

# AA203: Optimal and Learning-based Control

## Course Notes

James Harrison\*

April 20, 2019

### 3 The HJB and HJI Equations

In this section, we will extend the ideas of dynamic programming to the continuous time setting. Restating the continuous time optimal control problem, we assume dynamics

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

and cost

$$J(\mathbf{x}(0)) = c_f(\mathbf{x}(t_f), t_f) + \int_0^{t_f} c(\mathbf{x}(\tau), \mathbf{u}(\tau), \tau) d\tau. \quad (2)$$

where  $t_f$  is fixed.

#### 3.1 The Principle of Optimality in Continuous Time

##### 3.1.1 Hamilton-Jacobi-Bellman

As in the discrete time principle of optimality, consider the tail problem

$$J(\mathbf{x}(t), \{\mathbf{u}(\tau)\}_{\tau=t}^{t_f}, t) = c_f(\mathbf{x}(t_f), t_f) + \int_t^{t_f} c(\mathbf{x}(\tau), \mathbf{u}(\tau), \tau) d\tau \quad (3)$$

where  $t \leq t_f$  and  $\mathbf{x}(t)$  is an admissible state value. The optimal solution to this tail problem comes from the functional minimization

$$J^*(\mathbf{x}(t), t) = \min_{\{\mathbf{u}(\tau)\}_{\tau=t}^{t_f}} \left\{ c_f(\mathbf{x}(t_f), t_f) + \int_t^{t_f} c(\mathbf{x}(\tau), \mathbf{u}(\tau), \tau) d\tau \right\}. \quad (4)$$

---

\*Contact: jharrison@stanford.edu

Note, then, that due to the additivity of cost we can split the problem up over time,

$$J^*(\mathbf{x}(t), t) = \min_{\{\mathbf{u}(\tau)\}_{\tau=t}^{t_f}} \left\{ \int_t^{t+\Delta t} c(\mathbf{x}(\tau), \mathbf{u}(\tau), \tau) d\tau + c_f(\mathbf{x}(t_f), t_f) + \int_{t+\Delta t}^{t_f} c(\mathbf{x}(\tau), \mathbf{u}(\tau), \tau) d\tau \right\} \quad (5)$$

which by applying the principle of optimality to the tail cost,

$$J^*(\mathbf{x}(t), t) = \min_{\{\mathbf{u}(\tau)\}_{\tau=t}^{t+\Delta t}} \left\{ \int_t^{t+\Delta t} c(\mathbf{x}(\tau), \mathbf{u}(\tau), \tau) d\tau + J^*(\mathbf{x}(t+\Delta t), t+\Delta t) \right\}. \quad (6)$$

Let  $J_t^*(\mathbf{x}(t), t) = \nabla_t J^*(\mathbf{x}(t), t)$  and  $J_{\mathbf{x}}^*(\mathbf{x}(t), t) = \nabla_{\mathbf{x}} J^*(\mathbf{x}(t), t)$ . Taylor expanding, we have

$$\begin{aligned} J^*(\mathbf{x}(t), t) = \min_{\{\mathbf{u}(\tau)\}_{\tau=t}^{t+\Delta t}} \{ & c(\mathbf{x}(t), \mathbf{u}(t), t) \Delta t + J^*(\mathbf{x}(t), t) + (J_t^*(\mathbf{x}(t), t)) \Delta t \\ & + (J_{\mathbf{x}}^*(\mathbf{x}(t), t))^T (\mathbf{x}(t+\Delta t) - \mathbf{x}(t)) + o(\Delta t) \} \end{aligned} \quad (7)$$

for small  $\Delta t$ . The first term is a result of Taylor expanding the integral and applying the fundamental theorem of calculus. Note that we can pull  $J^*(\mathbf{x}(t), t)$  out of the minimization over cost, as this quantity will not vary under different choices of future actions. Dividing through by  $\Delta t$  and taking the limit  $\Delta t \rightarrow 0$ , we obtain the *Hamilton-Jacobi-Bellman* equation

$$0 = J_t^*(\mathbf{x}(t), t) + \min_{\mathbf{u}(t)} \{ c(\mathbf{x}(t), \mathbf{u}(t), t) + (J_{\mathbf{x}}^*(\mathbf{x}(t), t))^T f(\mathbf{x}(t), \mathbf{u}(t), t) \} \quad (8)$$

with terminal condition

$$J^*(\mathbf{x}(t_f), t_f) = c_f(\mathbf{x}(t_f), t_f). \quad (9)$$

For convenience, we will define the Hamiltonian

$$\mathcal{H}(\mathbf{x}(t), \mathbf{u}(t), J_{\mathbf{x}}^*, t) := c(\mathbf{x}(t), \mathbf{u}(t), t) + (J_{\mathbf{x}}^*(\mathbf{x}(t), t))^T f(\mathbf{x}(t), \mathbf{u}(t), t) \quad (10)$$

which allow us to compactly write the HJB equation as

$$0 = J_t^*(\mathbf{x}(t), t) + \min_{\mathbf{u}(t)} \{ \mathcal{H}(\mathbf{x}(t), \mathbf{u}(t), J_{\mathbf{x}}^*, t) \}. \quad (11)$$

