

LQR, iLQR, MPC

Reinforcement Learning

Roberto Capobianco



SAPIENZA
UNIVERSITÀ DI ROMA

Recap

Policy Iteration

— — —

- Outputs policies at every iteration: $\{\pi_0, \pi_1, \pi_2 \dots \pi_T\}$
- Different from Value Iteration that was outputting values

Procedure:

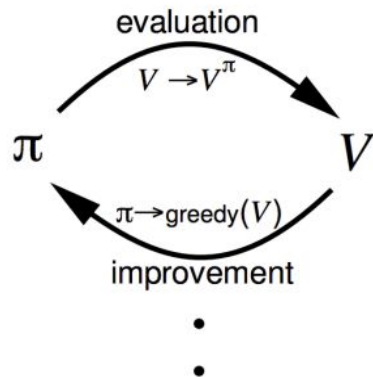
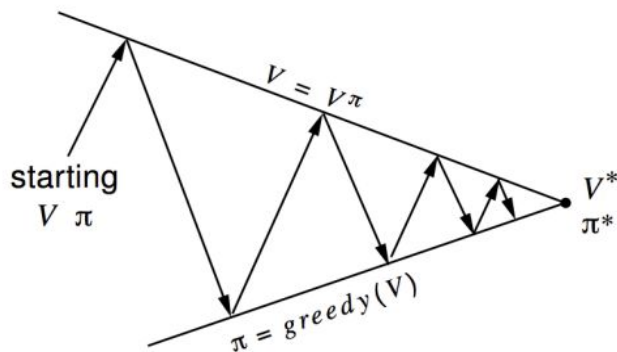
1. Start with a random guess π_0 (can be deterministic or stochastic)
2. For $t=0, \dots, T$:
 - a. Do **policy evaluation** and compute Q^{π^t} for all s, a
 - b. Do **policy improvement** as $\pi_{t+1} = \operatorname{argmax}_a Q^{\pi^t}(s, a)$ for all s

This algorithm only makes progress, and the performance progress of the policy is monotonic



Properties of Policy Iteration

- Monotonic improvement: $Q^{\pi^{t+1}} \geq Q^{\pi^t}$ for all s, a
- Convergence: $\|V^{\pi^i} - V^*\| \leq \gamma^{i+1} \|V^{\pi^0} - V^*\|$



Credits: David Silver



SAPIENZA
UNIVERSITÀ DI ROMA

Properties of Policy Iteration

- Monotonic improvement: $Q^{\pi^{t+1}} \geq Q^{\pi^t}$ for all s, a
- Convergence: $\|V^{\pi^i} - V^*\| \leq \gamma^{i+1} \|V^{\pi^0} - V^*\|$

Is there a max number of iterations of policy iteration?

$|A|^{|S|}$ since that is the maximum number of policies, and as the policy improvement step is monotonically improving, each policy can only appear in one round of policy iteration unless it is an optimal policy



Properties of Policy Iteration

— — —

- Monotonic improvement: $Q^{\pi^{t+1}} \geq Q^{\pi^t}$ for all s, a
- Convergence: $\|V^{\pi^i} - V^*\| \leq \gamma^{i+1} \|V^{\pi^0} - V^*\|$

When do we stop?

if the policy does not change anymore for any state



We Did Dynamic Programming!

Dynamic Programming can be applied if we have:

- *Optimal substructure*: Optimality exists and the optimal solution can be decomposed into subproblems
- *Overlapping subproblems*: Subproblems recur many times and the solutions can be cached and reused

MDPs satisfy both properties: thanks Bellman equation!



We Did Dynamic Programming!

We applied dynamic programming for **planning** as we assumed to know the MDP transition probabilities

| Problem | Bellman Equation | Algorithm |
|------------|--|-----------------------------|
| Prediction | Bellman Expectation Equation | Iterative Policy Evaluation |
| Control | Bellman Expectation Equation + Greedy Policy Improvement | Policy Iteration |
| Control | Bellman Optimality Equation | Value Iteration |

Credits: David Silver



SAPIENZA
UNIVERSITÀ DI ROMA

Primal Linear Program

As an alternative to VI and PI

Consider the Bellman optimality equation

$$V(s) = \max_a \{r_t + \gamma \mathbb{E}_{s' \sim p(\cdot | s, \pi(s))} [V(s')]\}$$

and write it as a linear program:

$$\min V(s)$$

such that $V(s) \geq r_t + \gamma \mathbb{E}_{s' \sim p(\cdot | s, \pi(s))} [V(s')]$ for all s, a



End Recap



SAPIENZA
UNIVERSITÀ DI ROMA

Finite-Horizon MDPs

Slightly different formulation:

$$(S, A, R, T, \mathbf{H}, \mu_0)$$

(time-horizon) $\mathbf{H} \geq 0$ and $\mathbf{s}_0 \sim \mu_0$ (initial state distribution)



Finite-Horizon MDPs

Slightly different formulation:

$$(S, A, R, T, \mathbf{H}, \mu_0)$$

(time-horizon) $\mathbf{H} \geq 0$ and $\mathbf{s}_0 \sim \mu_0$ (initial state distribution)

We consider time-dependent policies π

$$\pi = \{\pi_0, \pi_1, \pi_2 \dots \pi_{H-1}\}$$



Finite-Horizon MDPs

Slightly different formulation:

$$(S, A, R, T, \mathbf{H}, \mu_0)$$

(time-horizon) $\mathbf{H} \geq 0$ and $\mathbf{s}_0 \sim \mu_0$ (initial state distribution)

We consider time-dependent policies π

$$\pi = \{\pi_0, \pi_1, \pi_2 \dots \pi_{H-1}\}$$

Actions might be different for the same state depending on t



Finite-Horizon MDPs

Slightly different formulation:

$$(S, A, R, T, \mathbf{H}, \mu_0)$$

(time-horizon) $\mathbf{H} \geq 0$ and $\mathbf{s}_0 \sim \mu_0$ (initial state distribution)

We consider time-dependent policies π

$$\pi = \{\pi_0, \pi_1, \pi_2 \dots \pi_{H-1}\}$$

e.g., I could explore more at the beginning



Finite-Horizon MDPs

Slightly different formulation:

$$(S, A, R, T, \mathbf{H}, \mu_0)$$

(time-horizon) $\mathbf{H} \geq 0$ and $\mathbf{s}_0 \sim \mu_0$ (initial state distribution)

