

Lecture Schedule

Dynamical programming

- 1 The finite-horizon decision problem
31 January
- 2 Dynamical Programming
7 February
- 3 DP reformulations and introduction to Control
14 February

Control

- 4 Discretization and PID control
21 February
- 5 Direct methods and control by optimization
28 February
- 6 **Linear-quadratic problems in control**
7 March
- 7 Linearization and iterative LQR
14 March

Syllabus: <https://02465material.pages.compute.dtu.dk/02465public>
Help improve lecture by giving feedback on DTU learn

Reinforcement learning

- 8 Exploration and Bandits
21 March
- 9 Policy and value iteration
4 April
- 10 Monte-carlo methods and TD learning
11 April
- 11 Model-Free Control with tabular and linear methods
18 April
- 12 Eligibility traces and value-function approximations
25 April
- 13 Q-learning and deep-Q learning
2 May

Reading material:

- [Her24, Chapter 16]

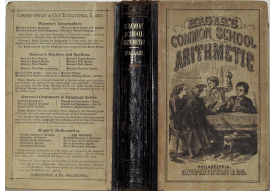
Learning Objectives

- Linear-quadratic regulator (LQR)
- Derivation of the LQR from DP
- Applications and variations

- Project evaluations will be ready in about a week
- Programming evaluations see <https://02465material.pages.compute.dtu.dk/02465public/projects/project1.html>
- Part 2:
 - Less programming
 - A bit more emphasis on linear algebra; don't be afraid to write short answers if they are correct.
 - Be inspired by existing examples

Recap

Useful linear algebra



- A matrix A is **positive semi-definite** if it is symmetric and $\mathbf{x}^T A \mathbf{x} \geq 0$ for all \mathbf{x}
 - This means A behaves like a positive number: $ax^2 \geq 0$.
- if A is a symmetric matrix then:

$$\frac{1}{2} \mathbf{x}^T A \mathbf{x} + \mathbf{b}^T \mathbf{x} = \frac{1}{2} (\mathbf{x} + A^{-1} \mathbf{b})^T A (\mathbf{x} + A^{-1} \mathbf{b}) - \frac{1}{2} \mathbf{b}^T A^{-1} \mathbf{b}$$

- This allows us to quickly find minimum

The Dynamical Programming algorithm

For every initial state x_0 , the optimal cost $J^*(x_0)$ is equal to $J_0(x_0)$, and optimal policy π^* is $\pi^* = \{\mu_0, \dots, \mu_{N-1}\}$, computed by the following algorithm, which proceeds backward in time from $k = N$ to $k = 0$ and for each $x_k \in S_k$ computes

$$J_N(x_N) = g_N(x_N) \quad (1)$$

$$J_k(x_k) = \min_{u_k \in \mathcal{A}_k(x_k)} \mathbb{E}_{w_k} \{g_k(x_k, u_k, w_k) + J_{k+1}(f_k(x_k, u_k, w_k))\} \quad (2)$$

$$\mu_k(x_k) = u_k^* \quad (u_k^* \text{ is the } u_k \text{ which minimizes the above expression}). \quad (3)$$

- For $k = 0, 1, \dots, N - 1$

$$\begin{aligned}x_{k+1} &= f_k(x_k, u_k, w_k) = A_k x_k + B_k u_k, \\g_k(x_k, u_k, w_k) &= \frac{1}{2} x_k^\top Q_k x_k + \frac{1}{2} u_k^\top R_k u_k, \\g_N(x_k) &= \frac{1}{2} x_N^\top Q_N x_N\end{aligned}$$

- **Note:** This is not the most general case, but will illustrate the main ideas

- Define $V_N \equiv Q_N$ and initialize:

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

- DP iteration (start at $k = N - 1$)

$$J_k(\mathbf{x}_k) = \min_{\mathbf{u}_k} \mathbb{E}_{w_k} \{g_k(\mathbf{x}_k, \mathbf{u}_k, w_k) + J_{k+1}(f_k(\mathbf{x}_k, \mathbf{u}_k, w_k))\}$$

