

# AA203: Optimal and Learning-based Control

## Course Notes

James Harrison

April 8, 2019

## 2 Dynamic Programming and the Linear Quadratic Regulator

### 2.1 The Optimal Control Problems

In this section, we will outline the deterministic continuous-time optimal control problem that we will aim to solve. We will denote the state at time  $t$  as  $\mathbf{x}(t) \in \mathbb{R}^n$ , and the control as  $\mathbf{u}(t) \in \mathbb{R}^m$ . We will also occasionally write these as  $\mathbf{x}_t$  and  $\mathbf{u}_t$ , respectively. We will write the continuous-time systems dynamics as

$$\dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t), t). \quad (1)$$

We will refer to a history of control input values during an interval  $[t_0, t_f]$  as a control history, and we will refer to a history of state values over this interval as a state trajectory.

Different control problems may call for various constraints. For example, we may constrain a quadrotor to only fly in space not occupied by obstacles. Examples of constraints we will see are

- Initial and final conditions,  $\mathbf{x}(t_0) = \mathbf{x}_0$ ,  $\mathbf{x}(t_f) = \mathbf{x}_f$
- Trajectory constraints,  $\underline{\mathbf{x}} \leq \mathbf{x}(t) \leq \bar{\mathbf{x}}$
- Control limits,  $\underline{\mathbf{u}} \leq \mathbf{u}(t) \leq \bar{\mathbf{u}}$ .

A state trajectory and control history that satisfy the constraints during the entire time interval  $[t_0, t_f]$  are called admissible trajectories and admissible controls, respectively.

Finally, we will define the performance measure,

$$J = c_f(\mathbf{x}(t_f), t_f) + \int_{t_0}^{t_f} c(\mathbf{x}(t), \mathbf{u}(t), t) dt \quad (2)$$

where  $c$  is the instantaneous cost function, and  $c_f$  is the terminal state cost. We are now able to state the continuous-time optimal control problem. We aim to find an admissible control,

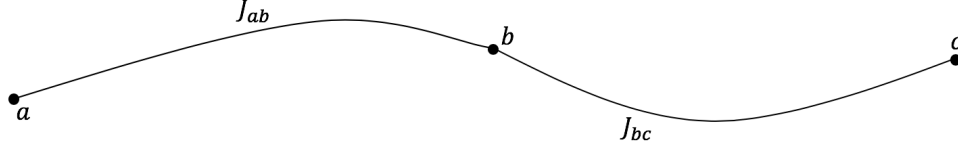


Figure 1: An optimal trajectory connecting point  $a$  to point  $c$ . There are no better (lower cost) trajectories than the sub-trajectory connecting  $b$  and  $c$ , by the principle of optimality.

$\mathbf{u}^*$ , which causes the system (1) to follow an admissible trajectory,  $\mathbf{x}^*$ , that minimizes the performance measure given by (2). The minimizer  $(\mathbf{x}^*, \mathbf{u}^*)$  is called an optimal trajectory-control pair.

Note, first of all, that this is an extremely general problem formulation. We have not fixed our system dynamics, cost function, or specific constraints. We can't, in general, guarantee the existence or uniqueness of the optimal solution.

There are two possible solution forms for the optimal control. The first,  $\mathbf{u}^* = e(\mathbf{x}(t_0), t)$  is referred to as an open-loop solution. This is an input function that is applied to the system, without using feedback. Practically, such solutions usually require augmentation with a feedback controller, as small model mismatch may lead to compounding errors. The second possible solution form is a feedback policy,  $\mathbf{u}^* = \pi(\mathbf{x}(t), t)$ . This feedback law maps all state-time pairs to an action and thus is usually more robust to possible model mismatch. However, depending on the particular problem formulation, open-loop solutions may be easier to compute.

## 2.2 Dynamic Programming and the Principle of Optimality

In this chapter we will outline the principle of optimality, and the method of dynamic programming (DP), one of two main approaches to solving the optimal control problem. The second, so-called variational approaches based on Pontryagin's Maximum Principle (PMP) will be discussed in future chapters. Dynamic programming has the strong advantage of yielding a feedback policy, however, exactly solving the dynamic programming problem is infeasible for many systems. We will address special cases in which the DP problem can be solved exactly, and approximate methods that work for a wide variety of systems.

Despite having just introduced the optimal control problem in continuous time, we will be operating in discrete time here, in which we aim to minimize

$$J_f(\mathbf{x}_0) = c_f(\mathbf{x}_N) + \sum_{k=0}^{N-1} c(\mathbf{x}_k, \mathbf{u}_k, k). \quad (3)$$

We will extend the methods we develop in this chapter to continuous time in the next chapter.

