# **Optimal Control for the Linear** Quadratic Regulator

### Lucas Janson and Sham Kakade **CS/Stat 184: Introduction to Reinforcement Learning** Fall 2022



- Feedback from last lecture
- Recap
- Derivation of optimal LQR policy
- Extensions



### Feedback from feedback forms

1. Thank you to everyone who filled out the forms! 2.





- Recap
- Derivation of optimal LQR policy
- Extensions

## Recap: LQR

### Problem Statement:



Value function for a policy  $\pi = (\pi_0, V_t^{\pi}(s)) = \mathbb{E} \left[ s_T^{\top} Q s_T + \sum_{i=t}^{T-1} (s_i^{\top} Q s_i) \right]$ 

And corresponding Q function: T-1

$$Q_t^{\pi}(s,a) = \mathbb{E}\left[s_T^{\top}Qs_T + \sum_{i=t}^{T-1} (s_i^{\top}Qs_i + a_i^{\top}Ra_i) \mid a_t = a, a_i = \pi_i(s_i) \; \forall i > t, \, s_t = s\right]$$



### Feedback from last lecture

- Recap
  - Derivation of optimal LQR policy
  - Extensions



## LQR Optimal Control

## $V_t^{\star}(s) = \min_{\pi} V_t^{\pi}(s) = \min_{\pi_t, \pi_{t+1}, \dots, \pi_{T-1}} \mathbb{E} \left[ s_T^{\top} Q s_T \right]$

### **Theorem**:

2. The optimal policy  $\pi_t^{\star}$  is linear, i.e.,  $\pi_t^{\star}(s) = -K_t s$  for some  $K_t \in \mathbb{R}^{k \times d}$ 3.  $P_t$ ,  $p_t$ , and  $K_t$  can be computed exactly

Today: prove the above theorem, deriving the optimal policy along the way

$$+ \sum_{i=t}^{T-1} (s_i^{\mathsf{T}} Q s_i + a_i^{\mathsf{T}} R a_i) \mid a_i = \pi_i(s_i) \; \forall i \ge t, \; s_t =$$

1.  $V_t^{\star}$  is a quadratic function, i.e.,  $V_t^{\star}(s) = s^{\top}P_t s + p_t$  for some  $P_t \in \mathbb{R}^{d \times d}$  and  $p_t \in \mathbb{R}^d$ 



## Key Steps in the Proof

1. Base case: Show that  $V_T^{\star}(s)$  is quadratic

- 2. Inductive hypothesis: Assuming  $V_{t+1}^{\star}(s)$  is quadratic,
  - a) Show that  $Q_t^{\star}(s, a)$  is quadratic (in both s and a)

  - c) Show  $V_{\tau}^{\star}(s)$  is quadratic

3. Conclusion:  $V_t^{\star}(s)$  is quadratic and  $\pi_t^{\star}(s)$  is linear and we'll have their formulas

- Dynamic programming (finite-horizon), stepping backwards in time from T to 0

b) Derive the optimal policy  $\pi_t^{\star}(s) = \arg \min Q_t^{\star}(s, a)$ , and show that it's linear



### Base case at *I*

$$V_t^{\pi}(s) = \mathbb{E}\left[s_T^{\top} Q s_T + \sum_{i=t}^{T-1} \left(s_i^{\top} Q s_i + a_i^{\top} R a_i\right) \mid a_i = \pi_i(s_i) \; \forall i \ge t, \; s_t = s\right]$$

For  $V_T^{\pi}$ , everything disappears except first term  $s_T^{\top}Qs_T = s^{\top}Qs$ :  $V_T^{\star}(s) = s^{\top} Q s$ 

> Denoting  $P_T := \zeta$  $V_T^{\star}(s) =$

 $(P_t \text{ and } p_t \text{ didn't do much here, but we're going to define them recursively in the next step)$ 

Recall the value function at a given t is:

Q and 
$$p_T := 0$$
, we get  
=  $s^T P_T s + p_T$ 

## Induction Step

### $Q_t^{\star}(s, a) =$

Assume  $V_{t+1}^{\star}(s) = s^{\top}P_{t+1}s + p_{t+1}$ , for all s, where  $P_{t+1} \in \mathbb{R}^{d \times d}$  and  $p_{t+1} \in \mathbb{R}^{d}$ 

## Induction Step (continued)

 $\begin{aligned} Q_t^{\star}(s,a) &= c(s,a) + \mathbb{E}_{s' \sim f(s,a,w_{t+1})} \left[ V_{t+1}^{\star}(s') \right] \\ &= s^{\top} \left( Q + A^{\top} P_{t+1} A \right) s + a^{\top} \left( R + B^{\top} P_{t+1} B \right) a + 2s^{\top} A^{\top} P_{t+1} B a + \text{tr} \left( \sigma^2 P_{t+1} \right) + p_{t+1} ds \end{aligned}$ 

 $\pi_t^{\star}(s) = \arg \min_a Q_t^{\star}(s, a)$ Set  $\nabla_a Q_t^{\star}(s, a) = 0$  and solve for *a*:  $\nabla_a Q_t^{\star}(s, a) =$ 

