

CMU 16-745 最优控制笔记

吕翔宇 (Xiangyu Lyu)

lvxiangyu11@gmail.com

TU Darmstadt, Germany

GitHub: lvxiangyu11

2025 年 5 月 9 日

目录

1	引言	3
2	动力学简介	4
2.1	连续时间的动力学系统	4
2.1.1	控制映射系统	6
2.1.2	Manipulator动力学	6
2.1.3	线性系统	7
2.1.4	平衡点 equilibria	7
2.1.5	平衡点的稳定性分析	8
2.2	离散时间的动力学系统	10
2.2.1	正向欧拉法	10
2.2.2	RK4方法	14
2.2.3	反向欧拉法	21
2.2.4	离散控制输入	22
3	数值优化基础	24
3.1	符号声明	24
3.2	Root Finding问题	25
3.2.1	Newton 方法	26
3.2.2	牛顿方法在最小化问题中的应用	29
3.3	最小值存在的充分条件	29
3.4	正则化 (Regularization) 与强凸性	30
	Appendix	32

1 引言

这里是 CMU 16-745 最优控制课程的笔记。我将涵盖课程中的重要概念、定理和例子。课程目标

1. 分析动力系统的稳定性。
2. 设计稳定平衡和轨迹的LQR控制器。
3. 使用离线轨迹优化设计非线性系统的轨迹。
4. 使用在线凸优化实现模型预测控制。
5. 了解随机性和模型不稳定性的影响。
6. 无最优模型时直接优化反馈策略。

写作风格备注：(TODO: 终稿注释掉!) 额外的公式推导使用浅绿色背景框额外的例子使用浅蓝色背景框代码使用浅灰色背景

2 动力学简介

这一节面向不了解控制学、动力学的读者，以一个单摆系统为例，快速的介绍连续和离散动力学系统的基础，包括稳定性分析、线性系统和非线性系统。

2.1 连续时间的动力学系统

Continuous-Time Dynamics System (CTDS) 是一个描述系统状态随时间变化的数学模型。它通常由以下方程表示：

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

其中， $x(t) \in \mathbb{R}^n$ 是系统状态， $u(t) \in \mathbb{R}^m$ 是控制输入， f 是一个描述系统的动态函数，展示了系统如何根据 u 而演化。

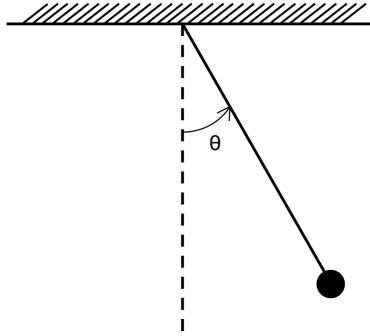


图 1: 简单单摆的物理模型。图中展示了摆球的受力情况，包括重力 mg 和拉力 \vec{T} ，以及摆球的运动方向和角度 θ 的关系。

对于1中的单摆系统，其状态 x 可以用角度 q 和速度 v 来表示。我们可以将其状态表示为：

$$x = \begin{bmatrix} q \\ v \end{bmatrix} \quad (2)$$

q 并不总是一个vector。configuration并不必须是一个vector，并且速度并不必须是configuration的导数。并且，仅在对系统的动力学描述是平滑的时候，才可以说这个系统的描述是有效的。

连续动态系统的简单单摆例子 考虑1中的单摆系统。我们可以用以下方程来描述它的动力学：

$$ml^2\ddot{\theta} + mgl\sin(\theta) = \tau \quad (3)$$

其推导过程如下：

从拉格朗日方程推导动力学方程推导* 在这一节中，我们使用拉格朗日方法推导单摆系统的动力学方程。首先，单摆的动能为：

$$T = \frac{1}{2}ml^2\dot{\theta}^2$$

势能为：

$$V = mgl(1 - \cos \theta)$$

其中， m 是质量， l 是摆长， θ 是角度， $\dot{\theta}$ 是角速度， g 是重力加速度。拉格朗日量 L 为动能 T 与势能 V 之差，即：

$$L = T - V = \frac{1}{2}ml^2\dot{\theta}^2 - mgl(1 - \cos \theta)$$

拉格朗日方程的一般形式为：

$$\frac{d}{dt} \left(\frac{\partial L}{\partial \dot{\theta}} \right) - \frac{\partial L}{\partial \theta} = 0$$

我们计算拉格朗日量的偏导数：

1. 对于 $\frac{\partial L}{\partial \dot{\theta}}$ ，有：

$$\frac{\partial L}{\partial \dot{\theta}} = ml^2\dot{\theta}$$

2. 对于 $\frac{\partial L}{\partial \theta}$ ，有：

$$\frac{\partial L}{\partial \theta} = mgl \sin \theta$$

将这些代入拉格朗日方程，得到：

$$\frac{d}{dt} (ml^2\dot{\theta}) - mgl \sin \theta = 0$$

即：

$$ml^2\ddot{\theta} + mgl \sin \theta = 0$$

如果引入控制输入 τ ，方程变为：

$$ml^2\ddot{\theta} + mgl \sin \theta = \tau$$

这个方程描述了单摆的动力学行为，其中 τ 是控制输入，能够影响单摆的角加速度。最终，我们得到了单摆的动力学方程：

$$ml^2\ddot{\theta} + mgl \sin \theta = \tau$$

则系统的配置 q 可以通过单摆的角度 θ 来表示，速度 v 可以通过角速度 $\dot{\theta}$ 来表示。并且考虑到控制量 u ，通过施加到系统中的力矩(torque) τ 来控制单摆的运动。则状态可

以表示为:

$$x = \begin{bmatrix} \theta \\ \dot{\theta} \end{bmatrix} \quad (4)$$

则其速度 \dot{x} 可以表示为:

$$\dot{x} = \begin{bmatrix} \dot{\theta} \\ \ddot{\theta} \end{bmatrix} = \begin{bmatrix} \dot{\theta} \\ -\frac{g}{l}\sin(\theta) + \frac{\tau}{ml^2} \end{bmatrix} \quad (5)$$

则系统动态方程可以写为 $f(x, u)$:

$$\dot{x} = f(x, u) = \begin{bmatrix} \dot{\theta} \\ -\frac{g}{l}\sin(\theta) + \frac{\tau}{ml^2} \end{bmatrix} \quad (6)$$

则这里状态是一个流形(manifold)的函数 $x \in \mathbb{S}^1 \times \mathbb{R}$ 是一个圆柱, 表示了系统的状态随时间的变化。

2.1.1 控制映射系统

很多系统可以定义为一个控制映射系统。这是一个上述动力系统的特殊形式, 其中控制输入通过一个映射矩阵 G 来影响系统。

$$\dot{x} = f_0(x) + G(x)u \quad (7)$$

对于上述简单单摆模型, 可以将其表示为:

$$\dot{x} = \begin{bmatrix} \dot{\theta} \\ -\frac{g}{l}\sin(\theta) + \frac{\tau}{ml^2} \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \theta \\ \dot{\theta} \end{bmatrix} + \begin{bmatrix} 0 \\ \frac{1}{ml^2} \end{bmatrix} \tau \quad (8)$$

2.1.2 Manipulator动力学

Manipulator是一个可以在空间中移动的物体。它的动力学可以用以下方程表示 (参考《现代机器人学机构、规划与控制》的第八章, 读者感兴趣可以自行阅读, 此处略):

$$M(q)\dot{v} + C(q, v) = B(q)u \quad (9)$$

其中, $M(q)$ 是质量矩阵, $C(q, v)$ 是科里奥利力矩阵和重力项, $B(q)$ 是控制输入矩阵。这个方程描述了Manipulator的动力学行为, 其中 u 是控制输入, 能够影响Manipulator的运动。则configuration的变化可以表示为:

$$\dot{q} = G(q)v \quad (10)$$

其中, $G(q)$ 是一个映射矩阵, 将速度 v 映射到configuration的变化 \dot{q} 。则系统的运动学可以描述为:

$$\dot{x} = f(x, y) = \begin{bmatrix} G(q)v \\ -M(q)^{-1}(B(q) - C) \end{bmatrix} \quad (11)$$

在简单单摆系统中 $M(q) = ml^2$, $C(q, v) = mgl \sin(\theta)$, $B = I$, $G = I$ 。

所有机械系统都可以表述为这个形式，这是因为这个形式是一个描述Euler-Lagerange方程的不同形式(Kinetic energy - potential energy $L = T - V$)。

$$L = \frac{v^T M(q) v}{2} - V(q) \quad (12)$$

2.1.3 线性系统

线性系统（Linear System）是表述控制问题和设计控制器的通用表述方法。我们知道可以相对简单的解决一个线性系统，线性系统描述为：

$$\dot{x} = A(t)x + B(t)u \quad (13)$$

其中， $A(t)$ 和 $B(t)$ 是时间变化的矩阵。线性系统的解可以通过求解常微分方程来获得。线性系统的稳定性分析通常使用特征值和特征向量的方法。线性系统的控制器设计通常使用状态反馈和观测器设计的方法。