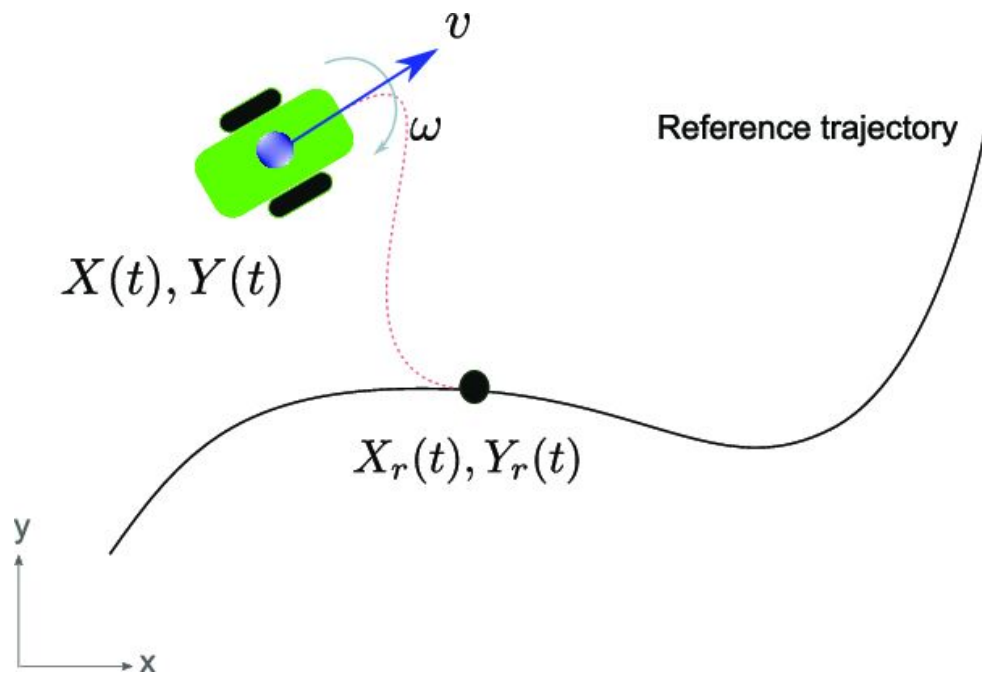
We consider time-dependent policies π

$$\pi = \{\pi_0, \pi_1, \pi_2 \dots \pi_{H-1}\}$$

Very common in control!



Finite-Horizon MDPs - Example



Finite-Horizon MDP & Policy Interaction

— — —

$$\text{MDP} = (S, A, R, T, H, \mu_0)$$

$$\pi = \{\pi_0, \pi_1, \pi_2 \dots \pi_{H-1}\}$$

$$s_0 \sim \mu_0$$

$$(s_0)$$



Finite-Horizon MDP & Policy Interaction

$$\text{MDP} = (S, A, R, T, H, \mu_0)$$

$$\pi = \{\pi_0, \pi_1, \pi_2 \dots \pi_{H-1}\}$$

$$s_0 \sim \mu_0$$

$$a_0 = \pi_0(s_0)$$

$$(s_0, a_0)$$



Finite-Horizon MDP & Policy Interaction

— — —

$$\text{MDP} = (S, A, R, T, H, \mu_0)$$

$$\pi = \{\pi_0, \pi_1, \pi_2 \dots \pi_{H-1}\}$$

$$s_1 \sim P(\cdot | s_0, a_0)$$

$$(s_0, a_0, s_1)$$



Finite-Horizon MDP & Policy Interaction

— — —

$$\text{MDP} = (S, A, R, T, H, \mu_0)$$

$$\pi = \{\pi_0, \pi_1, \pi_2 \dots \pi_{H-1}\}$$

$$s_1 \sim P(\cdot | s_0, a_0)$$

$$a_1 = \pi_1(s_1)$$

$$(s_0, a_0, s_1, a_1)$$



Finite-Horizon MDP & Policy Interaction

— — —

$$\text{MDP} = (S, A, R, T, H, \mu_0)$$

$$\pi = \{\pi_0, \pi_1, \pi_2 \dots \pi_{H-1}\}$$

$$(s_0, a_0, s_1, a_1, \dots s_{H-1}, a_{H-1})$$



Finite-Horizon MDP: V & Q

$$V_h^\pi(s) = \mathbb{E}_\pi[\sum_{\square=h}^{H-1} r(s_\square, a_\square)]$$

where $s_h=s$, $a_\square=\pi_\square(s_\square)$ and $s_{\square+1}\sim P(\cdot|s_\square, a_\square)$

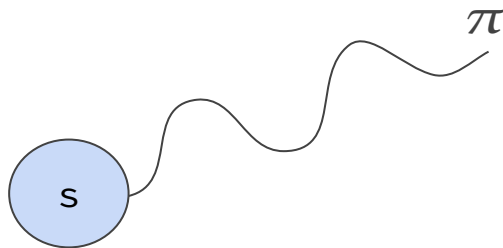
No discount factor!



Finite-Horizon MDP: V & Q

$$V_h^\pi(s) = \mathbb{E}_\pi[\sum_{\square=h}^{H-1} r(s_\square, a_\square)]$$

where $s_h=s$, $a_\square=\pi_\square(s_\square)$ and $s_{\square+1}\sim P(\cdot|s_\square, a_\square)$



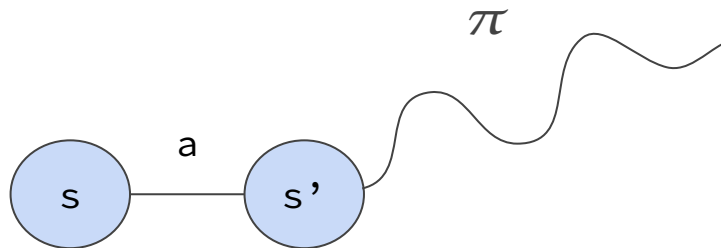
Finite-Horizon MDP: V & Q

$$V_h^\pi(s) = \mathbb{E}_\pi[\sum_{\square=h}^{H-1} r(s_\square, a_\square)]$$

where $s_h=s$, $a_\square=\pi_\square(s_\square)$ and $s_{\square+1}\sim P(\cdot|s_\square, a_\square)$

$$Q_h^\pi(s, a) = \mathbb{E}_\pi[\sum_{\square=h}^{H-1} r(s_\square, a_\square)]$$

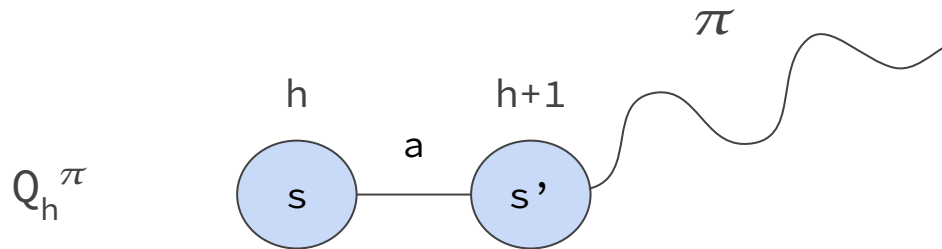
where $s_h=s$, $a_h=a$, $a_\square=\pi_\square(s_\square)$ and $s_{\square+1}\sim P(\cdot|s_\square, a_\square)$



Finite-Horizon MDP: Bellman Equation

— — —

$$Q_h^\pi(s, a) = r(s, a) + \mathbb{E}_{s' \sim p(\cdot | s, a)} [V_{h+1}^\pi(s')]$$



Finding the Optimal Policy

— — —

$$\pi^* = \{\pi_0^*, \pi_1^*, \pi_2^* \dots \pi_{H-1}^*\}$$

Easier problem than infinite horizon!