The principle of optimality is as follows. Figure 1 shows a trajectory from point  $a$  to  $c$ . If the cost of the trajectory,  $J_{ac} = J_{ab} + J_{bc}$ , is minimal, then  $J_{bc}$  is also a minimum cost trajectory connecting  $b$  and  $c$ . The proof of this principle, stated informally, is simple.

Assume there exists an alternative trajectory connecting  $b$  and  $c$ , for which we will write the cost as  $\tilde{J}_{bc}$ , that achieves  $\tilde{J}_{bc} < J_{bc}$ . Then, we have

$$\tilde{J}_{ac} = J_{ab} + \tilde{J}_{bc} \quad (4)$$

$$< J_{ab} + J_{bc} \quad (5)$$

$$= J_{ac}, \quad (6)$$

and thus  $J_{ac}$  isn't minimal. More formally,

**Theorem 2.1** (Discrete-time Principle of Optimality). *Let  $\pi^* = (\pi_0^*, \dots, \pi_{N-1}^*)$  be an optimal policy. Assume state  $\mathbf{x}_k$  is reachable. Consider the subproblem whereby we are at  $\mathbf{x}_k$  at time  $k$  and we wish to minimize the cost-to-go from time  $k$  to time  $N$ . Then the truncated policy  $(\pi_k^*, \dots, \pi_{N-1}^*)$  is optimal for the subproblem.*

Dynamic programming, intuitively, proceeds backwards in time, first solving simpler shorter horizon problems. If we have found the optimal policy for times  $k+1$  to  $N-1$ , along with the associated cost-to-go for each state, choosing the optimal policy for time  $k$  is a one step optimization problem. More concretely, we will assume we have dynamics of the form  $\mathbf{x}_{k+1} = f(\mathbf{x}_k, \mathbf{u}_k, k)$  with  $\mathbf{u}_k \in \mathcal{U}(\mathbf{x}_k)$ , and the cost given by (3). Then, dynamic programming iterates backward in time, from  $N-1$  to 0, with

$$J_N(\mathbf{x}_N) = c_T(\mathbf{x}_N) \quad (7)$$

$$J_k(\mathbf{x}_k) = \min_{\mathbf{u}_k \in \mathcal{U}(\mathbf{x}_k)} \{c_k(\mathbf{x}_k, \mathbf{u}_k, k) + J_{k+1}(f(\mathbf{x}_k, \mathbf{u}_k, k))\}. \quad (8)$$

Note that here we have considered only deterministic dynamical systems (there is no stochastic disturbance).

Practically, dynamic programming raises many practical issues if one were to attempt to apply it directly. To perform the recursion,  $J_{k+1}$  must be known for all  $\mathbf{x}_{k+1}$  (or more precisely, all  $\mathbf{x}_{k+1}$  that are reachable from  $\mathbf{x}_k$ ). If the state space is discrete (and relatively small), this is tractable as the cost-to-go may just be maintained in tabular form. In the next subsection, we will discuss an extremely important case in continuous space in which the cost-to-go can be computed exactly for all states. However, for general systems, we can not expect to be able to compute the cost-to-go for all states. Possible approaches to make the DP approach tractable are discretizing the state space, approximating the cost-to-go (i.e. restricting the family of functions that  $J_{k+1}$  may be in), or interpolating between cost-to-go computed for a finite set of states.

## 2.3 Discrete LQR

An important instance in which dynamic programming can be solved analytically for continuous state-action systems is the *linear quadratic regulator* problem. We will fix the dynamics of the system to be (possibly time-varying) linear,

$$\mathbf{x}_{k+1} = A_k \mathbf{x}_k + B_k \mathbf{u}_k \quad (9)$$

and the cost function is quadratic

$$c(\mathbf{x}_k, \mathbf{u}_k) = \frac{1}{2}(\mathbf{x}_k^T Q_k \mathbf{x}_k + \mathbf{u}_k^T R_k \mathbf{u}_k) \quad (10)$$

$$c_N(\mathbf{x}_k) = \frac{1}{2}\mathbf{x}_k^T Q_N \mathbf{x}_k \quad (11)$$

where  $Q_k \in \mathbb{R}^{n \times n}$  is positive semi-definite and  $R_k \in \mathbb{R}^{m \times m}$  is positive definite for all  $k = 0, \dots, N$ . Importantly, we assume  $\mathbf{x}_k$  and  $\mathbf{u}_k$  are unconstrained for all  $k$ . To perform DP recursion, we initialize

$$J_N^*(\mathbf{x}_N) = \frac{1}{2}\mathbf{x}_N^T Q_N \mathbf{x}_N := \frac{1}{2}\mathbf{x}_N^T V_N \mathbf{x}_N. \quad (12)$$