我们通常通过在现在状态附近的线性化来近似线性化非线性系统：

$$\dot{x} = f(x, u) \quad (14)$$

其中 $A = \frac{\partial f}{\partial x}$, $B = \frac{\partial f}{\partial u}$

2.1.4 平衡点 equilibria

equilibria指的是系统将会保持的位置，即一阶稳态。

$$\dot{x} = f(x, u) = 0 \quad (15)$$

在代数层面，这表示动力方程的根。对于简单单摆系统，平衡点可以表示为：

$$\dot{x} = f(x, u) = \begin{bmatrix} \dot{\theta} \\ -\frac{g}{l} \sin(\theta) + \frac{\tau}{ml^2} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad (16)$$

平衡点的求解 考虑一个简单的问题，我们如何找到一个平衡点？我们可以通过求解以下方程来找到平衡点：

$$\dot{x} = f(x, u) = \begin{bmatrix} \dot{\theta} \\ -\frac{g}{l} \sin(\theta) + \frac{\tau}{ml^2} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Rightarrow \frac{1}{ml^2} u = \frac{g}{l} \sin(\theta) \Rightarrow u = mgl \sin(\theta) \quad (17)$$

带入 $\theta = \pi/2$, 得到 $u = mgl$ 一般来说，可以将寻找平衡点的问题转化为求解一个非线性方程组的问题。

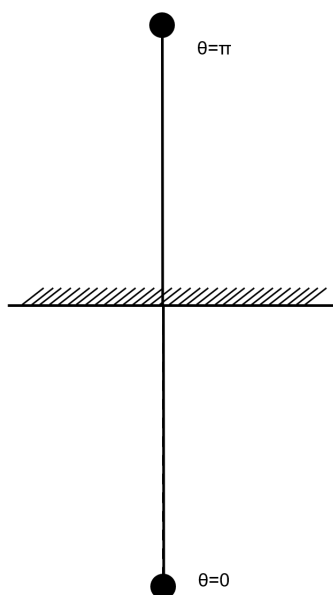


图 2: 图中展示了单摆系统的平衡点。平衡点是指系统在没有外部扰动的情况下, 能够保持静止或匀速运动的状态。

2.1.5 平衡点的稳定性分析

在知道了如何求平衡点后, 我们需要知道, 如果对系统施加一个小的扰动, 系统是否会回到平衡点。考虑一个1维系统, 在这个例子中, 当 $\frac{\partial f}{\partial x} < 0$ 时系统是稳定的。当 $\frac{\partial f}{\partial x} > 0$ 时系统是不稳定的。我们可以通过线性化来分析系统的稳定性。对于一个线性系统, 我们可以通过求解特征值来判断系统的稳定性。对于一个非线性系统, 我们可以通过求解雅可比矩阵的特征值来判断系统的稳定性, 因为从这个点越推越远。这个 $\frac{\partial f}{\partial x} < 0$ 的区域是一个吸引子(Attractor)又称basin of attraction of system(系统吸引盆地), 而 $\frac{\partial f}{\partial x} > 0$ 的区域是一个排斥子(Repeller)。

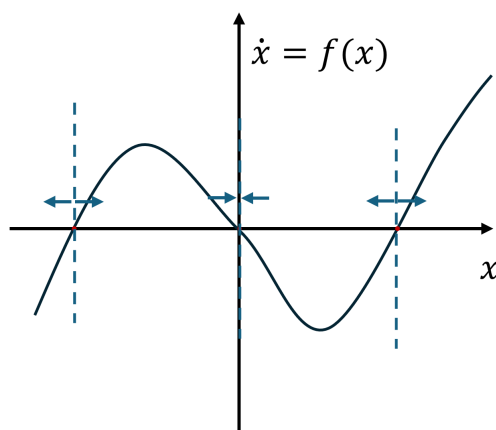


图 3: 图中展示了1维系统的稳定性分析。图中对于 $\dot{x} = 0$ 有三个解, 也就是说有三个平衡点, 但是只有中间的点满足 $\frac{\partial f}{\partial x} < 0$, 而左右两侧 $f(x)$ 对 x 的偏导都是大于0, 系统下一步将采取正向的行动, 导致系统不稳定。

推广到高维系统，这个 $\frac{\partial f}{\partial x} < 0$ 则为是一个 Jacobian matrix。为了研究系统的稳定性，我们可以对 Jacobian matrix 进行特征值分解。对于一个线性系统，我们可以通过求解特征值来判断系统的稳定性。即，如果每个特征值的实部都小于0，则系统是渐近稳定的；如果有一个特征值的实部大于0，则系统是不稳定的；如果所有特征值的实部都等于0，则系统是边界稳定的。可以用公式描述为：

$$Re(eig(\frac{\partial f}{\partial x})) = \begin{cases} < 0 & \text{渐近稳定} \\ > 0 & \text{不稳定} \\ = 0 & \text{边界稳定} \end{cases} \quad (18)$$

Jacobian矩阵定义为：

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \frac{\partial f_n}{\partial x_2} & \dots & \frac{\partial f_n}{\partial x_n} \end{bmatrix} \quad (19)$$

在上述简单单摆系统中

$$f(x) = \begin{bmatrix} \dot{\theta} \\ -\frac{g}{l}\sin(\theta) \end{bmatrix} \Rightarrow \frac{\partial f}{\partial x} = \begin{bmatrix} 0 & 1 \\ -\frac{g}{l}\cos(\theta) & 0 \end{bmatrix} \quad (20)$$

考虑平衡点 $\theta = 0$ ，我们可以得到 Jacobian 矩阵的特征值：

$$\frac{\partial f}{\partial x} = \begin{bmatrix} 0 & 1 \\ -\frac{g}{l} & 0 \end{bmatrix} \Rightarrow \lambda^2 + \frac{g}{l} = 0 \Rightarrow \lambda = \pm i\sqrt{\frac{g}{l}} \quad (21)$$

由于有一个特征值的实部大于0，表示系统在该点是一个鞍点，即系统在该点附近的运动是发散的。

考虑平衡点 $\theta = \pi$ ，我们可以得到 Jacobian 矩阵的特征值：

$$\frac{\partial f}{\partial x} = \begin{bmatrix} 0 & 1 \\ -\frac{g}{l} & 0 \end{bmatrix} \Rightarrow \lambda^2 - \frac{g}{l} = 0 \Rightarrow \lambda = 0 \pm \sqrt{\frac{g}{l}} \quad (22)$$

该点特征值实部为0，表示系统在该点是一个中心点，即系统在该点附近的运动是周期性的（不考虑空气摩擦力）。对于纯虚数特征值系统，这个系统叫做边缘稳定（marginally stable）。可以添加damping来使得系统稳定，如 $u = -K_d\dot{\theta}$ （摩擦力）。

2.2 离散时间的动力学系统

在这一节中，我们将介绍

- 连续的常微分方程（Ordinary Differential Equations, ODEs）
- 离散动力系统稳定性分析

在通常情况下，我们无法直接从 $\dot{x} = f(x)$ 中求出 $x(t)$ 的解析解，则无法直接求出系统的状态随时间的变化。我们可以通过在离散时间内通过数值的方法近似表示 $x(t)$ 。并且，离散时间模型可以捕获连续ODEs无法捕获的动态行为（比如，离散时间模型可以捕获系统的周期性行为，接触）。

2.2.1 正向欧拉法

离散时间动力学系统（Discrete-Time Dynamics System, DTDS）是一个描述系统状态随时间变化的数学模型。通常用于模拟计算，如模拟器等。它通常由以下方程表示：

$$x_{k+1} = f_{discrete}(x_k, u_k) \quad (23)$$

正向欧拉法* 这里将介绍正向欧拉法，并举一个例子。

用一个最简单的离散化表示，即正向欧拉法（Forward Euler Method）：

$$x_{k+1} = x_k + hf_{continuous}(x_k, u_k) \quad (24)$$

其中， h 是离散化的时间步长。这个方程表示在时间步长 h 内，系统状态 x_k 和控制输入 u_k 的关系。 $f_{continuous} = \frac{\partial f(x,u)}{\partial t}$

