nonlinear case: DDP/iterative LQR

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

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

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

$$\delta \mathbf{x}_t = \mathbf{x}_t - \hat{\mathbf{x}}_t$$

$$\delta \mathbf{u}_t = \mathbf{u}_t - \hat{\mathbf{u}}_t$$

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

Nonlinear case: DDP/iterative LQR

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 $\mathbf{u}_t = \mathbf{K}_t(\mathbf{x}_t - \hat{\mathbf{x}}_t) + \mathbf{k}_t + \hat{\mathbf{u}}_t$

Update $\hat{\mathbf{x}}_t$ and $\hat{\mathbf{u}}_t$ based on states and actions in forward pass

Nonlinear case: DDP/iterative LQR

Why does this work?

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



until convergence:

$$\begin{aligned}\mathbf{g} &= \nabla_{\mathbf{x}} g(\hat{\mathbf{x}}) \\ \mathbf{H} &= \nabla_{\mathbf{x}}^2 g(\hat{\mathbf{x}}) \\ \hat{\mathbf{x}} &\leftarrow \arg \min_{\mathbf{x}} \frac{1}{2}(\mathbf{x} - \hat{\mathbf{x}})^T \mathbf{H}(\mathbf{x} - \hat{\mathbf{x}}) + \mathbf{g}^T(\mathbf{x} - \hat{\mathbf{x}})\end{aligned}$$

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

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

Nonlinear case: DDP/iterative LQR

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

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

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

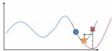
$$f(\mathbf{x}_t, \mathbf{u}_t) \approx f(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) + \nabla_{\mathbf{x}, \mathbf{u}} f(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) \begin{bmatrix} \delta \mathbf{x}_t \\ \delta \mathbf{u}_t \end{bmatrix} + \frac{1}{2} \left(\nabla_{\mathbf{x}, \mathbf{u}}^2 f(\hat{\mathbf{x}}_t,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

differential dynamic programming (DDP)

Nonlinear case: DDP/iterative LQR

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

why is this a bad idea?



until convergence:

$$\mathbf{F}_t = \nabla_{\mathbf{u}_t, \mathbf{u}_t} f(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t)$$

$$\mathbf{c}_t = \nabla_{\mathbf{x}_t, \mathbf{x}_t} c(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t)$$

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

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 $\mathbf{u}_t = \mathbf{K}_t(\mathbf{x}_t - \hat{\mathbf{x}}_t) + \hat{\mathbf{u}}_t$

Update $\hat{\mathbf{x}}_t$ and $\hat{\mathbf{u}}_t$ based on states and actions in forward pass.

Additional Reading

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