Finding the Optimal Policy

$$\pi^* = \{\pi_0^*, \pi_1^*, \pi_2^* \dots \pi_{H-1}^*\}$$

Let's reason backwards in time and apply dynamic programming:

$$Q_{H-1}^*(s, a)?$$



Finding the Optimal Policy

$$\pi^* = \{\pi_0^*, \pi_1^*, \pi_2^* \dots \pi_{H-1}^*\}$$

Let's reason backwards in time and apply dynamic programming:

$$Q_{H-1}^*(s,a) = r(s,a)$$



Finding the Optimal Policy

$$\pi^* = \{\pi_0^*, \pi_1^*, \pi_2^* \dots \pi_{H-1}^*\}$$

Let's reason backwards in time and apply dynamic programming:

$$Q_{H-1}^*(s, a) = r(s, a)$$

$$\pi_{H-1}^*(s) = \operatorname{argmax}_a Q_{H-1}^*(s, a)$$

$$V_{H-1}^*(s) = \max_a Q_{H-1}^*(s, a) = Q_{H-1}^*(s, \pi_{H-1}^*(s))$$



Finding the Optimal Policy

— — —

$$Q_{H-1}^*(s, a) = r(s, a)$$

$$\pi_{H-1}^*(s) = \operatorname{argmax}_a Q_{H-1}^*(s, a)$$

$$V_{H-1}^*(s) = \max_a Q_{H-1}^*(s, a) = Q_{H-1}^*(s, \pi_{H-1}^*(s))$$

Now we can reason about H-2!



Finding the Optimal Policy

$$Q_{H-1}^*(s, a) = r(s, a)$$

$$\pi_{H-1}^*(s) = \operatorname{argmax}_a Q_{H-1}^*(s, a)$$

$$V_{H-1}^*(s) = \max_a Q_{H-1}^*(s, a) = Q_{H-1}^*(s, \pi_{H-1}^*(s))$$

Now we can reason about H-2!

$$\text{Bellman Equation: } Q_h^\pi(s, a) = r(s, a) + \mathbb{E}_{s' \sim p(\cdot | s, a)} [V_{h+1}^\pi(s')]$$



Finding the Optimal Policy

$$Q_{H-1}^*(s,a) = r(s,a)$$

$$\pi_{H-1}^*(s) = \operatorname{argmax}_a Q_{H-1}^*(s,a)$$

$$V_{H-1}^*(s) = \max_a Q_{H-1}^*(s,a) = Q_{H-1}^*(s, \pi_{H-1}^*(s))$$

Now we can reason about H-2!

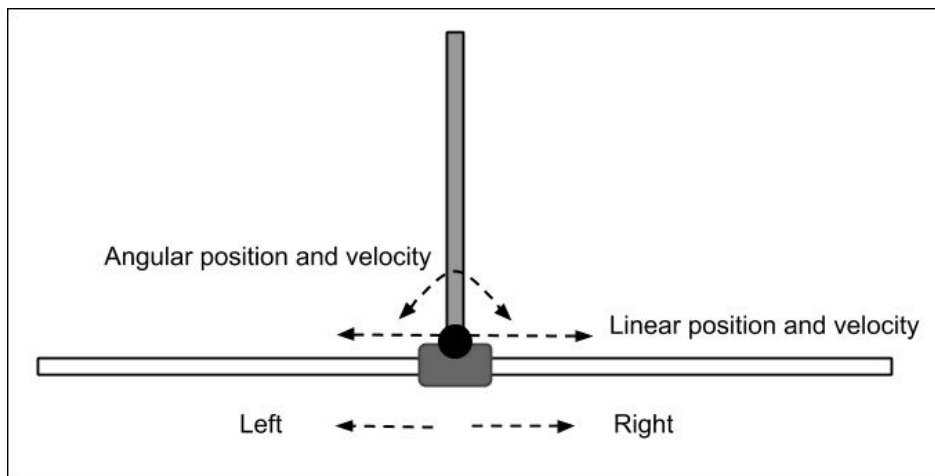
Bellman Equation: $Q_h^\pi(s,a) = r(s,a) + \mathbb{E}_{s' \sim p(\cdot|s,a)} [V_{h+1}^\pi(s')]$

$$\pi_h^*(s) = \operatorname{argmax}_a Q_h^*(s,a)$$



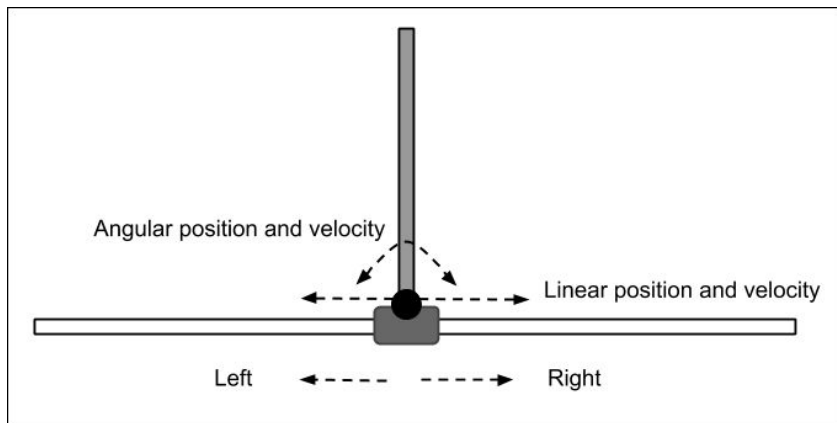
Control Problems

So far, we assumed discrete state and action spaces, but what about cartpole:



Control Problems

So far, we assumed discrete state and action spaces, but what about cartpole:



- **state:** angular pos & vel, linear pos & vel
- **action/control:** force applied on the cart
- **goal:** find the control policy which minimizes the long term cost c



Control Problems

— — —

More in general in control problems we have x in \mathbb{R}^d and u in \mathbb{R}^k

We denote the state as x and the action as u , as it's typical for control problems



Control Problems

— — —

More in general in control problems we have x in \mathbb{R}^d and u in \mathbb{R}^k

We also talk about cost instead of reward: we want to minimize the cost instead of maximizing the reward!