手算例子： 考虑一个简单的线性系统 $\dot{x} = -2x$ ，其离散化形式为：

$$x_{k+1} = x_k + h(-2x_k) = (1 - 2h)x_k \quad (25)$$

假设初始条件为 $x_0 = 1$ 且步长 $h = 0.1$ ，则计算前几个离散时间点的状态：

- 对于 $k = 0$ ， $x_0 = 1$ ；
- 对于 $k = 1$ ， $x_1 = (1 - 2 \times 0.1) \times 1 = 0.8$ ；
- 对于 $k = 2$ ， $x_2 = (1 - 2 \times 0.1) \times 0.8 = 0.64$ ；
- 对于 $k = 3$ ， $x_3 = (1 - 2 \times 0.1) \times 0.64 = 0.512$ 。

我们可以看到，系统状态随着时间步的增加逐渐减小。

用一个最简单的离散化表示，即正向欧拉法（Forward Euler Method）：

$$x_{k+1} = x_k + hf_{continuous}(x_k, u_k) \quad (26)$$

其中， h 是离散化的时间步长。这个方程表示在时间步长 h 内，系统状态 x_k 和控制输入 u_k 的关系。 $f_{continuous} = \frac{\partial f(x,u)}{\partial t}$

离散时间动力学系统的简单单摆例子* 继续单摆系统

单摆的连续时间动力学系统可以通过以下方程描述：

$$\dot{x} = f(x, u) = \begin{bmatrix} \dot{\theta} \\ \ddot{\theta} \end{bmatrix} = \begin{bmatrix} \dot{\theta} \\ -\frac{g}{l} \sin(\theta) + \frac{\tau}{ml^2} \end{bmatrix}$$

其中： θ 为单摆的角度（配置变量）， $\dot{\theta}$ 为角速度（速度变量）， $\ddot{\theta}$ 为角加速度。对于离散时间动力学系统，我们采用正向欧拉法进行离散化。设定离散时间步长为 h ，那么在每个时间步内，系统的状态 $x_k = [\theta_k, \dot{\theta}_k]$ 通过以下公式更新：

$$x_{k+1} = x_k + h \cdot f(x_k, u_k)$$

将连续时间的动力学方程 $f(x, u)$ 代入上式，我们得到：

$$x_{k+1} = \begin{bmatrix} \theta_k \\ \dot{\theta}_k \end{bmatrix} + h \cdot \begin{bmatrix} \dot{\theta}_k \\ -\frac{g}{l} \sin(\theta_k) + \frac{\tau_k}{ml^2} \end{bmatrix}$$

这个递推公式可以用来计算每个时间步的状态。

具体而言，我们可以分别更新角度和角速度：

$$\theta_{k+1} = \theta_k + h \cdot \dot{\theta}_k$$

$$\dot{\theta}_{k+1} = \dot{\theta}_k + h \left(-\frac{g}{l} \sin(\theta_k) + \frac{\tau_k}{ml^2} \right)$$

这个离散化公式描述了如何在每个时间步内，根据当前的角度和角速度，通过正向欧拉法更新单摆的状态。

针对简单单摆系统，设置 $l = m = 1, h = 0.1 \text{ or } 0.01$ 将会爆炸（blow up）。

以下实现了单摆系统的离散时间模拟代码，缩短模拟步长可以看到blow up被延缓了，但是仍会blow up。

```

1 # [discrete_time_dynamic_forward_euler.py]
2 def pendulum_dynamics(x):
3     l = 1.0
4     g = 9.81
5
6     theta = x[0] # 角度
7     theta_dot = x[1] # 角速度
8     theta_ddot = -(g/l) * np.sin(theta) # 角加速度，由重力引起的
9
10    return np.array([theta_dot, theta_ddot])
11
12 def pendulum_forward_euler(fun, x0, Tf, h):
13     t = np.arange(0, Tf + h, h)
```

```

14
15     x_hist = np.zeros((len(x0), len(t))) # 状态历史记录
16     x_hist[:, 0] = x0 # 初始状态
17
18     for k in range(len(t) - 1): # 迭代计算每个时间步的状态
19         x_hist[:, k+1] = x_hist[:, k] + h * fun(x_hist[:, k]) # 使用欧拉法更新状态
20     return x_hist, t
21
22 if __name__ == "__main__":
23     # 初始条件和参数
24     x0 = np.array([0.1, 0]) # 初始角度和角速度
25     Tf = 10.0 # 模拟总时间(秒)
26     h = 0.01 # 时间步长
27
28     # 计算单摆运动
29     x_hist, t_hist = pendulum_forward_euler(pendulum_dynamics, x0, Tf, h)
30
31     # 创建可视化对象并显示动画
32     visualizer = PendulumVisualizer(pendulum_length=1.0)
33     visualizer.visualize(x_hist, t_hist)

```

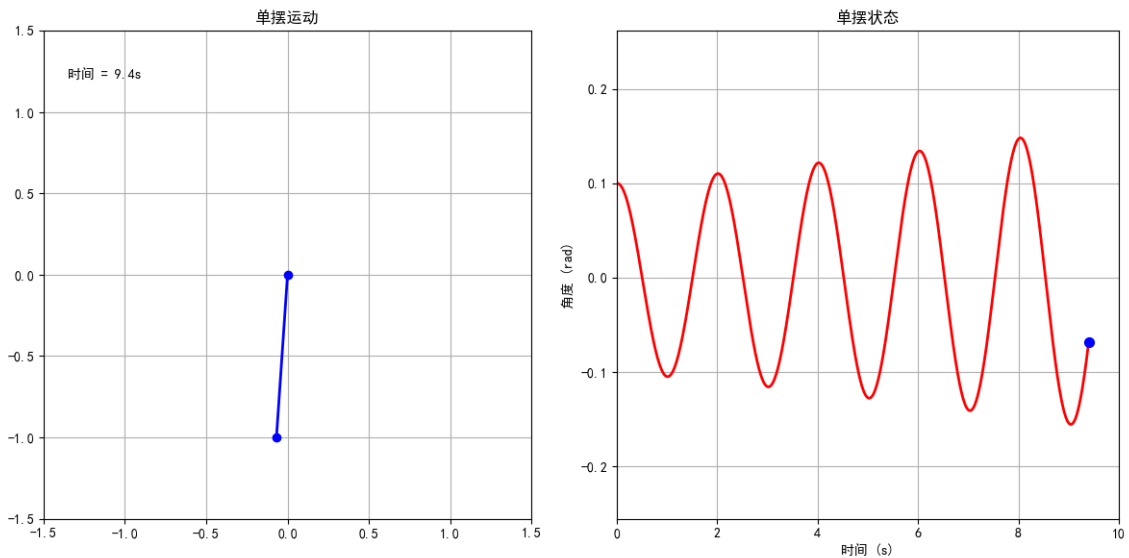


图 4: 图中展示了欧拉方法模拟的单摆系统的运动轨迹。可以看到，随着时间的推移，单摆逐渐blow up，说明该方法是数值不稳定的

正向欧拉法的稳定性分析 正向欧拉法的稳定性分析可以通过考虑离散时间系统的特征值来进行。对于线性系统，正向欧拉法的稳定性取决于离散化步长 h 和系统矩阵 A 的特征值。具体而言，如果所有特征值的模长小于1，则系统是渐近稳定的；如果有一个特征值的模长大于1，则系统是不稳定的。回顾之前判断系统是否稳定， $Re(eig(\frac{\partial f}{\partial x}))$ 与0的关系。对于离散时间，动力学系统可以迭代映射为：

$$x_N = f_d(f_d(f_d(\dots f_d(x_0))\dots)) \quad (27)$$

考虑线性和chain rule，有：

$$\frac{\partial x_N}{\partial x_0} = \frac{\partial f_d}{\partial x_0} \frac{\partial f_d}{\partial x_1} \cdots \frac{\partial f_d}{\partial x_{N-1}} \bigg|_{x_0} = A_d^N \quad (28)$$

其中 A_d 表示线性化的点的偏导 $\frac{\partial f_d}{\partial x}$ ，并且在迭代中，每个点的偏导数都是相同的。

稳定点是 $x = 0$ （可以通过修改坐标系达成），系统的稳定性意味着：

$$\lim_{N \rightarrow \infty} A_d^N x_0 = 0, \forall x_0 \in \mathbb{R}^n \Rightarrow \lim_{N \rightarrow \infty} A_d^N = 0 \Rightarrow |eig(A_d)| < 1 \quad (29)$$