The HJB equation is a partial differential equation that, for cost-to-go  $J^*(\mathbf{x}(t), t)$ , will satisfy all time-state pairs  $(\mathbf{x}(t), t)$ . The previous informal derivation assumed differentiability of  $J^*(\mathbf{x}(t), t)$ , which we do not know a priori. This assumption is rectified by the following theorem on solutions to the HJB equation.

**Theorem 3.1** (Sufficiency Theorem). *Suppose  $V(\mathbf{x}, t)$  is a solution to the HJB equation, that  $V(\mathbf{x}, t)$  is  $C^1$  in  $\mathbf{x}$  and  $t$ , and that*

$$\begin{aligned} 0 &= V_t(\mathbf{x}, t) + \min_{\mathbf{u} \in \mathcal{U}} \{ c(\mathbf{x}, \mathbf{u}, t) + (V_{\mathbf{x}}(\mathbf{x}, t))^T f(\mathbf{x}, \mathbf{u}, t) \} \\ V(\mathbf{x}, t_f) &= c_f(\mathbf{x}, t_f) \quad \forall \mathbf{x} \end{aligned}$$

Suppose also that  $\pi^*(\mathbf{x}, t)$  attains the minimum in this equation for all  $t$  and  $\mathbf{x}$ . Let  $\{\mathbf{x}^*(t) \mid t \in [t_0, t_f]\}$  be the state trajectory obtained from the given initial condition  $\mathbf{x}(0)$  when the control trajectory  $\mathbf{u}^*(t) = \pi^*(\mathbf{x}^*(t), t), t \in [t_0, t_f]$  is used. Then  $V$  is equal to the optimal cost-to-go function, i.e.,

$$V(\mathbf{x}, t) = J^*(\mathbf{x}, t) \quad \forall \mathbf{x}, t. \quad (12)$$

Furthermore, the control trajectory  $\{\mathbf{u}^*(t) \mid t \in [t_0, t_f]\}$  is optimal..

*Proof.* [Ber12], Volume 1, Section 3.2. □

### 3.1.2 Continuous-Time LQR

As a useful result of the HJB equations, we will derive LQR in continuous time. We aim to minimize

$$J(\mathbf{x}(0)) = \frac{1}{2} \mathbf{x}^T(t_f) Q_f \mathbf{x}(t_f) + \frac{1}{2} \int_0^{t_f} \mathbf{x}^T(t) Q(t) \mathbf{x}(t) + \mathbf{u}^T(t) R(t) \mathbf{u}(t) dt \quad (13)$$

subject to dynamics

$$\dot{\mathbf{x}}(t) = A(t) \mathbf{x}(t) + B(t) \mathbf{u}(t). \quad (14)$$

As in discrete LQR, we will assume  $Q_f, Q(t)$  are positive semidefinite, and  $R(t)$  is positive definite. We will also assume  $t_f$  is fixed, and the state and action are unconstrained.

We will write the Hamiltonian,

$$\mathcal{H} = \frac{1}{2} \mathbf{x}^T(t) Q(t) \mathbf{x}(t) + \frac{1}{2} \mathbf{u}^T(t) R(t) \mathbf{u}(t) + J_{\mathbf{x}}^*(\mathbf{x}(t), t)^T (A(t) \mathbf{x}(t) + B(t) \mathbf{u}(t)) \quad (15)$$

which yields necessary optimality conditions

$$0 = \nabla_{\mathbf{u}} \mathcal{H} = R(t) \mathbf{u}(t) + B^T(t) J_{\mathbf{x}}^*(\mathbf{x}(t), t). \quad (16)$$

