

Lecture 6: The Linear Quadratic Regulator

Last time, we discussed finite horizon MDPs with continuous state & action spaces. We also introduced linear dynamics (transitions).

today we consider a continuous MDP problem with linear dynamics and quadratic costs: the Linear Quadratic Regulator.

1) LQR LQR(A, B, Q, R)

For this continuous MDP,

$$\mathcal{S} = \mathbb{R}^{n_s}, \quad \mathcal{A} = \mathbb{R}^{n_a}$$

$$f(s_t, a_t, w_t) = A s_t + B a_t + w_t$$

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

$c(s, u) = s^T Q s + u^T R u$ quadratic costs, $Q, R \succ 0$ symmetric.
Horizon H & initial distribution μ_0 .

Putting this all together, Optimal Control Problem:

$$\min_{\pi} \mathbb{E} \left[s_H^T Q s_H + \sum_{t=0}^{H-1} s_t^T Q s_t + a_t^T R a_t \right]$$

$$\left. \begin{aligned} s_{t+1} &= A s_t + B a_t + w_t, \quad x_0 \sim \mu_0 \\ a_t &= \pi_t(s_t), \quad w_t \sim \mathcal{N}(0, \sigma^2 I) \end{aligned} \right\}$$

Example: 1d robot from last lecture

$$s_t = \begin{bmatrix} p_t \\ v_t \end{bmatrix} \begin{array}{l} \text{position} \\ \text{velocity} \end{array}$$

$$s_{t+1} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} s_t + \begin{bmatrix} 0 \\ 1/m \end{bmatrix} a_t$$

goal: move to goal position $p=0$ and be still $v=0$ without using too much force

Can use quadratic cost:

$$C(s, a) = \gamma_p p^2 + \gamma_v v^2 + \gamma_a a^2$$

$$= s^T \underbrace{\begin{bmatrix} \gamma_p & 0 \\ 0 & \gamma_v \end{bmatrix}}_Q s + \underbrace{\gamma_a}_{R} a^2$$

Depending on the relative weighting of γ_p , γ_v , and γ_a , optimal policy will be more or less aggressive.

Value and Q functions "cost to go"

$$V_t^\pi(s) = \mathbb{E} \left[s_H^T Q s_H + \sum_{k=t}^{H-1} s_k^T Q s_k + a_k^T R a_k \mid \begin{array}{l} s_{k+1} = A s_k + B a_k + w_k \\ a_k = \pi_k(s_k) \\ w_k \sim \mathcal{N}(0, \sigma^2 I) \\ s_t = s \end{array} \right]$$

$$Q_t^\pi(s, a) = \mathbb{E} \left[s_H^T Q s_H + \sum_{k=t}^{H-1} s_k^T Q s_k + a_k^T R a_k \mid \begin{array}{l} s_{k+1} = A s_k + B a_k + w_k \\ a_k = \pi_k(s_k) \\ w_k \sim \mathcal{N}(0, \sigma^2 I) \\ s_t = s, a_t = a \end{array} \right]$$

$$V_t^\pi(s) = Q_t^\pi(s, \pi(s))$$

Notice that due to terminal cost, $V_H(s)$ is nonzero

2) Optimal LQR Policy

We can derive the optimal LQR Policy via Dynamic Programming.

DP For optimal control:

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

Start: $V_H^*(s) = c_H(s)$

(due to terminal cost, initialize $V_H^*(s) \neq 0$ before, but the algorithm is the same).
unlike algorithm

for $t = H-1, H-2, \dots, 0$:

$$Q_t^*(s, a) = c_t(s, a) + \mathbb{E}_{s' \sim P(s, a)} [V_{t+1}^*(s')] = Q_{t+1}^*(s, \pi_{t+1}^*(s))$$

$$\pi_t^* = \underset{a}{\operatorname{argmin}} Q_t^*(s, a)$$

using "minimize cost" convention

Theorem (LQR optimal value fn. & Policy):

For LQR(A, B, Q, R), The optimal value function is quadratic:

$$V_t^*(s) = s^T P_t s + p_t$$

and the optimal policy is linear

$$\pi_t^*(s) = -K_t^* s$$

where (P_t, p_t, K_t^*) can be computed exactly from (A, B, Q, R) .

Proof: We prove by induction, using DP.

claim 1: (Base case) $V_H^*(s) = s^T P_H s + p_H$ is quadratic. $\forall s$

claim 2: (induction) Assume $V_{t+1}^*(s) = s^T P_{t+1} s + p_{t+1}$ $\forall s$. Then

1) $Q_t^*(s, a)$ is quadratic in s, a

2) $\pi_t^*(a) = \underset{a}{\operatorname{argmin}} Q_t^*(s, a)$ is linear in s

Therefore, $V_t^*(s) = s^T P_t s + p_t$ is quadratic.

Then by induction, V is quadratic & π linear.