如果 $|eig(A_d)| > 1$ ，则系统不稳定，迭代之后那些大于1的特征值将会占据主导地位，导致系统不稳定。

等价的， $eig(A_d)$ 必须在复数平面上位于单位圆内。

下面开始分析简单单摆系统在欧拉方法离散化后的稳定性， h 取0.1, $g=9.81$, $l=1.0$ 。

:

$$A_d = \frac{f_d}{x_k} = I + h A_{\text{continuous}} = I + h \begin{bmatrix} 0 & 1 \\ -\frac{g}{l} \cos(\theta) & 0 \end{bmatrix} \quad (30)$$

$$\Rightarrow eig(A_d|_{\theta=0}) = 1 \pm 0.313i \Rightarrow |eig(A_d)| = 1.313 > 1 \Rightarrow \text{不稳定} \quad (31)$$

我们可以将 h 和eigenvalue的关系绘制成图表，其代码如下：

```
1 # [euler_eigenvalue_stability.py]
2 # 此处仅展示关键代码
3
4 h_values = np.linspace(0.01, 0.5, 100)
5 eigenvalues = []
6
7 for h in h_values:
8     A_d = A_discrete(h, theta)
9     eigs = eigvals(A_d)
10    eigenvalues.append(eigs)
11
12 # 转换为numpy数组便于处理
13 eigenvalues = np.array(eigenvalues)
14
15 # 计算特征值的模
16 magnitudes = np.abs(eigenvalues)
```

可以从5中看出

- 欧拉法对简单单摆系统的离散化在任何正的时间步长下都是不稳定的，因为特征值的模总是大于1。
- 时间步长越大，特征值的模越大，系统越不稳定。
- 系统的两个特征值在实部上逐渐分离
- 所有点都在单位圆外，表明系统在所有正的 h 值下都是不稳定的

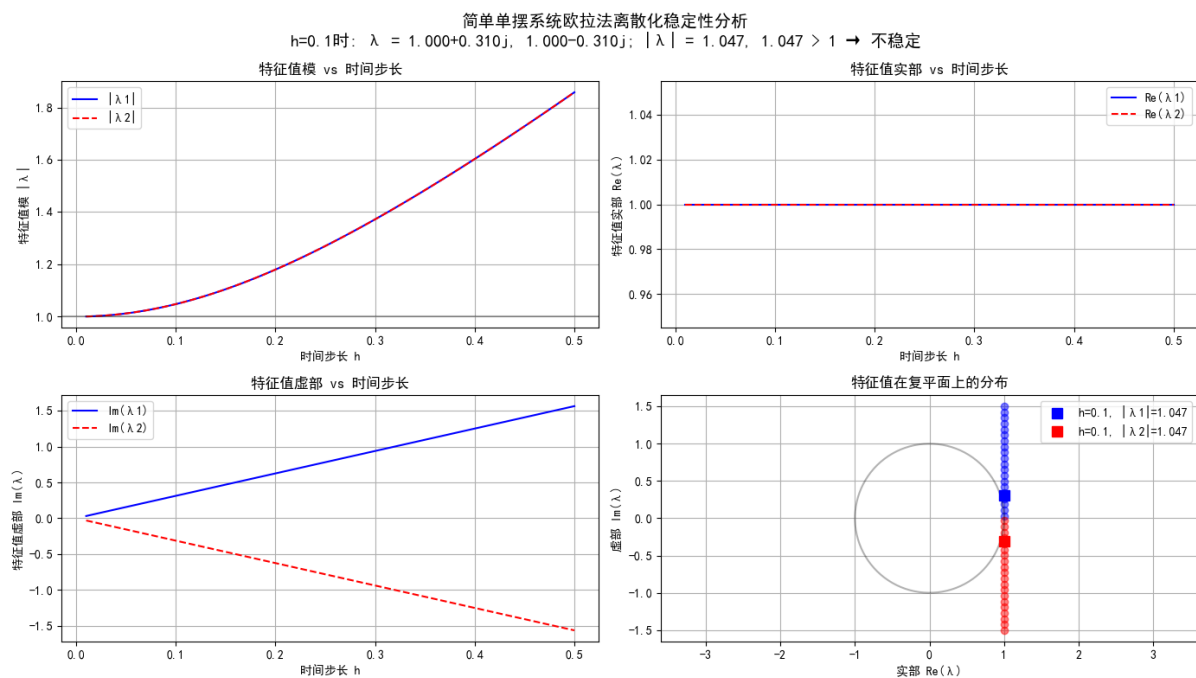


图 5: 图中展了简单单摆系统在欧拉法离散化后的稳定性分析, 通过观察不同时间步长下特征值的变化来判断系统的数值稳定性。

结论

- 谨慎使用正向欧拉法, 尤其是对于非线性系统。正向欧拉法在某些情况下可能会导致数值不稳定, 尤其是在系统具有强非线性或刚性时。对于这些情况, 可以考虑使用更高级的数值积分方法, 如RK4方法或隐式欧拉法等。
- 在稳定点时谨慎检测的能量行为, 是否是保守系统, 是否有耗散。
- 因为欧拉方法说到底是在 t 时刻的导数 $f(x_k, u_k)$, 相对于 x 总会有一个overshoot导致系统越来越blow up。所以避免使用正向欧拉法来模拟非线性系统的动力学行为。

2.2.2 RK4方法

4th order Runge-Kutta method (RK4) 是一种常用的数值积分方法, 用于求解常微分方程 (ODEs)。它通过在每个时间步内计算多个斜率来提高精度。RK4方法的基本思想是通过四个斜率的加权平均来近似ODE的解。

RK4方法* 这里将介绍RK4方法, 并举一个例子。

RK4方法 (四阶龙格-库塔法) 是一种常用的数值积分方法, 用于求解常微分方程。它的基本思想是通过在每个时间步内计算四个斜率来估算系统状态的变化,

进而得到更精确的解。

RK4方法的基本步骤如下：

1. 计算四个斜率：

$$k_1 = hf(x_k, u_k) \quad (\text{evaluate at beginning})$$

$$k_2 = hf\left(x_k + \frac{1}{2}k_1, u_k\right) \quad (\text{evaluate at midpoint 1})$$

$$k_3 = hf\left(x_k + \frac{1}{2}k_2, u_k\right) \quad (\text{evaluate at midpoint 2})$$

$$k_4 = hf(x_k + k_3, u_k) \quad (\text{evaluate at end})$$

其中， h 是时间步长， $f(x_k, u_k)$ 是系统的动力学方程， x_k 和 u_k 分别表示在时刻 k 时的状态和控制变量。

2. 更新状态：

$$x_{k+1} = x_k + \frac{1}{6}(k_1 + 2k_2 + 2k_3 + k_4)$$

通过这一步，RK4方法将四个斜率加权平均，得到更精确的状态更新公式。

RK4方法通过计算四个斜率来提高精度，相比于正向欧拉法，RK4方法在每个时间步内计算了更多的信息，从而能够更准确地逼近常微分方程的解。RK4方法的误差是 $O(h^4)$ ，比欧拉法的 $O(h^2)$ 更高效，适合用于高精度要求的数值计算。

手算例子： 现在，我们通过一个简单的例子来手动计算RK4方法的应用。假设我们要求解以下简单的常微分方程：

$$\frac{dx}{dt} = -2x$$

初始条件为 $x(0) = 1$ ，我们选择时间步长 $h = 0.1$ 。我们将使用RK4方法计算 x 在 $t = 0.1$ 时刻的值。

1. 计算四个斜率：

$$k_1 = 0.1 \cdot (-2 \cdot 1) = -0.2$$

$$k_2 = 0.1 \cdot (-2 \cdot (1 + \frac{1}{2} \cdot (-0.2))) = -0.19$$

$$k_3 = 0.1 \cdot (-2 \cdot (1 + \frac{1}{2} \cdot (-0.19))) = -0.19$$

$$k_4 = 0.1 \cdot (-2 \cdot (1 + (-0.19))) = -0.162$$

2. 更新状态:

$$x_1 = 1 + \frac{1}{6}(-0.2 + 2(-0.19) + 2(-0.19) + (-0.162)) = 0.818$$

因此, 使用RK4方法得到的 $x(0.1)$ 的近似值为0.818。

通过这个手算例子, 我们可以看到RK4方法比简单的欧拉法能够提供更精确的结果, 尤其是在步长较大的情况下。

```
1 # [discrete_time_dynamic_RK4.py]
2 import numpy as np
3 from src.PendulumVisualizer import *
4 def pendulum_dynamics(x):
5     l = 1.0
6     g = 9.81
7
8     theta = x[0]
9     theta_dot = x[1]
10    theta_ddot = -(g/l) * np.sin(theta)
11
12    return np.array([theta_dot, theta_ddot])
13
14 def pendulum_rk4(fun, x0, Tf, h):
15     t = np.arange(0, Tf + h, h)
16
17     x_hist = np.zeros((len(x0), len(t)))
18     x_hist[:, 0] = x0
19
20     for k in range(len(t) - 1):
21         x = x_hist[:, k]
22
23         k1 = fun(x)
24         k2 = fun(x + h/2 * k1)
25         k3 = fun(x + h/2 * k2)
26         k4 = fun(x + h * k3)
27
28         x_hist[:, k+1] = x + h/6 * (k1 + 2*k2 + 2*k3 + k4)
29
30     return x_hist, t
31
32 if __name__ == "__main__":
33     x0 = np.array([0.1, 0])
34     Tf = 10.0
35     h = 0.01
36
37     x_hist, t_hist = pendulum_rk4(pendulum_dynamics, x0, Tf, h)
38
39     visualizer = PendulumVisualizer(pendulum_length=1.0)
40     visualizer.visualize(x_hist, t_hist)
```