Optimal Control

Given a dynamical system with a non-linear transition function f , state x in \mathbb{R}^d and control u in \mathbb{R}^k , we want to find a control policy π such that

$$\text{minimize } \mathbb{E}_{\pi}[c_H(x_H) + \sum_{h=0}^{H-1} c_h(x_h, u_h)]$$

$$\text{where } u_h = \pi(x_h) \text{ and } x_0 \sim \mu_0$$



Optimal Control

Given a dynamical system with a non-linear transition function f , state x in \mathbb{R}^d and control u in \mathbb{R}^k , we want to find a control policy π such that

$$\text{minimize } \mathbb{E}_{\pi}[c_H(x_H) + \sum_{h=0}^{H-1} c_h(x_h, u_h)]$$

$$\text{where } u_h = \pi(x_h) \text{ and } x_0 \sim \mu_0$$

Now this seems very familiar! Can we treat it as a Finite-Horizon MDP and use value iteration?



Optimal Control

Given a dynamical system with a non-linear transition function f , state x in \mathbb{R}^d and control u in \mathbb{R}^k , we want to find a control policy π such that

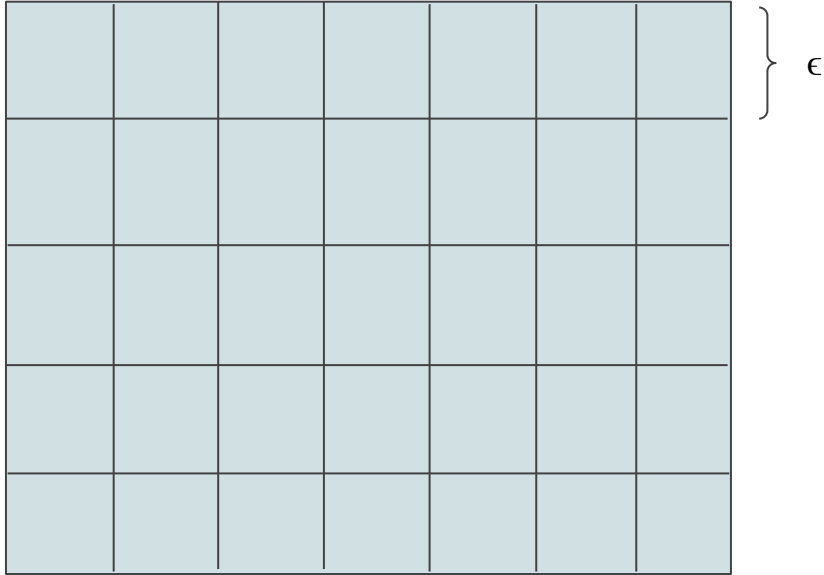
$$\text{minimize } \mathbb{E}_{\pi}[c_H(x_H) + \sum_{h=0}^{H-1} c_h(x_h, u_h)]$$

$$\text{where } u_h = \pi(x_h) \text{ and } x_0 \sim \mu_0$$

Now this seems very familiar! Can we treat it as a Finite-Horizon MDP and use value iteration? **YES if we can discretize**



Discretization



x

x in \mathbb{R}^d and u in \mathbb{R}^k

Number of total points on the discretized grid increases exponentially and becomes

$$(1/\epsilon)^d + (1/\epsilon)^k$$



SAPIENZA
UNIVERSITÀ DI ROMA

Bellman's Curse of Dimensionality

- n -dimensional (discrete) state space
- The number of states grows exponentially in n

In practice discretization is useful, but it is only computationally feasible up to 5 or 6 dimensional state spaces



Bellman's Curse of Dimensionality

- n -dimensional (discrete) state space
- The number of states grows exponentially in n

In practice discretization is useful, but it is only computationally feasible up to 5 or 6 dimensional state spaces

Let's try to work directly in continuous space, starting from simplified problems



Linear Systems

Consider a system of this kind:

$$x_{t+1} = Ax_t + Bu_t$$

- x_t state at time t
- u_t control (i.e., action) at time t

A in $\mathbb{R}^{d \times d}$, B in $\mathbb{R}^{d \times k}$



Linear Systems

Consider a system of this kind:

$$x_{t+1} = Ax_t + Bu_t$$

This is our
transition function!

- x_t state at time t
- u_t control (i.e., action) at time t

A in $\mathbb{R}^{d \times d}$, B in $\mathbb{R}^{d \times k}$



Quadratic Cost Function

Consider a cost function of this kind

$$c(x_t, u_t) = x_t^T Q x_t + u_t^T R u_t$$

alternative notation
 $g(x_t, u_t)$

- Q in $\mathbb{R}^{d \times d}$ and R in $\mathbb{R}^{k \times k}$ square matrices
- Q and R positive definite

As a result, there is a non-zero cost for any non-zero state with all-zero control



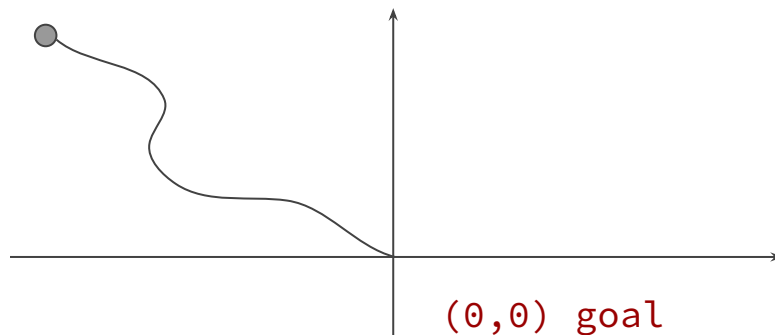
Quadratic Cost Function

— — —

Consider a cost function of this kind

$$c(x_t, u_t) = x_t^T Q x_t + u_t^T R u_t$$

As a result, there is a non-zero cost for any non-zero state with all non-zero control



Q Functions for Linear Systems and Quadratic Cost

— — —

$$Q_h^*(x, u) = c(x, u) + \mathbb{E}_{x' \sim p(\cdot | x, u)} [V_{h+1}^*(x')]$$



Q Functions for Linear Systems and Quadratic Cost

— — —

$$Q_h^*(x, u) = c(x, u) + \mathbb{E}_{x' \sim p(\cdot | x, u)} [V_{h+1}^*(x')]$$

$$c(x_t, u_t) = x_t^T Q x_t + u_t^T R u_t$$



Q Functions for Linear Systems and Quadratic Cost

— — —

$$Q_h^*(x, u) = c(x, u) + \mathbb{E}_{x' \sim p(\cdot | x, u)} [V_{h+1}^*(x')] =$$
$$x_t^T Q x_t + u_t^T R u_t + \mathbb{E}_{x' \sim p(\cdot | x, u)} [V_{h+1}^*(x')]$$