Then, applying (8), we have

$$J_{N-1}^*(\mathbf{x}_{N-1}) = \frac{1}{2} \min_{\mathbf{u}_{N-1} \in \mathbb{R}^m} \{ \mathbf{x}_{N-1}^T Q_{N-1} \mathbf{x}_{N-1} + \mathbf{u}_{N-1}^T R_{N-1} \mathbf{u}_{N-1} + \mathbf{x}_N^T V_N \mathbf{x}_N \} \quad (13)$$

which, applying the dynamics,

$$\begin{aligned} J_{N-1}^*(\mathbf{x}_{N-1}) = \frac{1}{2} \min_{\mathbf{u}_{N-1} \in \mathbb{R}^m} \{ & \mathbf{x}_{N-1}^T Q_{N-1} \mathbf{x}_{N-1} + \mathbf{u}_{N-1}^T R_{N-1} \mathbf{u}_{N-1} \\ & + (A_{N-1} \mathbf{x}_{N-1} + B_{N-1} \mathbf{u}_{N-1})^T V_N (A_{N-1} \mathbf{x}_{N-1} + B_{N-1} \mathbf{u}_{N-1}) \}. \end{aligned} \quad (14)$$

Rearranging, we have

$$\begin{aligned} J_{N-1}^*(\mathbf{x}_{N-1}) = \frac{1}{2} \min_{\mathbf{u}_{N-1} \in \mathbb{R}^m} \{ & \mathbf{x}_{N-1}^T (Q_{N-1} + A_{N-1}^T V_N A_{N-1}) \mathbf{x}_{N-1} \\ & + \mathbf{u}_{N-1}^T (R_{N-1} + B_{N-1}^T V_N B_{N-1}) \mathbf{u}_{N-1} \\ & + 2\mathbf{u}_{N-1}^T (B_{N-1}^T V_N A_{N-1}) \mathbf{x}_{N-1} \}. \end{aligned} \quad (15)$$

Note that this optimization problem is convex in  $\mathbf{u}_{N-1}$  as  $R_{N-1} + B_{N-1}^T V_N B_{N-1} > 0$ . Therefore, any local minima is a global minima, and therefore we can simply apply the first order optimality conditions. Differentiating,

$$\frac{\partial J_{N-1}^*(\mathbf{x}_{N-1})}{\partial \mathbf{u}_{N-1}} = (R_{N-1} + B_{N-1}^T V_N B_{N-1}) \mathbf{u}_{N-1} + (B_{N-1}^T V_N A_{N-1}) \mathbf{x}_{N-1} \quad (16)$$

and setting this to zero yields

$$\mathbf{u}_{N-1}^* = -(R_{N-1} + B_{N-1}^T V_N B_{N-1})^{-1} (B_{N-1}^T V_N A_{N-1}) \mathbf{x}_{N-1} \quad (17)$$

which we write

$$\mathbf{u}_{N-1}^* = L_{N-1} \mathbf{x}_{N-1} \quad (18)$$

which is a time-varying linear feedback policy. Plugging this feedback policy into (14),

$$\begin{aligned} J_{N-1}^*(\mathbf{x}_{N-1}) = & \mathbf{x}_{N-1}^T (Q_{N-1} + L_{N-1}^T R_{N-1} L_{N-1}) \\ & + (A_{N-1} + B_{N-1} L_{N-1})^T V_N (A_{N-1} + B_{N-1} L_{N-1}) \mathbf{x}_{N-1}. \end{aligned} \quad (19)$$

Critically, this implies that the cost-to-go is always a positive semi-definite quadratic function of the state. Because the optimal policy is always linear, and the optimal cost-to-go is always quadratic, the DP recursion may be recursively performed backward in time and the minimization may be performed analytically.

Following the same procedure, we can write the DP recursion for the discrete-time LQR controller:

1.  $V_N = Q_N$
2.  $L_k = -(R_k + B_k^T V_{k+1} B_k)^{-1} (B_k^T V_{k+1} A_k)$
3.  $V_k = Q_k + L_k^T R_k L_k + (A_k + B_k L_k)^T V_{k+1} (A_k + B_k L_k)$
4.  $\mathbf{u}_k^* = L_k \mathbf{x}_k$
5.  $J_k^*(\mathbf{x}_k) = \frac{1}{2} \mathbf{x}_k^T V_k \mathbf{x}_k$

There are several implications of this recurrence relation. First, even if  $A, B, Q, R$  are all constant (not time-varying), the policy is still time-varying. Why is this the case? Control effort invested early in the problem will yield dividends over the remaining length of the horizon, in terms of lower state cost for all future time steps. However, as the remaining length of the episode becomes shorter, this tradeoff is increasingly imbalanced, and the control effort will decrease. However, for a linear time-invariant system, if  $(A, B)$  controllable, the feedback gain  $L_k$  approach a constant as the episode length approaches infinity. This time-invariant policy is practical for long horizon control problems, and may be approximately computed by running the DP recurrence relation until approximate convergence.

## 2.4 Iterative LQR and Differential Dynamic Programming

## 2.5 LQG

## References