RK4稳定性分析 通过计算RK4的雅可比矩阵的特征值我们可以得到:

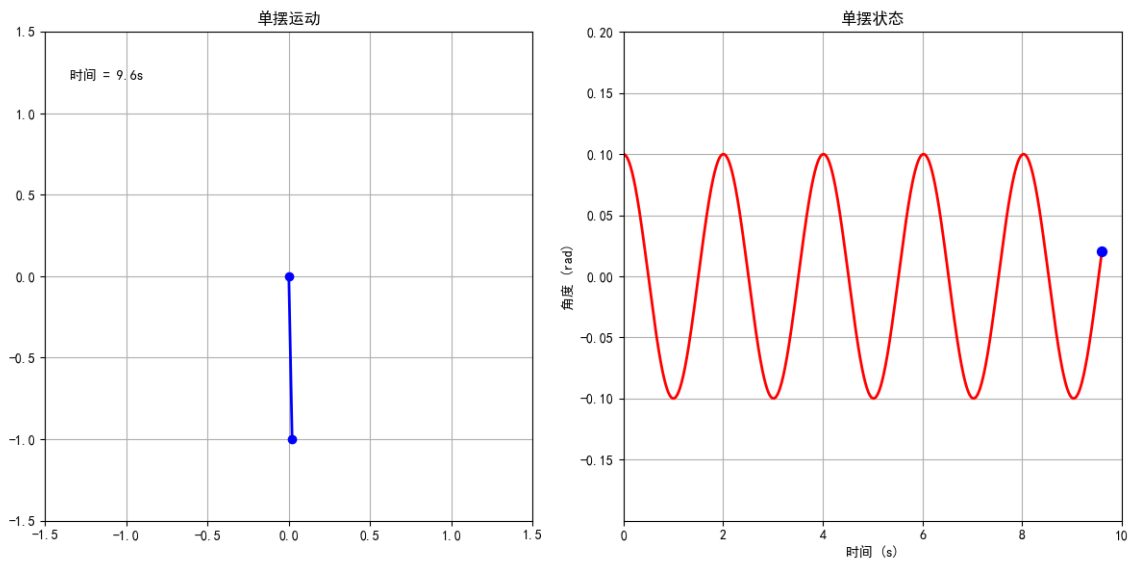


图 6: 图中展示了RK4模拟的单摆系统的运动轨迹。可以看到，随着时间的推移，单摆系统的运动轨迹更加平滑且稳定。

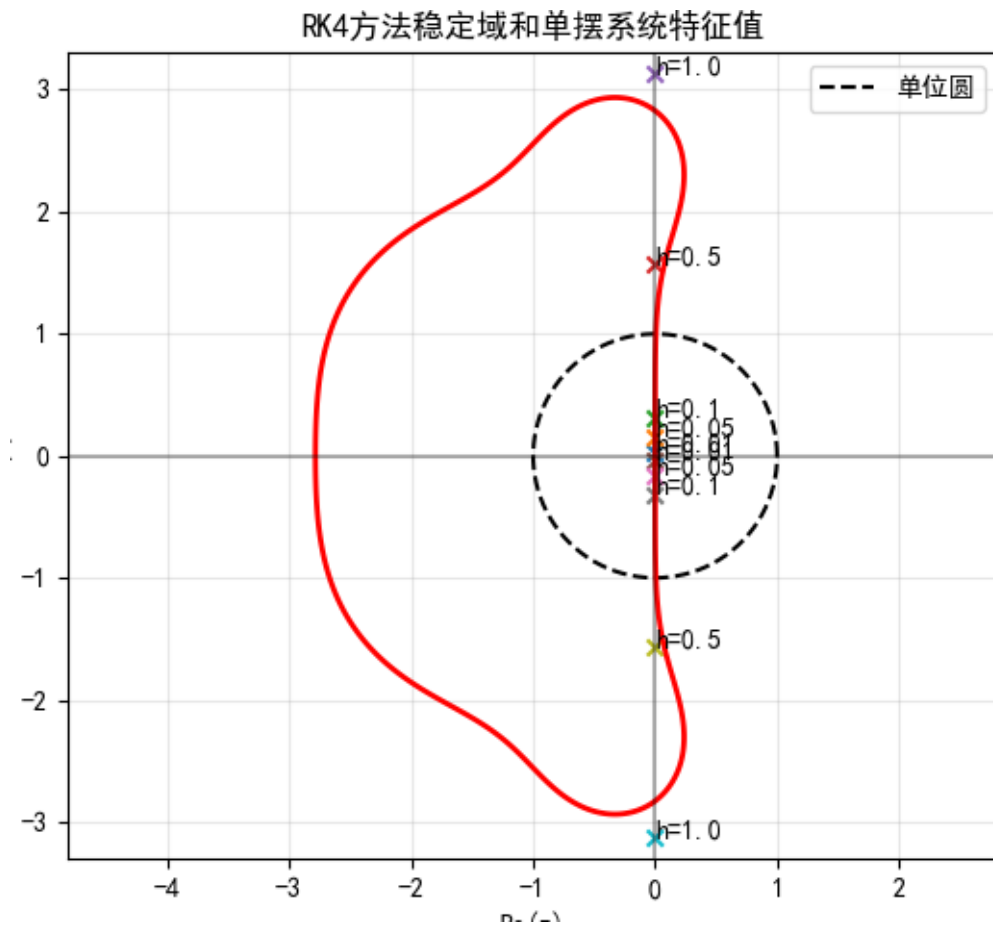


图 7: 图中展示了不同步长下，RK4方法的特征值，可以看到落在单位圆中的量。

RK4 方法稳定性分析

RK4 方法稳定性分析实验结果

- 单摆线性化系统的特征值:

$$\lambda = \begin{bmatrix} 0 + 3.1321j \\ 0 - 3.1321j \end{bmatrix}$$

- 步长 $h = 0.01$ 下, RK4 方法的雅可比矩阵特征值:

$$\lambda_J = \begin{bmatrix} 0.9995 + 0.0313j \\ 0.9995 - 0.0313j \end{bmatrix}$$

- 放大因子 (特征值绝对值):

$$|\lambda_J| = \begin{bmatrix} 0.999999999993 \\ 0.999999999993 \end{bmatrix}$$

不同步长下的稳定性分析

- 对于不同的步长 h , RK4 方法的稳定性如下所示:

$$h = 0.01$$

$$\lambda = \begin{bmatrix} 0 + 3.1321j \\ 0 - 3.1321j \end{bmatrix}, \quad \lambda_J = \begin{bmatrix} 0.9995 + 0.0313j \\ 0.9995 - 0.0313j \end{bmatrix}, \quad |\lambda_J| = \begin{bmatrix} 0.999999999993 \\ 0.999999999993 \end{bmatrix}$$

$$h = 0.05$$

$$\lambda = \begin{bmatrix} 0 + 3.1321j \\ 0 - 3.1321j \end{bmatrix}, \quad \lambda_J = \begin{bmatrix} 0.9878 + 0.1560j \\ 0.9878 - 0.1560j \end{bmatrix}, \quad |\lambda_J| = \begin{bmatrix} 0.999999897875 \\ 0.999999897875 \end{bmatrix}$$

$$h = 0.1$$

$$\lambda = \begin{bmatrix} 0 + 3.1321j \\ 0 - 3.1321j \end{bmatrix}, \quad \lambda_J = \begin{bmatrix} 0.9514 + 0.3081j \\ 0.9514 - 0.3081j \end{bmatrix}, \quad |\lambda_J| = \begin{bmatrix} 0.999993524289 \\ 0.999993524289 \end{bmatrix}$$

$$h = 0.5$$

$$\lambda = \begin{bmatrix} 0 + 3.1321j \\ 0 - 3.1321j \end{bmatrix}, \quad \lambda_J = \begin{bmatrix} 0.0244 + 0.9259j \\ 0.0244 - 0.9259j \end{bmatrix}, \quad |\lambda_J| = \begin{bmatrix} 0.9262 \\ 0.9262 \end{bmatrix}$$

$$h = 1.0$$

$$\lambda = \begin{bmatrix} 0 + 3.1321j \\ 0 - 3.1321j \end{bmatrix}, \quad \lambda_J = \begin{bmatrix} 0.1048 \pm 1.9889j \end{bmatrix}, \quad |\lambda_J| = \begin{bmatrix} 1.9916 \\ 1.9916 \end{bmatrix}$$