- Remember to store optimal u_k^* as $\pi_k(x_k) = u_k^*$

DP solution gives the controller:

$$\textcircled{1} V_N = Q_N$$

$$\textcircled{2} L_k = -(R_k + B_k^T V_{k+1} B_k)^{-1} (B_k^T V_{k+1} A_k)$$

$$\textcircled{3} V_k = Q_k + L_k^T R_k L_k + (A_k + B_k L_k)^T V_{k+1} (A_k + B_k L_k)$$

$$\textcircled{4} \mathbf{u}_k^* = L_k \mathbf{x}_k$$

$$\textcircled{5} J_k^*(\mathbf{x}_k) = \frac{1}{2} \mathbf{x}_k^T V_k \mathbf{x}_k$$

Linear Quadratic Regulator

Double Integrator Example

- True dynamics

$$\dot{\mathbf{x}}(t) = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \mathbf{x}(t) + \begin{bmatrix} 0 \\ 1 \end{bmatrix} \mathbf{u}(t) \quad (4)$$

- **Euler discretization** using $\Delta = 1$ System evolves according to:

$$\mathbf{x}_{k+1} = \underbrace{\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}}_{=A} \mathbf{x}_k + \underbrace{\begin{bmatrix} 0 \\ 1 \end{bmatrix}}_{=B} \mathbf{u}_k$$

- Cost function:

$$J(\mathbf{x}_0) = \sum_{k=0}^N \frac{1}{2\rho} x_{k,1}^2 + \sum_{k=0}^{N-1} \frac{1}{2} u_k^2$$

- Can be put into standard form using matrices/start position:

$$Q_k = Q_N = \begin{bmatrix} \frac{1}{\rho} & 0 \\ 0 & 0 \end{bmatrix} \quad R = 1$$

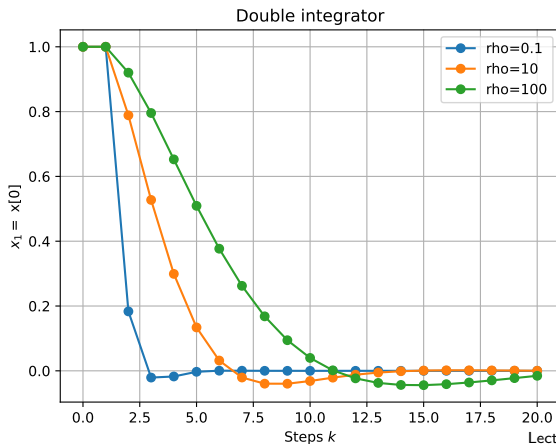
Linear Quadratic Regulator

Exponential integrator

- Apply discrete LQR
- Simulate starting in $\mathbf{x}_0 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ using policy

$$\pi_k(\mathbf{x}_k) = L_k \mathbf{x}_k$$

- What about the true system $\dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}, \mathbf{u})$?



Consider a (generic) LQR problem of the form:

$$\mathbf{x}_{k+1} = A\mathbf{x}_k + B\mathbf{u}_k \quad (5)$$

$$\text{cost} = \sum_{k=0}^{N-1} \frac{1}{2} \mathbf{x}_k^\top Q \mathbf{x}_k + \frac{1}{2} R_0 \mathbf{u}_k^\top \mathbf{u}_k \quad (6)$$

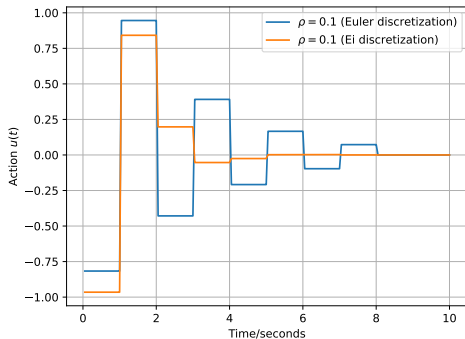
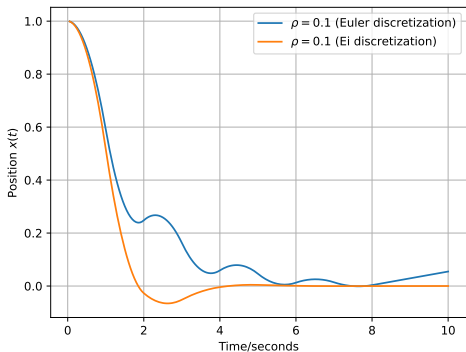
Where $R_0 > 0$ is a constant. After LQR, the controller selects actions using $\mathbf{u}_k = L_k \mathbf{x}_k$. What do you think typically happens with the matrix L_k when $R_0 \rightarrow \infty$ (**very big** R_0)

- a. The entries in L_k becomes very small, negative numbers
- b. The entries in L_k becomes very big, positive numbers
- c. It is not possible to say anything about the typical case
- d. The entries in L_k gets closer to zero
- e. Don't know.