Q Functions for Linear Systems and Quadratic Cost

$$Q_h^*(x, u) = c(x, u) + \mathbb{E}_{x' \sim p(\cdot | x, u)} [V_{h+1}^*(x')] = \\ x_t^T Q x_t + u_t^T R u_t + \mathbb{E}_{x' \sim p(\cdot | x, u)} [V_{h+1}^*(x')]$$

$$x_{t+1} = Ax_t + Bu_t$$

(deterministic for now)



Q Functions for Linear Systems and Quadratic Cost

— — —

$$\begin{aligned} Q_h^*(x, u) &= c(x, u) + \mathbb{E}_{x' \sim p(\cdot | x, u)} [V_{h+1}^*(x')] = \\ &= x_t^T Q x_t + u_t^T R u_t + \mathbb{E}_{x' \sim p(\cdot | x, u)} [V_{h+1}^*(x')] = \\ &= x_t^T Q x_t + u_t^T R u_t + V_{h+1}^*(Ax_t + Bu_t) \end{aligned}$$



Inductive Step: Assumption

— — —

$$\begin{aligned} Q_h^*(x, u) &= c(x, u) + \mathbb{E}_{x' \sim p(\cdot | x, u)} [V_{h+1}^*(x')] = \\ &= x_t^T Q x_t + u_t^T R u_t + \mathbb{E}_{x' \sim p(\cdot | x, u)} [V_{h+1}^*(x')] = \\ &= x_t^T Q x_t + u_t^T R u_t + V_{h+1}^*(Ax_t + Bu_t) \end{aligned}$$

Assume $V_{h+1}^*(x_t) = x_t^T P_{h+1} x_t$ with P in $\mathbb{R}^{d \times d}$



Inductive Step: Question

— — —

$$\begin{aligned} Q_h^*(x, u) &= c(x, u) + \mathbb{E}_{x' \sim p(\cdot | x, u)} [V_{h+1}^*(x')] = \\ &= x_t^T Q x_t + u_t^T R u_t + \mathbb{E}_{x' \sim p(\cdot | x, u)} [V_{h+1}^*(x')] = \\ &= x_t^T Q x_t + u_t^T R u_t + V_{h+1}^*(Ax_t + Bu_t) \end{aligned}$$

Can we show that also $V_h^*(x_t) = x_t^T P_h x_t$ with P in $\mathbb{R}^{d \times d}$?



Inductive Step: Implement the Assumption

— — —

$$\begin{aligned} Q_h^*(x, u) &= c(x, u) + \mathbb{E}_{x' \sim p(\cdot | x, u)} [V_{h+1}^*(x')] = \\ &x_t^T Q x_t + u_t^T R u_t + \mathbb{E}_{x' \sim p(\cdot | x, u)} [V_{h+1}^*(x')] = \\ &x_t^T Q x_t + u_t^T R u_t + V_{h+1}^*(Ax_t + Bu_t) = \\ &x_t^T Q x_t + u_t^T R u_t + (Ax_t + Bu_t)^T P_{h+1} (Ax_t + Bu_t) \end{aligned}$$



Finding the Best Action

We are now interested in finding the $\pi_h^*(x) = \operatorname{argmin}_u Q_h^*(x, u)$

$$Q_h^*(x_t, u_t) = x_t^T Q x_t + u_t^T R u_t + (A x_t + B u_t)^T P_{h+1} (A x_t + B u_t)$$



Finding the Best Action

We are now interested in finding the $\pi_h^*(x) = \operatorname{argmin}_u Q_h^*(x, u)$

$$Q_h^*(x_t, u_t) = x_t^T Q x_t + u_t^T R u_t + (A x_t + B u_t)^T P_{h+1} (A x_t + B u_t)$$

How do you find minima?



Finding the Best Action

We are now interested in finding the $\pi_h^*(x) = \operatorname{argmin}_u Q_h^*(x, u)$

$$Q_h^*(x_t, u_t) = x_t^T Q x_t + u_t^T R u_t + (A x_t + B u_t)^T P_{h+1} (A x_t + B u_t)$$

Set $\nabla_u Q_h^*(x, u) = 0$ and solve for u



Inductive Step: Proof

— — —

Assume $V_{h+1}^* = x_t^T P_{h+1} x_t^T$, can we show that also $V_h^*(x_t) = x_t^T P_h x_t$

Inductive Step: Proof

— — —

Assume $V_{h+1}^* = x_t^T P_{h+1} x_t$, can we show that also $V_h^*(x_t) = x_t^T P_h x_t$

1. Set $Q_h^*(x_t, u_t) = x_t^T Q x_t + u_t^T R u_t + (A x_t + B u_t)^T P_{h+1} (A x_t + B u_t)$



Inductive Step: Proof

— — —

Assume $V_{h+1}^* = x_t^T P_{h+1} x_t$, can we show that also $V_h^*(x_t) = x_t^T P_h x_t$

1. Set $Q_h^*(x_t, u_t) = x_t^T Q x_t + u_t^T R u_t + (A x_t + B u_t)^T P_{h+1} (A x_t + B u_t)$
2. $\nabla_u Q_h^*(x_t, u_t) = 2R u_t + 2B^T P_{h+1} (A x_t + B u_t) = 0$



Inductive Step: Proof

Assume $V_{h+1}^* = x_t^T P_{h+1} x_t$, can we show that also $V_h^*(x_t) = x_t^T P_h x_t$

1. Set $Q_h^*(x_t, u_t) = x_t^T Q x_t + u_t^T R u_t + (A x_t + B u_t)^T P_{h+1} (A x_t + B u_t)$
2. $\nabla_u Q_h^*(x_t, u_t) = 2R u_t + 2B^T P_{h+1} (A x_t + B u_t) = 0$
3. Solve for $u = \pi_h^*(x_t) = -(R + B^T P_h B)^{-1} B^T P_{h+1} A x_t = K_h^* x_t$



Inductive Step: Proof

Assume $V_{h+1}^* = x_t^T P_{h+1} x_t$, can we show that also $V_h^*(x_t) = x_t^T P_h x_t$

- $u = \pi_h^*(x_t) = -(R+B^T P_h B)^{-1} B^T P_h A x_t = K_h^* x_t$
- We know that $V_h^*(x) = \max_u Q_h^*(x, u) = Q_h^*(x, \pi_h^*(x)) = Q_h^*(x, K_h^* x)$