Since  $\nabla_{\mathbf{u}\mathbf{u}}^2 \mathcal{H} = R(t) > 0$ , the control that satisfies the necessary conditions is the global minimizer. Rearranging, we have

$$\mathbf{u}^*(t) = -R^{-1}(t) B^T(t) J_{\mathbf{x}}^*(\mathbf{x}(t), t) \quad (17)$$

which we can plug back into the Hamiltonian to yield

$$\mathcal{H} = \frac{1}{2} \mathbf{x}^T(t) Q(t) \mathbf{x}(t) + \frac{1}{2} J_{\mathbf{x}}^*(\mathbf{x}(t), t)^T B(t) R^{-1}(t) B^T(t) J_{\mathbf{x}}^*(\mathbf{x}(t), t) \quad (18)$$

$$\begin{aligned} &+ J_{\mathbf{x}}^*(\mathbf{x}(t), t)^T A(t) \mathbf{x}(t) - J_{\mathbf{x}}^*(\mathbf{x}(t), t)^T B(t) R^{-1}(t) B^T(t) J_{\mathbf{x}}^*(\mathbf{x}(t), t) \\ &= \frac{1}{2} \mathbf{x}^T(t) Q(t) \mathbf{x}(t) - \frac{1}{2} J_{\mathbf{x}}^*(\mathbf{x}(t), t)^T B(t) R^{-1}(t) B^T(t) J_{\mathbf{x}}^*(\mathbf{x}(t), t) + J_{\mathbf{x}}^*(\mathbf{x}(t), t)^T A(t) \mathbf{x}(t). \end{aligned} \quad (19)$$

This gives the HJB equation

$$0 = J_t^*(\mathbf{x}(t), t) + \frac{1}{2}\mathbf{x}^T(t)Q(t)\mathbf{x}(t) - \frac{1}{2}J_{\mathbf{x}}^*(\mathbf{x}(t), t)^T B(t)R^{-1}(t)B^T(t)J_{\mathbf{x}}^*(\mathbf{x}(t), t) + J_{\mathbf{x}}^*(\mathbf{x}(t), t)^T A(t)\mathbf{x}(t) \quad (20)$$

with boundary condition

$$J^*(\mathbf{x}(t_f), t_f) = \frac{1}{2}\mathbf{x}^T(t_f)Q_f\mathbf{x}(t_f). \quad (21)$$

It may appear as if we are stuck here, as this form of the HJB doesn't immediately yield  $J^*(\mathbf{x}(t), t)$ . Armed with the knowledge that the discrete time LQR problem has a quadratic cost-to-go, we will cross our fingers and guess a solution of the form

$$J^*(\mathbf{x}(t), t) = \frac{1}{2}\mathbf{x}^T(t)V(t)\mathbf{x}(t). \quad (22)$$

Substituting, we have

$$0 = \frac{1}{2}\mathbf{x}^T(t)\dot{V}(t)\mathbf{x}(t) + \frac{1}{2}\mathbf{x}^T(t)Q(t)\mathbf{x}(t) - \frac{1}{2}\mathbf{x}^T(t)V(t)B(t)R^{-1}(t)B^T(t)V(t)\mathbf{x}(t) + \mathbf{x}^T(t)V(t)A(t)\mathbf{x}(t) \quad (23)$$

Note that we will decompose

$$\mathbf{x}^T(t)V(t)A(t)\mathbf{x}(t) = \frac{1}{2}\mathbf{x}^T(t)V(t)A(t)\mathbf{x}(t) + \frac{1}{2}\mathbf{x}^T(t)A^T(t)V(t)\mathbf{x}(t) \quad (24)$$

which yields

$$0 = \frac{1}{2}\mathbf{x}^T(t) \left( \dot{V}(t) + Q(t) - V(t)B(t)R^{-1}(t)B^T(t)V(t) + V(t)A(t) + A^T(t)V(t) \right) \mathbf{x}(t). \quad (25)$$

This equation must hold for all  $\mathbf{x}(t)$ , so

$$-\dot{V}(t) = Q(t) - V(t)B(t)R^{-1}(t)B^T(t)V(t) + V(t)A(t) + A^T(t)V(t) \quad (26)$$

with boundary condition  $V(t_f) = Q_f$ .

Therefore, the HJB PDE has been reduced to a set of matrix ordinary differential equations (the Riccati equation). This is integrated backwards in time to find the full control policy as a function of time. Once we have found  $V(t)$ , the control policy is

$$\mathbf{u}^*(t) = -R^{-1}(t)B^T(t)V(t)\mathbf{x}(t). \quad (27)$$

Similarly to the discrete case, the feedback gains tend toward constant in the limit of the infinite horizon problem, under some technical assumptions.