Linear Quadratic Regulator

Double integrator example

- **Blue:** LQR using Euler $\mathbf{x}_{k+1} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \mathbf{x}_k + \begin{bmatrix} 0 \\ 1 \end{bmatrix} \mathbf{u}_k$
- **Red:** LQR using Exponential $\mathbf{x}_{k+1} = e^{A\Delta} \mathbf{x}_k + A^{-1} (e^{A\Delta} - I) B \mathbf{u}_k$



- LQR is optimal in discrete problem
- Discrete controller can be bad in real problem (always check!)
- Always use EI for linear dynamics

Linear Quadratic Regulator

Example: The locomotive



Steer locomotive (starting at $x = -1$) to goal ($x^* = 0$)

$$\ddot{x}(t) = \frac{1}{m}u(t) \quad (7)$$

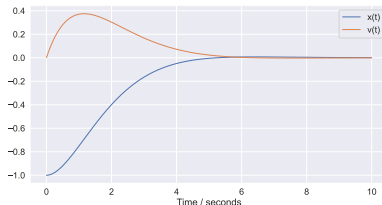
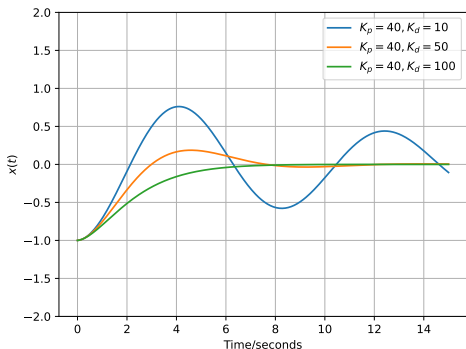
Can be re-written as:

$$\dot{\mathbf{x}} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \mathbf{x} + \begin{bmatrix} 0 \\ \frac{1}{m} \end{bmatrix} u \quad (8)$$



Discretized to $\mathbf{x}_{k+1} = A\mathbf{x}_k + B\mathbf{u}_k$.

$$e_k = x^* - x_k$$

$$u_k = e_k K_p + K_d \frac{e_k - e_{k-1}}{\Delta}$$



- Alternatively: Use a cost function $\sum_k \mathbf{x}_k^\top Q \mathbf{x}_k + \mathbf{u}_k^\top \mathbf{u}_k$ and use LQR!

 `lecture_04_pid_d.py`  `lecture_06_lqr_locomotive.py`

Recall LQR has the form:

$$1 \quad V_N = Q_N$$

$$2 \quad L_k = -(R_k + B_k^T V_{k+1} B_k)^{-1} (B_k^T V_{k+1} A_k)$$

$$3 \quad V_k = Q_k + L_k^T R_k L_k + (A_k + B_k L_k)^T V_{k+1} (A_k + B_k L_k)$$

$$4 \quad \mathbf{u}_k^* = L_k \mathbf{x}_k$$

$$5 \quad J_k^*(\mathbf{x}_k) = \frac{1}{2} \mathbf{x}_k^T V_k \mathbf{x}_k$$

- What happens if we repeat step 2 and 3 *many* times?
- The method will converge: $L_k \rightarrow L$
 - Select actions $\mathbf{u}_k = L \mathbf{x}_k$ ("**plan until convergence**")
- If you think about it, this corresponds to planning on $N \rightarrow \infty$ horizon.
- **This is quite popular in control theory; what we will do in RL.**