We can replace it in

$$Q_h^* = x_t^T Q x_t + x_t^T K_h^{*T} R K_h^* x_t + (A x_t + B K_h^* x_t)^T P_{h+1} (A x_t + B K_h^* x_t)$$



Inductive Step: Proof

Assume $V_{h+1}^* = x_t^T P_{h+1} x_t$, can we show that also $V_h^*(x_t) = x_t^T P_h x_t$

- $u = \pi_h^*(x_t) = -(R + B^T P_h B)^{-1} B^T P_h A x_t = K_h^* x_t$
- We know that $V_h^*(x) = \max_u Q_h^*(x, u) = Q_h^*(x, \pi_h^*(x)) = Q_h^*(x, K_h^* x)$

We can replace it in

$$\begin{aligned} Q_h^* &= x_t^T Q x_t + x_t^T K_h^{*T} R K_h^* x_t + (A x_t + B K_h^* x_t)^T P_{h+1} (A x_t + B K_h^* x_t) = \\ x_t^T (Q + K_h^{*T} R K_h^* + (A + B K_h^*)^T P_{h+1} (A + B K_h^*)) x_t &= x_t^T P_h x_t \end{aligned}$$



Cost at H

Recall our goal is to minimize $\mathbb{E}_{\pi}[c_H(x_H) + \sum_{h=0}^{H-1} c_h(x_h, u_h)]$

So, at H, we only have $c_H(x_H)$. This is typically 0 or $x_H^T Q x_H$



Cost at H

Recall our goal is to minimize $\mathbb{E}_{\pi}[c_H(x_H) + \sum_{h=0}^{H-1} c_h(x_h, u_h)]$

So, at H, we only have $c_H(x_H)$. This is typically 0 or $x_H^T Q x_H$

This is our base-case and we set P_H to 0 or Q
or Q_{final}



LQR algorithm

- Initialize P_H (at 0 or Q)
- Starting from $h = H-1$, backwards
 - Set $K_h^* = -(R+B^T P_{h+1} B)^{-1} B^T P_{h+1} A$
 - Compute $u = \pi_h^*(x_t) = K_h^* x_t$
 - Set $P_h = (Q + K_h^{*T} R K_h^* + (A + B K_h^*)^T P_{h+1} (A + B K_h^*))$
 - Set $V_h^* = x_t^T P_h x_t$



LQR algorithm

- Initialize P_H (at 0 or Q)
- Starting from $h = H-1$, backwards
 - Set $K_h^* = -(R+B^T P_{h+1} B)^{-1} B^T P_{h+1} A$
 - Compute $u = \pi_h^*(x_t) = K_h^* x_t$
 - Set $P_h = (Q + K_h^{*T} R K_h^* + (A + B K_h^*)^T P_{h+1} (A + B K_h^*))$
 - Set $V_h^* = x_t^T P_h x_t$

This is the Value Iteration update done in closed form: it is always the same and solves this particular continuous-state system with a quadratic cost



LQR algorithm

- Initialize P_H (at 0 or Q)
- Starting from $h = H-1$, backwards
 - Set $K_h^* = -(R+B^T P_{h+1} B)^{-1} B^T P_{h+1} A$
 - Compute $u = \pi_h^*(x_t) = K_h^* x_t$
 - Set $P_h = (Q + K_h^{*T} R K_h^* + (A + B K_h^*)^T P_{h+1} (A + B K_h^*))$
 - Set $V_h^* = x_t^T P_h x_t$

Riccati Equation



LQR algorithm

- Initialize P_H (at 0 or Q)
- Starting from $h = H-1$, backwards
 - Set $K_h^* = -(R+B^T P_{h+1} B)^{-1} B^T P_{h+1} A$
 - Compute $u = \pi_h^*(x_t) = K_h^* x_t$
 - Set $P_h = (Q + K_h^{*T} R K_h^* + (A + B K_h^*)^T P_{h+1} (A + B K_h^*))$
 - Set $\mathbf{J}_h^* = x_t^T P_h x_t$

Since V is the value, we will rename it to J , for this kind of problems, as the “cost-to-go”



LQR Extensions

Extensions to the LQR make it more generally applicable to:

- Affine systems
- Systems with stochasticity
- Regulation around non-zero fixed point for non-linear systems
- Trajectory following for non-linear systems
- ...



LQR for Affine Systems

Affine system: $Ax_t + Bu_t + c$

Cost: $c(x_t, u_t) = x_t^T Q x_t + u_t^T R u_t$

Solve it **in the same way** by simply re-defining the state and transition function as

$$z_t = [x_t; 1]$$

$$z_{t+1} = \begin{bmatrix} x_{t+1} \\ 1 \end{bmatrix} = \begin{bmatrix} A & c \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_t \\ 1 \end{bmatrix} + \begin{bmatrix} B \\ 0 \end{bmatrix} u_t = A' z_t + B' u_t$$



LQR for Stochastic Systems

Stochastic system: $Ax_t + Bu_t + w_t$ with w_t **zero-mean Gaussian noise**

$$w_t \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$$

$$\text{Cost: } c(x_t, u_t) = x_t^T Q x_t + u_t^T R u_t$$

The optimal control policy is the same; cost-to-go has an extra term depending on the variance of the system and that cannot be controlled: $J_h^* = x_t^T P_h x_t + p_h$

$$p_h = \text{tr}(\sigma^2 P_{h+1}) + p_{h+1}$$

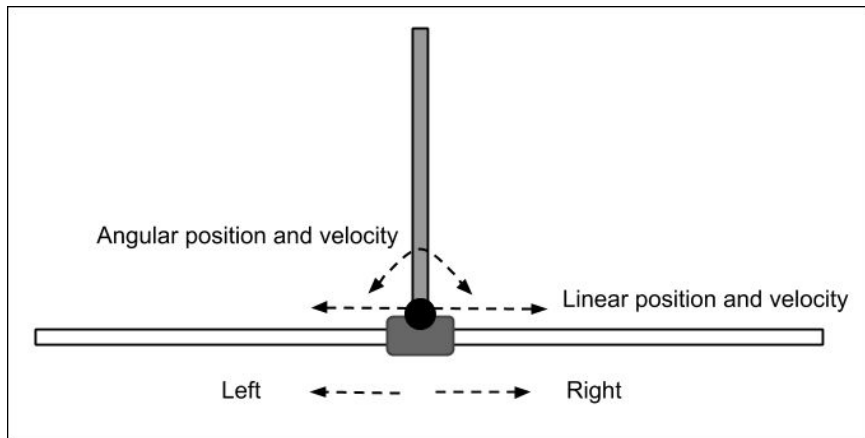
$$\text{base-case } p_H = 0$$



LQR for Non-Linear Systems

Non-linear system: $\mathbf{f}(\mathbf{x}_t, \mathbf{u}_t)$

Cost: $c(\mathbf{x}_t, \mathbf{u}_t) = \mathbf{x}_t^T \mathbf{Q} \mathbf{x}_t + \mathbf{u}_t^T \mathbf{R} \mathbf{u}_t$