### Concluding the Induction step:

 $Q_{t}^{\star}(s,a) = s^{\top} \left( Q + A^{\top} P_{t+1} A \right) s + a^{\top} \left( R + B^{\top} P_{t+1} B \right) a + 2s^{\top} A^{\top} P_{t+1} B a + \text{tr} \left( \sigma^{2} P_{t+1} \right) + p_{t+1} B$  $\pi_t^{\star}(s) = -(R + B^{\top} P_{t+1} B)^{-1} B^{\top} P_{t+1} A s$  $=K_t$  $\mathbf{V}_{\mathbf{X}}$ 

$$V_t^{\star}(s) = Q_t^{\star}(s, \pi_t^{\star}(s))$$
  
=  $s^{\top} \left( Q + A^{\top} P_{t+1} A \right) s + s^{\top} K_t^{\top} \left( R + A^{\top} P_{t+1} A \right) s$ 

$$P_t = Q + A^{\mathsf{T}} P_{t+1} A - A^{\mathsf{T}} P_{t+1} B (R + B^{\mathsf{T}} P_t)$$
$$p_t = \operatorname{tr} \left( \sigma^2 P_{t+1} \right) + p_{t+1}$$

- $+ B^{T}P_{t+1}B K_{t}s 2s^{T}A^{T}P_{t+1}BK_{t}s + tr (\sigma^{2}P_{t+1}) + p_{t+1}$
- Collecting the quadratic and constant terms together,  $V_t^{\star}(s) = s^{\top}P_t s + p_t$ , where:



$$V_T^{\star}(s) = s^{\top} Q s,$$

### We have shown that $P_t = Q + A^{\mathsf{T}} P_{t+1} A - A^{\mathsf{T}}$ $p_t = \text{tr}(\sigma^2 P_{t+1}) + p_{t+1}$

 $K_t = (R + B)$ 

Optimal policy has nothing to do with initial distribution  $\mu_0$  or the noise  $\sigma^2$ !

### Summary:

define 
$$P_T = Q, p_T = 0$$
,

$$V_t^{\star}(s) = s^{\top} P_t s + p_t, \text{ where:}$$
$$P_{t+1}B(R + B^{\top} P_{t+1}B)^{-1}B^{\top} P_{t+1}A$$

Along the way, we also have shown that  $\pi_t^{\star}(s) = -K_t s$ , where:

$$^{\mathsf{T}}P_{t+1}B)^{-1}B^{\mathsf{T}}P_{t+1}A$$

- Feedback from last lecture
- Recap
- Derivation of optimal LQR policy
  - Extensions



### Time-Dependent Costs and Dynamics



Exact same derivation, only thing that changes is the Ricatti equation:  $P_t = Q_t + A_t^{\top} P_{t+1} A_t - A_t^{\top} P_{t+1} B_t (R_t + B_t^{\top} P_{t+1} B_t)^{-1} B_t^{\top} P_{t+1} A_t$ 

$$\begin{bmatrix} -1 \\ S_t & T \\ Q_t & S_t + a_t^T & R_t \\ S_t & S_t + a_t^T & R_t \\ S_t & S_t & T \\ S_t & S_t & S_t \\ S_t$$

### More General Quadratic Cost Function

# $\arg\min_{\pi_0,\ldots,\pi_{T-1}:\mathbb{R}^d\to\mathbb{R}^k} \mathbb{E} \left| s_T^{\mathsf{T}}Q_T s_T + s_T^{\mathsf{T}}q_T + c_T + s_T^{\mathsf{T}}q_T + c_T \right|$ such that $s_{t+1} = A_t s_t + B_t a_t + v_t + w_t$ , $s_0 \sim$

$$-\sum_{t=0}^{T-1} \left( s_t^{\mathsf{T}} Q_t s_t + a_t^{\mathsf{T}} R_t a_t + a_t^{\mathsf{T}} M_t s_t + s_t^{\mathsf{T}} q_t + a_t^{\mathsf{T}} r_t + a_t^{\mathsf{T}} r_t \right)$$

$$\sim \mu_0$$
,  $a_t = \pi_t(s_t)$ ,  $w_t \sim N(0, \sigma^2 I)$ 

### Derivation is similar—you will work it out on HW3



### Tracking a Predefined Trajectory



### Expanding all the quadratic terms produces a special case of the previous slide!

$$+ \sum_{t=0}^{T-1} \left( (s_t - s_t^{\star})^\top Q_t (s_t - s_t^{\star}) + (a_t - a_t^{\star})^\top R_t (a_t$$





- Feedback from last lecture
- Recap
- Derivation of optimal LQR policy
- Extensions



## Today's summary:

LQR optimal policy/controller

- Used dynamic programming / inductive argument to derive  $\pi_{t}^{\star}$ Same argument applies to extensions to some more complicated situations

Next time:

 Applying LQR approximation separately at each time point to get a locally-optimal solution to a nonlinear control problem

1-minute feedback form: <u>https://bit.ly/3RHtlxy</u>