- The cost term $\frac{1}{2}\mathbf{x}^\top Q\mathbf{x} + \frac{1}{2}\mathbf{u}^\top R\mathbf{u}$ is **smallest** when $\mathbf{x} = \mathbf{u} = \mathbf{0}$
- Implies that LQR will control system to state $\mathbf{x} = \mathbf{u} = \mathbf{0}$
- Suppose we want to drive system towards $\mathbf{x}_g, \mathbf{u}_g$?
 - Use $c(\mathbf{x}, \mathbf{u}) = \frac{1}{2}(\mathbf{x} - \mathbf{x}_g)^\top Q(\mathbf{x} - \mathbf{x}_g) + \frac{1}{2}(\mathbf{u} - \mathbf{u}_g)^\top R(\mathbf{u} - \mathbf{u}_g)$
- more generally assume

$$c_k(\mathbf{x}_k, \mathbf{u}_k) = \frac{1}{2}\mathbf{x}_k^\top Q_k \mathbf{x}_k + \frac{1}{2}\mathbf{u}_k^\top R_k \mathbf{u}_k + \mathbf{u}_k^\top H_k \mathbf{x}_k + \mathbf{q}_k^\top \mathbf{x}_k + \mathbf{r}_k^\top \mathbf{u}_k + q_k \quad (9)$$

$$c_N(\mathbf{x}_k) = \frac{1}{2}\mathbf{x}_k^\top Q_N \mathbf{x}_k + \mathbf{q}_N^\top \mathbf{x}_k + q_N \quad (10)$$

and dynamics

$$\mathbf{x}_{k+1} = A_k \mathbf{x}_k + B_k \mathbf{u}_k + \mathbf{d}_k$$

Linear Quadratic Regulator

General discrete LQR algorithm

How to start living in luxury and never work again!

$$\cdots (V_{k+1} + \mu I) \cdots$$

1. $V_N = Q_N; \mathbf{v}_N = \mathbf{q}_N; v_N = q_N$

2.
$$L_k = -S_{\mathbf{u}\mathbf{u},k}^{-1} S_{\mathbf{u}\mathbf{x},k} \quad S_{\mathbf{u},k} = \mathbf{r}_k + B_k^T \mathbf{v}_{k+1} + B_k^T V_{k+1} \mathbf{d}_k$$

$$\mathbf{l}_k = -S_{\mathbf{u}\mathbf{u},k}^{-1} S_{\mathbf{u},k} \quad S_{\mathbf{u}\mathbf{u},k} = R_k + B_k^T V_{k+1} B_k$$

$$S_{\mathbf{u}\mathbf{x},k} = H_k + B_k^T V_{k+1} A_k.$$

3.
$$V_k = Q_k + A_k^T V_{k+1} A_k - L_k^T S_{\mathbf{u}\mathbf{u},k} L_k$$

$$\mathbf{v}_k = \mathbf{q}_k + A_k^T (\mathbf{v}_{k+1} + V_{k+1} \mathbf{d}_k) + S_{\mathbf{u}\mathbf{x},k}^T \mathbf{l}_k$$

$$v_k = v_{k+1} + q_k + \mathbf{d}_k^T \mathbf{v}_{k+1} + \frac{1}{2} \mathbf{d}_k^T V_{k+1} \mathbf{d}_k + \frac{1}{2} \mathbf{l}_k^T S_{\mathbf{u},k}$$

4. $\mathbf{u}_k^* = \mathbf{l}_k + L_k \mathbf{x}_k$

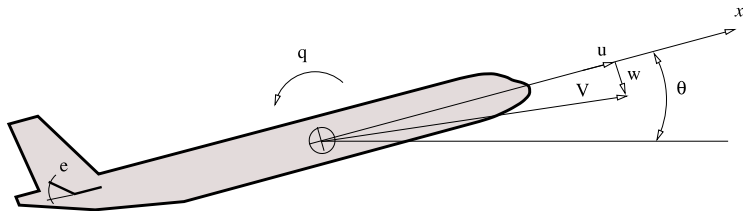
5. $J_k(\mathbf{x}_k) = \frac{1}{2} \mathbf{x}_k^T V_k \mathbf{x}_k + \mathbf{v}_k^T \mathbf{x}_k + v_k.$

Doctors hate this one weird trick!