We have a goal: $\mathbf{x}^*, \mathbf{u}^*$

We want to stabilize the system around it, such that:

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



LQR for Non-Linear Systems

Non-linear system: $\mathbf{f}(\mathbf{x}_t, \mathbf{u}_t)$

Cost: $c(\mathbf{x}_t, \mathbf{u}_t) = \mathbf{x}_t^T \mathbf{Q} \mathbf{x}_t + \mathbf{u}_t^T \mathbf{R} \mathbf{u}_t$

We can linearize the dynamics around $\mathbf{f}(\mathbf{x}^*, \mathbf{u}^*)$ using Taylor expansion:

$$\mathbf{x}_{t+1} \approx \mathbf{f}(\mathbf{x}^*, \mathbf{u}^*) + \nabla_{\mathbf{x}} \mathbf{f}(\mathbf{x}^*, \mathbf{u}^*) (\mathbf{x}_t - \mathbf{x}^*) + \nabla_{\mathbf{u}} \mathbf{f}(\mathbf{x}^*, \mathbf{u}^*) (\mathbf{u}_t - \mathbf{u}^*)$$



LQR for Non-Linear Systems

Non-linear system: $\mathbf{f}(\mathbf{x}_t, \mathbf{u}_t)$

Cost: $c(\mathbf{x}_t, \mathbf{u}_t) = \mathbf{x}_t^T \mathbf{Q} \mathbf{x}_t + \mathbf{u}_t^T \mathbf{R} \mathbf{u}_t$

We can linearize the dynamics around $\mathbf{f}(\mathbf{x}^*, \mathbf{u}^*)$ using Taylor expansion:

$$\mathbf{x}_{t+1} \approx \mathbf{f}(\mathbf{x}^*, \mathbf{u}^*) + \nabla_{\mathbf{x}} \mathbf{f}(\mathbf{x}^*, \mathbf{u}^*) (\mathbf{x}_t - \mathbf{x}^*) + \nabla_{\mathbf{u}} \mathbf{f}(\mathbf{x}^*, \mathbf{u}^*) (\mathbf{u}_t - \mathbf{u}^*)$$

$$\mathbf{x}_{t+1} - \mathbf{f}(\mathbf{x}^*, \mathbf{u}^*) \approx \mathbf{A} (\mathbf{x}_t - \mathbf{x}^*) + \mathbf{B} (\mathbf{u}_t - \mathbf{u}^*)$$



LQR for Non-Linear Systems

Non-linear system: $\mathbf{f}(\mathbf{x}_t, \mathbf{u}_t)$

Cost: $g(\mathbf{x}_t, \mathbf{u}_t)$

We can linearize the dynamics around $\mathbf{f}(\mathbf{x}^*, \mathbf{u}^*)$ using Taylor expansion:

$$\mathbf{x}_{t+1} \approx \mathbf{f}(\mathbf{x}^*, \mathbf{u}^*) + \nabla_{\mathbf{x}} \mathbf{f}(\mathbf{x}^*, \mathbf{u}^*) (\mathbf{x}_t - \mathbf{x}^*) + \nabla_{\mathbf{u}} \mathbf{f}(\mathbf{x}^*, \mathbf{u}^*) (\mathbf{u}_t - \mathbf{u}^*)$$

$$\mathbf{x}_{t+1} - \mathbf{x}^* \approx \mathbf{A}(\mathbf{x}_t - \mathbf{x}^*) + \mathbf{B}(\mathbf{u}_t - \mathbf{u}^*)$$

$$\text{Let } \mathbf{z}_t = \mathbf{x}_t - \mathbf{x}^* \text{ and } \mathbf{v}_t = \mathbf{u}_t - \mathbf{u}^*$$



LQR for Non-Linear Systems

Non-linear system: $\mathbf{f}(\mathbf{x}_t, \mathbf{u}_t)$

Assume cost is: $c(\mathbf{x}_t, \mathbf{u}_t) = \mathbf{x}_t^T \mathbf{Q} \mathbf{x}_t + \mathbf{u}_t^T \mathbf{R} \mathbf{u}_t$

Let $\mathbf{z}_t = \mathbf{x}_t - \mathbf{x}^*$ and $\mathbf{v}_t = \mathbf{u}_t - \mathbf{u}^*$

$$\mathbf{z}_{t+1} = \mathbf{A} \mathbf{z}_t + \mathbf{B} \mathbf{v}_t$$

Solve it the usual way: $\mathbf{v}_t = \mathbf{K} \mathbf{z}_t \rightarrow \mathbf{u}_t = \mathbf{u}^* + \mathbf{K}(\mathbf{x}_t - \mathbf{x}^*)$



LQR for Non-Linear Systems

Non-linear system: $\mathbf{f}(\mathbf{x}_t, \mathbf{u}_t)$

Cost: $g(\mathbf{x}_t, \mathbf{u}_t)$

We can also linearize the cost using 2nd order Taylor expansion (do it at home), but it becomes something like:

$$\mathbf{x}_t^T \mathbf{Q} \mathbf{x}_t + \mathbf{u}_t^T \mathbf{R} \mathbf{u}_t + \mathbf{u}_t^T \mathbf{M} \mathbf{x}_t + \mathbf{x}_t^T \mathbf{q} + \mathbf{u}_t^T \mathbf{r} + c$$



Penalize for Change in Controls

Standard system: $Ax_t + Bu_t$

Cost: $c(x_t, u_t) = x_t^\top Q x_t + u_t^\top R u_t$

$$\underbrace{\begin{bmatrix} x_{t+1} \\ u_t \end{bmatrix}}_{x'_{t+1}} = \underbrace{\begin{bmatrix} A & B \\ 0 & I \end{bmatrix}}_{A'} \underbrace{\begin{bmatrix} x_t \\ u_{t-1} \end{bmatrix}}_{x'_t} + \underbrace{\begin{bmatrix} B \\ I \end{bmatrix}}_{B'} \underbrace{\Delta u_t}_{u'_t}$$

$$\text{cost} = -(x'^\top Q' x' + \Delta u^\top R' \Delta u) \quad Q' = \begin{bmatrix} Q & 0 \\ 0 & R \end{bmatrix}$$



LQR for Linear Time Varying (LTV) Systems

LTV system: $A_t x_t + B_t u_t$ (Note that A , B , Q , R also have t)