符号说明

- λ : 线性化系统的特征值，反映系统本身的动态特性；
- J : RK4 离散系统下的雅可比矩阵，用于描述线性近似；
- λ_J : RK4 方法下的雅可比矩阵特征值；
- $|\lambda_J|$: 特征值的模（放大因子），表示数值解是否随时间增长或衰减。

```
1  # [discrete_time_dynamic_RK4.py]
2  # 此处仅展示关键代码
3  def compute_rk4_jacobian_eigenvalues(h):
4      # 线性系统  $x' = Ax$  的情况下，A的特征值为 lambda
5      # 对于单摆在小角度近似下，线性化后的A矩阵为
6      l = 1.0
7      g = 9.81
8      A = np.array([[0, 1], [-g/l, 0]]) # 线性化的单摆动力学
9
10     # 计算A的特征值
11     eigen_A = np.linalg.eigvals(A)
12     print("单摆线性化系统的特征值:", eigen_A)
13
14     # 计算RK4方法的放大矩阵特征值
15     #  $R(z) = 1 + z + z^2/2 + z^3/6 + z^4/24$ 
16
17     # 计算不同特征值的放大因子
18     amplification_factors = []
19     eigenvalues = []
20
21     # 对每个特征值计算放大因子
22     for lam in eigen_A:
23         z = h * lam
24         R = 1 + z + z**2/2 + z**3/6 + z**4/24
25         amplification_factors.append(np.abs(R))
26         eigenvalues.append(R)
27
28     print(f'步长 h = {h} 下RK4方法的雅可比矩阵特征值:', eigenvalues)
29     print(f'放大因子(特征值绝对值):', amplification_factors)
30
31     return eigen_A, eigenvalues, amplification_factors
```

稳定性判据与特征值分析 设系统矩阵 A 的特征值为 λ_i ，步长为 h ，则有：

$$z_i = h \cdot \lambda_i$$

RK4 迭代中相应的放大因子为：

$$\rho_i = R(z_i)$$

其模长：

$$|\rho_i| = |R(z_i)|$$

判断稳定性标准：

- 若 $|\rho_i| < 1$ ，则该模态衰减（数值稳定）；
- 若 $|\rho_i| = 1$ ，则该模态保持（边界稳定）；
- 若 $|\rho_i| > 1$ ，则该模态发散（不稳定）。

这是因为在数值方法（如 Runge-Kutta 方法）中，我们并不是直接求解原始的微分方程 $\dot{x} = Ax$ ，而是通过一种离散时间的迭代形式来逼近它。

我们逐步解释为何会有：

$$z_i = h \cdot \lambda_i, \quad \rho_i = R(z_i)$$

1. 原始系统的连续解形式 设系统为：

$$\dot{x}(t) = Ax(t)$$

其解析解为：

$$x(t) = e^{At}x(0)$$

也就是说，系统演化由指数矩阵 e^{At} 控制。

2. 数值方法的离散迭代 RK4 等数值方法的本质是逼近这个指数矩阵的一种方式。对于一个离散步长 h ，我们用某个近似演化矩阵 $R(hA)$ 来代替 e^{hA} ，从而实现一步迭代：

$$x_{n+1} = R(hA)x_n$$

3. 为什么用特征值来简化？ 矩阵 A 的特征值 λ_i 描述了系统在各特征方向上的增长或衰荡速率。在分析数值稳定性时，我们关注每个模态（也就是特征值对应的方向）在一步后是否被放大或抑制。

由于特征值和矩阵函数的可交换性（在 A 可对角化时）：

$$\text{如果 } Av_i = \lambda_i v_i, \text{ 则 } R(hA)v_i = R(h\lambda_i)v_i$$

也就是说，数值迭代对每个特征模态的作用可以简化为一个复数放大因子：

$$\rho_i = R(h\lambda_i)$$

其中：

- $h\lambda_i = z_i$: 当前模态对应的无量纲时间步;
- $R(z_i)$: RK4 稳定性函数, 在 z_i 的值, 给出该模态的“放大比例”;
- $|R(z_i)|$: 判断该模态是否稳定 (小于1表示衰减, 等于1表示保持, 大于1表示爆炸)。

4. 直觉类比 你可以将 λ_i 理解为一个“频率”或“增长率”, 而 h 是我们离散模拟的“采样间隔”。那么:

$$z_i = h\lambda_i$$

用这个“频率率 (时间步长)”来捕捉这个模态, 数值方法能否稳定逼近它? 如果 z_i 落在 RK4 的稳定区域, 则该模态可以被良好模拟; 否则, 就会出现误差爆炸。

2.2.3 反向欧拉法

pybullet、mujoco使用的正是反向欧拉方法。对于隐式ODEs, 例如:

$$f_d(x_{k+1}, x_k, u_k) = 0 \quad (32)$$

这表示在每个时间步内, 系统的状态 x_{k+1} 和控制输入 u_k 之间的关系是隐式的。我们可以通过迭代的方法来求解这个方程。反向欧拉法 (Backward Euler Method) 是一种常用的数值积分方法, 用于求解隐式ODEs。它的基本思想是通过在每个时间步内计算一个隐式方程来估算系统状态的变化, 进而得到更精确的解。

$$x_{k+1} = x_k + hf(x_{k+1}) \quad (33)$$

相较于forward euler, 这里的backward euler计算 $f(x)$ 是用的未来的时间 x_{k+1} 而不是当前的时间 x_k 。于是我们可以通过反向欧拉法来模拟系统:

$$f_d(x_{k+1}, x_k, u_k) = x_{k+1} - x_k - hf(x_{k+1}, u_k) = 0 \quad (34)$$

可以将这个问题转换为一个 x_{k+1} 时刻的root-finding问题, 这个问题将在之后介绍。类似于机器人学中计算inverse dynamic的Newton-Euler方法? (TODO: 添加引用)

反向欧拉法直觉理解 backward euler 通过离散化向系统添加了damping, 即产生了undershoot (相对于euler方法的overshoot), 使得系统在每个时间步内都能保持稳定。通过在每个时间步内计算一个隐式方程, 反向欧拉法能够更好地捕捉系统的动态行为, 尤其是在处理刚性系统时。反向欧拉法的误差是 $O(h)$, 比欧拉法的 $O(h^2)$ 更高效, 适合用于高精度要求的数值计算。但是这导致了由此构建的模拟器与现实存在很大的差距! 值得注意的是:

- 显式方法通常比隐式方法更稳定。
- 对于forward simulation，解决隐式ODEs代价很高。
- 对于很多直接的轨迹优化方法，不再如此高代价。

2.2.4 离散控制输入

对于系统，我们前文论述了状态通过 $x(t)$ 来描述，并考虑了连续与离散时间的状态。现在我们来考虑系统的控制输入 $u(t)$ 的离散化表示。

$$u(t) = u_k, t_k \leq t < t_{k+1} \quad (35)$$

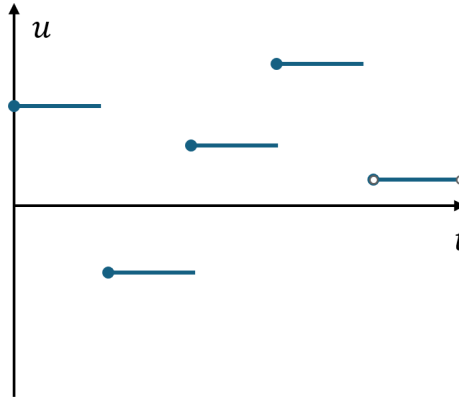


图 8: 这是一个zero-order hold的示意图。它表示在每个时间步内，控制输入 $u(t)$ 保持不变，直到下一个时间步到来。这个方法在数字控制系统中非常常见，因为它可以简化控制器的设计和实现。

零阶保持 $u(t) = u_n, t_n \leq t \leq T_{n+1}$

而一阶保持 $u(t) = u_n + (\frac{u_{k+1} - u_k}{h}(t - t_n))$ ，其中 h 是时间步长。一阶保持是一个更好的代替方法。它表示在每个时间步内，控制输入 $u(t)$ 是一个线性函数，直到下一个时间步到来。这个方法在数字控制系统中也很常见，因为它可以更好地捕捉系统的动态行为。