$$V_k \leftarrow \frac{1}{2} (V_k^T + V_k)$$

Linear Quadratic Regulator

Boing 747 Example



$$\begin{bmatrix} \dot{u} \\ \dot{w} \\ \dot{q} \\ \dot{\theta} \end{bmatrix} = \underbrace{\begin{bmatrix} -0.003 & 0.039 & 0. & -0.322 \\ -0.065 & -0.319 & 7.74 & 0. \\ 0.02 & -0.101 & -0.429 & 0. \\ 0. & 0. & 1. & 0. \end{bmatrix}}_A \underbrace{\begin{bmatrix} u - u_w \\ w - w_w \\ q \\ \theta \end{bmatrix}}_x + \underbrace{\begin{bmatrix} 0.01 & 1. \\ -0.18 & -0.04 \\ -1.16 & 0.598 \\ 0. & 0. \end{bmatrix}}_B \underbrace{\begin{bmatrix} e \\ t \end{bmatrix}}_u$$

$$\begin{bmatrix} y_1(t) \\ y_2(t) \end{bmatrix} = \underbrace{\begin{bmatrix} 1. & 0. & 0. & 0. \\ 0. & -1. & 0. & 7.74 \end{bmatrix}}_{=P} \begin{bmatrix} u(t) - u_w(t) \\ w(t) - w_w(t) \\ q(t) \\ \theta(t) \end{bmatrix}$$

- y_1 and y_2 corresponds to the airspeed and climb rate.

- **Start:** $x = 0$ (steady flight)

- **Want:** Air speed of 10: $y^* = \begin{bmatrix} 10 \\ 0 \end{bmatrix}$

- Write dynamics as $\dot{\mathbf{x}} = A\mathbf{x} + B\mathbf{u}$
- Introduce cost function:

$$\int_0^{t_F} \left(\frac{1}{2}(\mathbf{y} - \mathbf{y}^*)^\top (\mathbf{y} - \mathbf{y}^*) + \frac{1}{2}\mathbf{u}^\top \mathbf{u} \right) dt$$

- Discretize dynamics using Exponential Integration to get $\mathbf{x}_{k+1} = \bar{A}\mathbf{x}_k + \bar{B}\mathbf{u}_k$
- Discretize cost to get one of the form

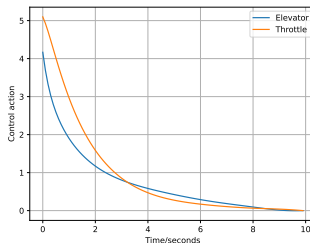
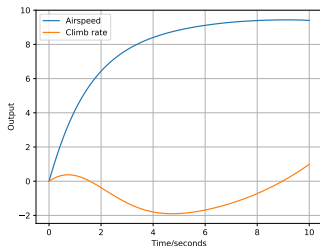
$$\sum_{k=0}^{\infty} \frac{1}{2}\mathbf{x}_k^\top Q\mathbf{x}_k + \mathbf{q}\mathbf{x}_k + q_0 + \frac{1}{2}\mathbf{u}_k^\top R\mathbf{u}_k$$

- Apply LQR!

Linear Quadratic Regulator

Outcome and a Quiz

- Control law $\mathbf{u}_k = L\mathbf{x}_k$



Left: airspeed and climb rate. **Right:** Elevator and throttle

Why does the output adjust quickly but fail to get entirely to the goal y^* ?

- Something bad happened to the dynamics with the exponential integration
- The explanation has to do with planning on a finite horizon
- The explanation is that R in $\mathbf{u}_k^\top R \mathbf{u}_k$ should be bigger

d. Don't know.

- Consider the case where there is additive Gaussian noise:



$$\mathbf{x}_{k+1} = A_k \mathbf{x}_k + B_k \mathbf{u}_k + \boldsymbol{\omega}_k$$

- We can still solve the problem, and (amazingly!) the noise has **no influence** on the control law

$$\mathbf{u}_k = L_k \mathbf{x}_k$$

- LQR is robust to noise

- Stability/controllability of LQR?
 - **Important subject which we ignore**
- What if matrices A_k, B_k are random?
 - **This too can be solved[Ber05, Chapter 4]**
- What about partial observation?
 - **I.e. assume we observe $o_k = D_k x_k$ [Ber05, Chapter 4]**
- What about constraints? What if we know $u_L \leq \mathbf{u}_k \leq u_B$?
- Euler integration is often not ideal.
 - **Alternatives including error analysis**

-  D.P. Bertsekas.
Dynamic Programming and Optimal Control.
Number v. 1 in Athena Scientific optimization and computation series.
Athena Scientific, 2005.
-  Tue Herlau.
Sequential decision making.
(Freely available online), 2024.