Cost: $c(x_t, u_t) = x_t^T Q_t x_t + u_t^T R_t u_t$



LQR for Linear Time Varying (LTV) Systems

LTV system: $A_t x_t + B_t u_t$ (Note that A , B , Q , R also have t)

Cost: $c(x_t, u_t) = x_t^T Q_t x_t + u_t^T R_t u_t$

- Initialize P_H (at 0 or Q)
- Going backwards (with $i=1,2,\dots$)
 - Set $K_{H-i}^* = -(R_{H-i} + B_{H-i}^T P_{H-i+1} B_{H-i})^{-1} B_{H-i}^T P_{H-i+1} A_{H-i}$
 - Compute $u = \pi_{H-i}^*(x_t) = K_{H-i}^* x_t$
 - Set $P_{H-i} = (Q_{H-i} + K_{H-i}^{*T} R_{H-i} K_{H-i}^* + (A_{H-i} + B_{H-i} K_{H-i}^*)^T P_{H-i+1} (A_{H-i} + B_{H-i} K_{H-i}^*))$
 - Set $V_{H-i}^* = x_t^T P_{H-i} x_t$



LQR for Trajectory Following in Non-Linear Systems

Non-linear system: $\mathbf{f}(\mathbf{x}, \mathbf{u})$

Trajectory: $\mathbf{x}_0^*, \dots, \mathbf{x}_T^*$

We want controls such that $\mathbf{x}_{t+1}^* = \mathbf{f}(\mathbf{x}_t^*, \mathbf{u}_t^*)$ for all the trajectory states

LQR for Trajectory Following in Non-Linear Systems

Non-linear system: $\mathbf{f}(\mathbf{x}, \mathbf{u})$

Trajectory: $\mathbf{x}_0^*, \dots, \mathbf{x}_T^*$

We want controls such that $\mathbf{x}_{t+1}^* = \mathbf{f}(\mathbf{x}_t^*, \mathbf{u}_t^*)$ for all the trajectory states

In other words we are doing:

$$\begin{aligned} \min_{u_0, u_1, \dots, u_{H-1}} \sum_{t=0}^{H-1} (\mathbf{x}_t - \mathbf{x}_t^*)^\top Q (\mathbf{x}_t - \mathbf{x}_t^*) + (\mathbf{u}_t - \mathbf{u}_t^*)^\top R (\mathbf{u}_t - \mathbf{u}_t^*) \\ \text{s.t. } \mathbf{x}_{t+1} = \mathbf{f}(\mathbf{x}_t, \mathbf{u}_t) \end{aligned}$$



LQR for Trajectory Following in Non-Linear Systems

Non-linear system: $\mathbf{f}(\mathbf{x}, \mathbf{u})$

Trajectory: x_0^*, \dots, x_T^*

We want controls such that $x_{t+1}^* = f(x_t^*, u_t^*)$ for all the trajectory states

We use the same trick as before but linearizing in x_{t+1}^* and u_{t+1}^* , then treat it as LTV system (the linearization point changes!)

$$x_{t+1} - x_{t+1}^* \approx \mathbf{A}_t(x_t - x_t^*) + \mathbf{B}_t(u_t - u_t^*)$$

$$u_t - u_t^* = K_i(x_{t+1} - x_{t+1}^*)$$



Iterative LQR

Attempts to solve the most general case:

$$\min_u \sum_{h=0}^H g_h(x_h, u_h)$$

$$\text{subject to } x_{t+1} = f(x_t, u_t)$$

by iteratively approximating it and leveraging the lqr formulation



Iterative LQR - Algorithm

— — —

Initialization: pick either a control policy $\pi^{(0)}$ (start at step 1) or a sequence of states and controls $\{x^{(0)}_t, u^{(0)}_t\}$, $t = 0, \dots, H$

Steps, for $i = 0, \dots, I$ iterations:

1. Execute $\pi^{(i)}$ and collect $\{x^{(i)}_t, u^{(i)}_t\}$, $t = 0, \dots, H$
2. Obtain the LQ approximation of the problem, by linearizing around $\{x^{(i)}_t, u^{(i)}_t\}$ using 1st order Taylor for dynamics and 2nd for cost
3. Compute $\pi^{(i+1)}$ by using the LQR-LTV formulation
4. Iterate



Iterative LQR - LTV

$$x_{t+1} - x_{t+1}^{(i)} = f(x_t^{(i)}, u_t^{(i)}) - x_{t+1}^{(i)} + \nabla_x f(x_t^{(i)}, u_t^{(i)}) (x_t - x_t^{(i)}) + \nabla_u f(x_t^{(i)}, u_t^{(i)}) (u_t - u_t^{(i)})$$

$$z_t = [x_t - x_t^{(i)} \quad 1]^\top$$

$$v_t = (u_t - u_t^{(i)})$$

$$A_t = \begin{bmatrix} \frac{\partial f}{\partial x}(x_t^{(i)}, u_t^{(i)}) & f(x_t^{(i)}, u_t^{(i)}) - x_{t+1}^{(i)} \\ 0 & 1 \end{bmatrix}$$

$$B_t = \begin{bmatrix} \frac{\partial f}{\partial u}(x_t^{(i)}, u_t^{(i)}) \\ 0 \end{bmatrix}$$

And similar derivation for Q and R (you can do it yourselves)



Iterative LQR - Convergence

— — —

As it is, it might not converge, because the optimal policy for the LQ approximation might end up not staying close to the linearization points



Iterative LQR - Convergence

As it is, it might not converge, because the optimal policy for the LQ approximation might end up not staying close to the linearization points

Solution: set the cost function as

$$(1 - \alpha)g(x_t, u_t) + \alpha(\|x_t - x_t^{(i)}\|_2^2 + \|u_t - u_t^{(i)}\|_2^2)$$

and if g is bounded and α close enough to 1, 2nd term will dominate and ensure linearizations are good approximations of the solution trajectory



Iterative LQR - Tricks

— — —

- Good initialization matters, as you can get stuck in local optima
 - **Alternative solution:** use simulated annealing (if you don't have computation constraints) or line-search
- if g is non-convex then LQ approximation might lead to non positive-definite Q and R
 - **Solution:** add more cost for deviating from linearization points until they get positive-definite



Receding Horizon Control

— — —

At iLQR convergence, the system could execute the policy and not be at the expected trajectory

e.g., due to inaccurate dynamics, noise, etc.

Solution-idea: at every timestep t , replan from t to H

Practical solution: replanning at every timestep until H is too costly, but replan over shorter horizon and use cost-to-go $J^{(t+h)}$

This is in practice a Model-Predictive Controller, with an iLQR solver