## 3.2 Differential Games

We have so far addressed the case in which we aim to solve the optimal control problem for a single agent. We will now consider an adversarial game setting, in which there exists another player that aims to maximally harm the first agent. In particular, we will consider zero sum games in which the second agent aims to maximize the cost of the first agent. While the differential game setting is not restricted to this case — agents may have separate cost functions that partially interfere or aid each other — the zero-sum case lends itself to useful analytical tools.

### 3.2.1 Differential Games and Information Patterns

We consider the two player differential game with dynamics

$$\dot{\mathbf{x}}(t) = f(\mathbf{x}(t), \mathbf{u}(t), \mathbf{d}(t)) \quad (28)$$

where the first player takes action  $\mathbf{u}(t)$  at time  $t$ , and the second player takes action  $\mathbf{d}(t)$ . The state  $\mathbf{x}(t)$  is the joint state of both players. We write the cost as

$$J(\mathbf{x}(t)) = c_f(\mathbf{x}(0)) + \int_t^0 c(\mathbf{x}(\tau), \mathbf{u}(\tau), \mathbf{d}(\tau)) d\tau \quad (29)$$

which the first agent aims to maximize, and the second agent aims to minimize.

To fully specify the differential game, we must specify what each agent knows, and when. This is referred to as the *information pattern* of the game. In addition to capturing the knowledge of the state available to each agent, the information pattern also captures the knowledge of each other agents' strategies available to each agent.

### 3.2.2 Hamilton-Jacobi-Isaacs

The key idea in building the multi-agent equivalent of the HJB equation will again be to apply the principle of optimality. We consider the information pattern in which the adversary has access to the instantaneous control action of the first agent, so the cost takes the form

$$J(\mathbf{x}(t), t) = \min_{\Gamma(\mathbf{u})(\cdot)} \max_{\mathbf{u}(\cdot)} \left\{ \int_t^0 c(\mathbf{x}(\tau), \mathbf{u}(\tau), \mathbf{d}(\tau)) d\tau + c_f(\mathbf{x}(0)) \right\}. \quad (30)$$

Applying the dynamic programming principle, we have

$$J(\mathbf{x}(t), t) = \min_{\Gamma(\mathbf{u})(\cdot)} \max_{\mathbf{u}(\cdot)} \left\{ \int_t^{t+\Delta t} c(\mathbf{x}(\tau), \mathbf{u}(\tau), \mathbf{d}(\tau)) d\tau + J(\mathbf{x}(t + \Delta t), t + \Delta t) \right\}. \quad (31)$$

We can take the same strategy as with the informal derivation of the HJB equation, and Taylor expand both terms to yield

$$\begin{aligned} J(\mathbf{x}(t), t) = \min_{\Gamma(\mathbf{u})(\cdot)} \max_{\mathbf{u}(\cdot)} \{ & c(\mathbf{x}(t), \mathbf{u}(t), \mathbf{d}(t)) \Delta t + J(\mathbf{x}(t), t) \\ & + (J_{\mathbf{x}}(\mathbf{x}(t), t))^T f(\mathbf{x}(t), \mathbf{u}(t), \mathbf{d}(t)) \Delta t + J_t(\mathbf{x}(t), t) \Delta t \}. \end{aligned} \quad (32)$$

Note that we are optimizing over instantaneous actions, and so we are optimizing over finite dimensional quantities as opposed to functions. Dividing through by  $\Delta t$  and removing redundant terms, we get the *Hamilton-Jacobi-Isaacs* (HJI) equation

$$0 = J_t(\mathbf{x}, t) + \max_{\mathbf{u}} \min_{\mathbf{d}} \{c(\mathbf{x}, \mathbf{u}, \mathbf{d}) + (J_{\mathbf{x}}(\mathbf{x}, \mathbf{u}, \mathbf{d}))^T f(\mathbf{x}, \mathbf{d}, \mathbf{u})\} \quad (33)$$

with boundary condition

$$J(\mathbf{x}, 0) = c_f(\mathbf{x}). \quad (34)$$

Note that we have switched the order of the min/max.

### 3.2.3 Reachability

Differential games have applications in multi-agent modeling (both in the context of autonomous systems engineering and, e.g., economics and operations research). One concrete application in engineering is reachability analysis. In this setting, an agent aims to compute the set of states in which there exists a policy that either avoids a target set or enters a target set, subject to adversarial disturbances. The former case, in which we would like to avoid a target set, is useful for safety verification. If we are able to, even in the worst case, guarantee e.g. collision avoidance, we have guarantees on safety (subject of course to our system assumptions). The latter case is useful for task satisfaction. For example, we would like a quadrotor to reach a set of safe hovering poses, even under adversarial disturbances. Finding the backward reachable set in this case would find all states such that there exists a policy that succeeds in reaching the target set.