Proof of claim 1:

$$V_H^*(s) = C_H(s) = s^T Q s. \quad \text{So } P_H = Q, \quad p_H = 0. \checkmark$$

(symmetric)

Proof of claim 2:

Part 1) $Q_t^*(s, a) = s^T A s + a^T R a + \mathbb{E}_{s' \sim P(s, a)} [V_{t+1}^*(s')]$

$$\mathbb{E}_{s'} [V_{t+1}^*(s')] = \mathbb{E}_{w \sim N(0, I)} [V_{t+1}^*(As + Ba + w)]$$

$$\begin{aligned} V_{t+1}^*(As + Ba + w) &= (As)^T P_{t+1} (As) + (As)^T P_{t+1} Ba + (As)^T P_{t+1} w \\ &+ (Ba)^T P_{t+1} As + (Ba)^T P_{t+1} Ba + (Ba)^T P_{t+1} w \\ &+ w^T P_{t+1} As + w^T P_{t+1} Ba + w^T P_{t+1} w + p_{t+1} \end{aligned}$$

also P_{t+1} symmetric

once we take expectation, many terms = 0 because $\mathbb{E}w = 0$.

$$\begin{aligned} \mathbb{E}_{s'} [V_{t+1}^*(s')] &= s^T A^T P_{t+1} A s + 2s^T A^T P_{t+1} Ba + a^T B^T P_{t+1} Ba \\ &+ \mathbb{E}_w [w^T P_{t+1} w] + p_{t+1} \end{aligned}$$

to simplify the remaining expectation, recall
 $W^T P W = \text{Tr}(W^T P W) = \text{Tr}(P W W^T)$ ← cyclic property of trace
 $\mathbb{E}[W^T P W] = \mathbb{E}[\text{Tr}(P W W^T)] = \text{Tr}(P \mathbb{E}[W W^T])$ ← linearity of expectation
 $= \sigma^2 \text{Tr}(P)$

Finally,

$$Q_t^*(s, a) = s^T (Q + A^T P_{t+1} A) s + a^T (R + B^T P_{t+1} B) a + 2 s^T A^T P_{t+1} B a + \sigma^2 \text{tr}(P) + P_{t+1} \quad \checkmark$$

Done with part 1 because this is a quadratic function.

Part 2

$$\pi_t^*(s) = \underset{a}{\text{argmin}} Q_t^*(s, a)$$

First let's derive the minimization for a generic quadratic function.
 $Q(s, a) = s^T M_1 s + a^T \overset{\text{symmetric}}{M_2} a + 2 s^T \overset{\text{not symmetric}}{M_3} a + c$

minimum must occur at a critical point.

$$\begin{aligned} \nabla_a Q(s, a) &= \nabla_a (s^T M_1 s) + \nabla_a (a^T M_2 a) + \nabla_a (s^T M_3 a) \\ &= 0 + 2 M_2 a + 2 M_3^T s \end{aligned}$$

Then $\nabla_a Q(s, a) = 0$ when $M_2 a = -M_3^T s$

for now assuming invertibility, $a = \underbrace{-M_2^{-1} M_3^T s}_{\text{linear function of } s}$

Going back to $Q_t^*(s, a)$
 $M_2 = \underbrace{R + B^T P_{t+1} B}_{\substack{\geq 0 \\ \geq 0 \\ \text{invertible}}} \quad \& \quad M_3 = A^T P_{t+1} B$

Therefore, $\pi_t^*(s) = - \underbrace{(R + B^T P_{t+1} B)^{-1} B^T P_{t+1} A}_{\text{Define } K_t^*}$

Last piece: check that $V_t^*(s) = Q_t^*(s, \pi_t^*(s))$ is quadratic & derive equations for P_t and p_t

Rather than plug in directly, recall our **general quadratic function**

$$\begin{aligned} Q(s, a^*) &= s^T M_1 s + s^T M_3 M_2^{-1} M_2 M_2^{-1} M_3^T s \\ &\quad - 2 s^T M_3 M_2^{-1} M_3^T s + C \\ &= s^T (M_1 - M_3 M_2^{-1} M_3^T) s + C \end{aligned}$$

Therefore, this is the form of $V_t^*(s)$.

Plugging M_1, M_2, M_3, C in, we have

$$\begin{aligned} V_t^*(s) &= s^T (Q + A^T P_{t+1} A - A^T P_{t+1} B (R + B^T P_{t+1} B)^{-1} B^T P_{t+1} A) s \\ &\quad + \sigma^2 \text{tr}(P_{t+1}) + p_{t+1} \quad P_t \end{aligned}$$

✓

This concludes the proof of Claim 2 and therefore the proof of Theorem. □

Collecting the iterative Definitions together:

$$P_H = Q, \quad P_H = 0$$

for $t = H-1, \dots, 0$:

$$P_t = Q + A^T P_{t+1} A - A^T P_{t+1} B (R + B^T P_{t+1} B)^{-1} B^T P_{t+1} A$$

$$P_t = P_{t+1} + \sigma^2 \text{tr}(P_{t+1})$$

$$K_t^* = (R + B^T P_{t+1} B)^{-1} B^T P_{t+1} A$$

Some straight forward Extensions:

1) time-varying costs / dynamics

e.g. $s_{t+1} = A_t s_t + B_t a_t + w_t$

$$c_t(s, a) = s_t^T Q_t s_t + a_t^T R_t a_t$$

2) non-stochastic disturbance

$$s_{t+1} = A s_t + B a_t + w_t + v_t$$

where v_t is known a priori

3) trajectory tracking

$$c_t(s, a) = (s - s_t^*)^T Q (s - s_t^*) + (a - a_t^*)^T R (a - a_t^*)$$

for desired trajectory (s_0^*, a_0^*, \dots) known a priori.

(this case can be reduced to case 2 if substitute $s \leftarrow s - s_t^*$ and $a \leftarrow a - a_t^*$)