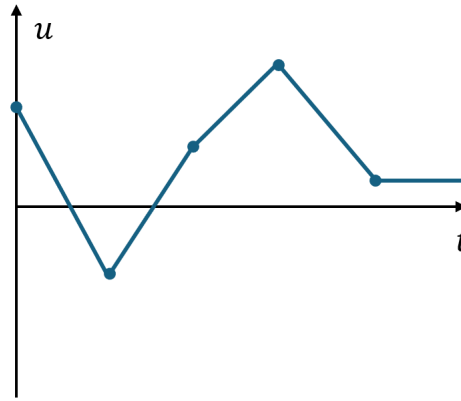


图 9: 这是一个first-order hold的示意图。它表示在每个时间步内，控制输入 $u(t)$ 是一个线性函数，直到下一个时间步到来。这个方法在数字控制系统中也很常见，因为它可以更好地捕捉系统的动态行为。通过使用一阶保持方法，我们可以更好地捕捉系统的动态行为，尤其是在处理快速变化的系统时。与零阶保持相比，一阶保持方法能够更准确地模拟系统的响应，并减少控制输入的突变，从而提高系统的稳定性和性能。

高阶保持，使用高阶的多项式来近似控制输入。然而在Bang-Bang控制中（比如说击球的这一瞬间，球换向快速弹回），两个极端值快速变换，其最优控制率不是光滑的，此时高阶保持甚至不如零阶保持

3 数值优化基础

这一节将介绍常用的数值优化方法。从一个 root-finding 问题开始，然后介绍

- 不动点迭代方法
- 牛顿方法
 - 最小化
 - 充分条件
 - 正则化
 - overshooting 问题

3.1 符号声明

这一节介绍多个符号规定，记住就好。

一元函数导数维度 考虑 $f(x): \mathbb{R}^n \rightarrow \mathbb{R}$ ，则

$$\frac{\partial f}{\partial x} \in \mathbb{R}^{1 \times n} \Rightarrow \text{梯度是行向量。}$$

因为 $\frac{\partial f}{\partial x}$ 是将 Δx 映射到 Δf 的线性算子。并且通过泰勒展开我们可以得到：

$$f(x + \Delta x) \approx f(x) + \frac{\partial f}{\partial x} \Delta x,$$

因为 Δx 为列向量，所以如果要想 $f(x + \Delta x)$ 合法，则 $\frac{\partial f}{\partial x}$ 必然是行向量，

多元函数偏导（雅可比）维度 类似地，设 $g(y): \mathbb{R}^m \rightarrow \mathbb{R}^n$ ，其雅可比为

$$J := \frac{\partial g(y)}{\partial y} \in \mathbb{R}^{n \times m}, \quad J = \begin{pmatrix} \frac{\partial g_1}{\partial y_1} & \cdots \\ \vdots & \ddots \end{pmatrix}$$

因为：

$$g(y + \Delta y) \approx g(y) + \frac{\partial g}{\partial y} \Delta y.$$

中 Δy 是 $m \times 1$ ，所以雅可比必然是 $n \times m$ 。

链式法则（Chain Rule）：下面公式演示了复合函数的一阶泰勒展开

$$f(g(y + \Delta y)) \approx f(g(y)) + \frac{\partial f}{\partial x} \cdot \frac{\partial g}{\partial y} \cdot \Delta y.$$

nabla 记号 ∇

$$\nabla f(x) = \left(\frac{\partial f}{\partial x} \right)^\top \in \mathbb{R}^{n \times 1} \quad \text{列向量,}$$

$$\nabla^2 f(x) = \frac{\partial}{\partial x}(\nabla f(x)) = \frac{\partial^2 f}{\partial x^2} \in \mathbb{R}^{n \times n}.$$

3.2 Root Finding问题

考虑根问题 (Root Finding): 给定函数 $f(x)$, 希望找到 x^* 使得

$$f(x^*) = 0 \quad \Rightarrow \quad x^* \text{ 是根.}$$

例如: 连续时间动力系统的平衡点

$$\dot{x} = f(x), \quad \text{令 } f(x^*) = 0 \Rightarrow x^* \text{ 是平衡.}$$

这个问题与不动点问题密切相关, 即求解满足

$$f(x^*) = x^*$$

的 x^* , 例如离散时间动力系统的平衡点。

不动点迭代 (Fixed-Point Iteration) 最朴素的做法是不断迭代:

$$x_{k+1} = f(x_k),$$

若动态系统是稳定的, 则不断前向模拟该系统, 将会收敛到某个平衡点。

不动点收敛性分析 考虑映射 $x_{k+1} = f(x_k)$, 若存在不动点 x^* 满足 $f(x^*) = x^*$, 我们希望该迭代能收敛到 x^* 。

定理 (Banach 不动点定理): 设 $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ 是压缩映射, 即存在 $0 < L < 1$, 使得对任意 x, y , 有

$$\|f(x) - f(y)\| \leq L\|x - y\|,$$

则 f 有唯一不动点 x^* , 且对任意初始点 x_0 , 迭代 $x_{k+1} = f(x_k)$ 收敛于 x^* 。

收敛性证明: 记 x^* 为不动点, 考虑误差 $e_k := x_k - x^*$, 则

$$e_{k+1} = x_{k+1} - x^* = f(x_k) - f(x^*) \Rightarrow \|e_{k+1}\| \leq L\|e_k\| \leq L^2\|e_{k-1}\| \leq \cdots \leq L^k\|e_0\|.$$

因此有:

$$\|x_k - x^*\| \leq L^k\|x_0 - x^*\| \rightarrow 0 \quad \text{当 } k \rightarrow \infty.$$

这说明只要映射满足 Lipschitz 常数 $L < 1$, 就能保证收敛性, 且为线性收敛。

梯度意义下的不动点 在梯度下降的意义下，若将梯度流离散化：

1. 前向欧拉：

$$x_{k+1} = x_k - h \nabla f(x_k),$$

这就是梯度下降（Gradient Descent）的离散形式。

注：可添加动画或直观演示。

2. 后向欧拉：类似于近似梯度法，也可看作是一种近似的「梯度流」方法（continuous limit）。

为什么梯度下降在大规模问题中更受欢迎？ 在大规模机器学习中，梯度下降（Gradient Descent）相比牛顿方法（Newton's Method）更常用，原因包括：

- **计算更便宜：** 牛顿法需要计算 Hessian 并求逆，代价为 $\mathcal{O}(n^3)$ ，而梯度下降只需计算一阶导数。（将在后面的公式中看到这个复杂度）
- **高维不可接受：** 在高维情况下，牛顿方法由于计算复杂度过高，不适合使用。
- 此外，对于非凸函数（non-convex function）：
 - **多个平衡点（equilibria）：** 若希望找到全局最小值，需要在多个平衡点之间选择，等价于 NP 难问题。
 - **几乎需要穷举搜索（exhaustive search）：** 没有启发式或凸性结构支持时，只有遍历所有不动点才能确定最优解。
 - **仅适用于稳定的不动点：** 牛顿迭代法仅在初始点足够靠近稳定的不动点时才有效；否则会发散。
 - **收敛速度慢：** 特别在函数曲率变化大或存在鞍点时，牛顿法的收敛性无法保证。

3.2.1 Newton 方法

Newton 方法的基本思想是对 $f(x)$ 作线性近似，并设为零从而求解 Δx ，用于迭代更新。

由线性近似，在当前点 x 附近，采用一阶泰勒展开：

$$f(x + \Delta x) \approx f(x) + \frac{\partial f}{\partial x} \Delta x.$$

然后立即可以得到牛顿方法的两个步骤

- 令近似为零，转为 root-finding 问题，得到第一步，计算

$$f(x) + \frac{\partial f}{\partial x} \Delta x = 0 \quad \Rightarrow \quad \Delta x = - \left(\frac{\partial f}{\partial x} \right)^{-1} f(x)$$

注：这里求雅可比逆是 $O(n^3)$

- 第二步，更新变量：

$$x \leftarrow x + \Delta x,$$

- 反复迭代直到收敛。

注：若导数矩阵不可逆或病态，可采用正则化策略：

如添加零平方项以求伪逆 (pseudo-inverse)：即 Tikhonov regularization。

Backward Euler 步的固定点迭代示例 考虑使用 backward Euler 步来构造一个固定点迭代过程。目标是求解如下形式的不动点：

$$x_n = x_0 + hf(x_n),$$

其中 h 是步长参数。

- 初始化： $x_n = x_0$
- 误差定义：

$$\text{err} = \|x_0 + hf(x_n) - x_n\|_2$$

- 迭代至收敛：

$$\begin{aligned} \text{while err} > \epsilon : \quad & x_n \leftarrow x_0 + hf(x_n) \\ & \text{err} = \|x_0 + hf(x_n) - x_n\|_2^2 \end{aligned}$$