More concretely, the first case aims to find a set

$$\mathcal{A}(t) = \{\bar{\mathbf{x}} : \exists \Gamma(\mathbf{u})(\cdot), \forall \mathbf{u}(\cdot), \dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{u}, \mathbf{d}), \mathbf{x}(t) = \bar{\mathbf{x}}, \mathbf{x}(0) \in \mathcal{T}\} \quad (35)$$

where  $\mathcal{T}$  is the unsafe set which we aim to avoid. Breaking this down,  $\mathcal{A}(t)$  is the set of states at time  $t$  such that there exists  $\Gamma(\mathbf{u})$  that maps action  $\mathbf{u}$  to a disturbance such that, following the dynamics induced by the disturbance and the action sequence, the state is in  $\mathcal{T}$  at time 0 (note that we are considering  $0 \leq t$ ).

The second case aims to find a set

$$\mathcal{R}(t) = \{\bar{\mathbf{x}} : \forall \Gamma(\mathbf{u})(\cdot), \exists \mathbf{u}(\cdot), \dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{u}, \mathbf{d}), \mathbf{x}(t) = \bar{\mathbf{x}}, \mathbf{x}(0) \in \mathcal{T}\}, \quad (36)$$

where in this case  $\mathcal{T}$  is the set that we wish to reach. In this setting, we wish to find all states that, no matter what strategy the disturbance takes, there exist control actions that can steer the system to the goal state. Because the disturbance is adversarial (we reason over all adversary strategies), this is an extremely conservative form of safety analysis.

Computation of the backward reachable set results from solving a differential *game of kind* in which the outcome is Boolean (i.e. whether or not  $\mathbf{x}(0) \in \mathcal{T}$ ). This boolean outcome can be encoded by removing the running cost and choosing a particular form for the final cost. In particular, we can choose a final cost where

$$\mathbf{x} \in \mathcal{T} \iff c_f(\mathbf{x}) \leq 0. \quad (37)$$

As a result, the agent should aim to maximize  $c_f$  to avoid  $\mathcal{T}$ , whereas the disturbance should aim to minimize it. The two settings then take the following forms:

- Set avoidance:  $J(\mathbf{x}, t) = \min_{\Gamma(\mathbf{u})} \max_{\mathbf{u}} c_f(\mathbf{x}(0))$
- Set reaching:  $J(\mathbf{x}, t) = \max_{\Gamma(\mathbf{u})} \min_{\mathbf{u}} c_f(\mathbf{x}(0))$

**Sets vs. Tubes.** We have so far considered avoidance or reachability problems for which we care about set membership at time  $t = 0$ . However, for something like collision avoidance, we would like to stay collision free at every time as opposed to a particular time. *Backward reachable sets* capture the case in which only the final time set membership matters, and states for times  $t < 0$  do not matter. *Backward reachable tubes* capture the entire time duration of the problem. Any state that passes through the target at any time in the problem duration is included. This yields a modified value function of the form

$$J(\mathbf{x}, t) = \min_{\Gamma(\mathbf{u})} \max_{\mathbf{u}} \min_{\tau \in [t, 0]} c_f(\mathbf{x}(\tau)). \quad (38)$$

If the target set membership holds at any time  $\tau'$ , then  $\min_{\tau \in [t, 0]} c_f(\mathbf{x}(\tau)) \leq c_f(\mathbf{x}(\tau')) \leq 0$ .

### 3.3 Further Reading

Our coverage of reachability analysis is based on the [MBT05], which is an important early work in the field, in addition to being a relatively comprehensive coverage of the method. For a review of differential games with a (slight) emphasis on economics and management science, we refer the reader to [Bre10]. For a review of HJB and continuous time LQR, we refer the reader to [Ber12] and [Kir12].

## References

- [Ber12] Dimitri P Bertsekas. *Dynamic programming and optimal control*. Number 1. 4 edition, 2012.
- [Bre10] Alberto Bressan. Noncooperative differential games. a tutorial. 2010.
- [Kir12] Donald E Kirk. *Optimal control theory: an introduction*. Courier Corporation, 2012.
- [MBT05] Ian M Mitchell, Alexandre M Bayen, and Claire J Tomlin. A time-dependent hamilton-jacobi formulation of reachable sets for continuous dynamic games. *IEEE Transactions on Automatic Control*, 2005.