该迭代形式对应于隐式欧拉法或非线性方程的固定点迭代形式，可看作是牛顿方法的近似替代方案。

Backward Euler + Newton 方法联合迭代 (backward_euler_step_newton_method)

目标：求解隐式方程

$$x_n = x_0 + hf(x_n)$$

这是 Backward Euler 步的形式，我们使用 Newton 方法求解该不动点方程。

- 定义残差函数：

$$r(x_n) := x_0 + hf(x_n) - x_n$$

- 迭代步骤：

$$x_n \leftarrow x_0$$

while $\|r(x_n)\| > \epsilon$ do:

 计算残差 $r = x_0 + hf(x_n) - x_n$

 计算雅可比 $\partial r = \frac{\partial r}{\partial x_n} = h \frac{\partial f}{\partial x}(x_n) - I$

 牛顿更新: $x_n \leftarrow x_n - (\partial r)^{-1} r$

这个过程将 Backward Euler 构造的非线性方程，通过 Newton 方法迭代求解，适用于刚性系统或更稳定收敛的时间积分方法。

Backward Euler 与 Newton 方法的比较

- Backward Euler 方法虽然具有 **人工阻尼 (artificial damping)** 效果，但是收敛很慢（ $-\text{Re}(\text{eig}(R))$ —可以看出来），较稳定。
- 固定点迭代方法（如 $x_{k+1} = f(x_k)$ ）通常 **不能直接找到极小值**，其收敛速度为 **线性收敛 (linear convergence)**，如图所示呈指数衰减但较慢。
- 相比之下，Newton 方法可在 **约 3 步迭代内** 快速接近最优解，具有 **二次收敛率 (quadratic convergence)**，即误差下降速度为：

$$\|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2.$$

如图所示，误差曲线呈指数平方下降趋势。（TODO:添加图，添加代码）

Newton 方法的总结与优势

- Newton 方法配合 backward Euler 等技术，**收敛非常快**。
- 拥有 **二次收敛 (Quadratic Convergence)**，即误差下降速度极快。
- 可以达到 **机器精度 (machine precision)**，误差极小。
- **主要代价**在于对 Jacobian 的因式分解 / 求逆操作，复杂度约为：

$$\mathcal{O}(n^3)$$

- 若能 **利用问题结构 (problem structure)**，可进一步降低复杂度，后续会展开说明。

3.2.2 牛顿方法在最小化问题中的应用

考虑优化问题

$$\min_x f(x), \quad f(x): \mathbb{R}^n \rightarrow \mathbb{R}$$

- 若 f 是光滑的，则要满足 $\frac{\partial f}{\partial x}|_{x^*} = 0$ 在局部最小值成立 (KKT)。
- 可将该问题转化为 root-finding 问题: $\nabla f(x) = 0$

应用 Newton 方法:

1. 计算更新量

$$\begin{aligned}\nabla f(x + \Delta x) &\approx \nabla f(x) + \nabla^2 f(x) \Delta x = 0 \\ \Rightarrow \Delta x &= -[\nabla^2 f(x)]^{-1} \nabla f(x)\end{aligned}$$

2. 更新

$$x \leftarrow x + \Delta x$$

解释: 牛顿方法 = 二次近似最小化 牛顿方法可以理解为在当前点 x 对 $f(x)$ 作局部二次近似:

$$f(x + \Delta x) \approx f(x) + \nabla f(x)^\top \Delta x + \frac{1}{2} \Delta x^\top \nabla^2 f(x) \Delta x$$

然后最小化该二次函数，得到的最优 Δx 就是 Newton 步。

例子:

$$f(x) = x^4 + x^3 - x^2 - x, \quad \nabla f(x) = 4x^3 + 3x^2 - 2x - 1, \quad \nabla^2 f(x) = 12x^2 + 6x - 2$$

图像显示 $f(x)$ 有多个局部极小值，例如 $x \approx 1.0, -1.5, 0.0$ (TODO:添加图片)

结论: Newton 方法是 **局部 root-finding** 方法，对 $\nabla f(x) = 0$ 求解，其收敛结果取决于初始点:

- 可能收敛到极小值、极大值或鞍点。
- 要求初始点足够接近目标点。

3.3 最小值存在的充分条件

- 一阶必要条件: $\nabla f(x^*) = 0$ 是最小值的必要条件，但不是充分条件!
- 看一个标量情况:

$$\Delta x = -(\nabla^2 f)^{-1} \cdot \nabla f$$

梯度方向决定前进方向，Hessian 决定“步长”或“学习率”。

- 若 Hessian 为正:

$$\nabla^2 f > 0 \Rightarrow \text{descent (最小化)}$$

若 Hessian 为负:

$$\nabla^2 f < 0 \Rightarrow \text{ascent (最大化)}$$

- 更进一步, 在 \mathbb{R}^n 中, 若

$$\nabla^2 f \succ 0 \quad (\text{严格正定}) \Rightarrow \nabla^2 f \in \mathcal{S}_{++}$$

即所有特征值均为正, 是局部极小的充分条件。

因此, 局部极小的二阶充分条件是:

在极小值点, Hessian 的所有特征值为正, 即

\Rightarrow 方向为下降方向 (descent)

3.4 正则化 (Regularization) 与强凸性

- 若 $\nabla^2 f > 0$ 在整体上成立 $\Leftrightarrow f(x)$ 是强凸函数 (strongly convex), 则可用牛顿法全局求解。
- 但这对非凸或非线性问题通常不成立, 因此不能总是依赖牛顿法。
- 正则化的动机: 为了让 Hessian 始终正定, 从而确保下降方向, 可以进行如下处理:

– 设 $H = \nabla^2 f$, 若 $H \not\succ 0$, 则执行:

$$\text{while } H \not\succ 0: \quad H \leftarrow H + \beta I$$

其中 $\beta > 0$, I 为单位矩阵, β 可逐步增大。

– 更新步长:

$$\Delta x = -H^{-1} \nabla f, \quad x \leftarrow x + \Delta x$$

– 这种方法称为 **Tikhonov 正则化 (Tikhonov regularization)**, 在最小二乘 (least squares) 等问题中也常见, 体现为给权重加上二次惩罚项。

- **damped Newton** 方法, 通过调节步长或对 Hessian 添加扰动来控制更新方向, 使其总是一个下降步 (descent step):

– 它保证是下降方向, 同时还能让步长变得更小 (make step size smaller)。

– 给 Hessian 加上单位阵 (positive identity) 本质上等价于在优化目标中添加一个二次惩罚项 (quadratic penalty):

例如 $\beta \Delta x^\top \Delta x$ 被加到目标函数上 \Rightarrow 惩罚步长过大, 避免 overshoot

为什么给 Hessian 加上 βI 等价于加惩罚项？ 在牛顿方法中，我们在 x 附近用二阶泰勒展开近似目标函数：

$$f(x + \Delta x) \approx f(x) + \nabla f(x)^\top \Delta x + \frac{1}{2} \Delta x^\top \nabla^2 f(x) \Delta x$$

此时求解的是使该近似最小的 Δx 。

为了增强数值稳定性，或者避免 Hessian 非正定带来的问题，我们会对 Hessian 做如下正则化：

$$\tilde{H} = \nabla^2 f(x) + \beta I$$

此时对应的优化问题变为最小化：

$$\begin{aligned} & f(x) + \nabla f(x)^\top \Delta x + \frac{1}{2} \Delta x^\top (\nabla^2 f(x) + \beta I) \Delta x \\ &= f(x) + \nabla f(x)^\top \Delta x + \frac{1}{2} \Delta x^\top \nabla^2 f(x) \Delta x + \frac{1}{2} \beta \|\Delta x\|^2 \end{aligned}$$

因此，给 Hessian 加 βI 等价于在原函数上添加一项二次惩罚项：

$$\frac{1}{2} \beta \|\Delta x\|^2$$

这样做的意义是：

- * 控制 Δx 的幅度，防止步长过大（overshooting）
- * 保证 Hessian 正定，从而确保下降方向
- * 提高病态 Hessian 的数值稳定性

– 总结：这种方法的效果是 **抑制步长过大**，尤其适合 Hessian 非正定或发散风险高的情形。

- *TODO*: 加入工程上的实际例子来说明此策略效果。

Appendix