

ODE-Constrained Optimization with Automatic Differentiation

Tianju Xue

December 11, 2022

1 Background

This tutorial discusses the formulation of solving ODE-constrained optimization (ODE-CO) problems with the adjoint method. The tutorial comes with an example about the famous [Lorenz system](#), implemented in this [Jupyter Notebook](#). The code is based on [JAX](#) with its handy Automatic Differentiation feature that makes life much easier.

2 ODE-constrained optimization

Assume $\mathbf{u}(t, \mathbf{a}, \mathbf{b}) \in \mathbb{R}^k$ is the variable we want to solve with $\mathbf{a} \in \mathbb{R}^m$ being the ODE parameters and $\mathbf{b} \in \mathbb{R}^k$ being the initial condition. Note that t , \mathbf{a} , and \mathbf{b} are independent variables, and determine the value of \mathbf{u} . The forward problem is defined as

$$\frac{d}{dt}\mathbf{u}(t, \mathbf{a}, \mathbf{b}) = \mathbf{r}(\mathbf{u}, t, \mathbf{a}), \quad (1a)$$

$$\mathbf{u}(t_0, \mathbf{a}, \mathbf{b}) = \mathbf{b}, \quad (1b)$$

which is an ODE system. In some cases, this ODE system is just the semi-discretized system of the Finite Element/Finite difference method with spatial discretization performed but time discretization not yet performed (e.g., [method of lines](#)).

We are interested in solving the inverse problem or the design problem. In plain explanation, an inverse problem is when you have certain desired requirement for the solution \mathbf{u} and you want to figure out what values of \mathbf{a} and \mathbf{b} fulfill the requirement. Mathematically, we have an optimization problem of the following form:

$$\min_{\mathbf{a} \in \mathbb{R}^m, \mathbf{b} \in \mathbb{R}^k} \mathcal{J}(\mathbf{u}) = h(\mathbf{u}(t_f, \mathbf{a}, \mathbf{b})) + \int_{t_0}^{t_f} g(\mathbf{u}(t, \mathbf{a}, \mathbf{b}))dt, \quad (2a)$$

$$\text{s.t. } \frac{d}{dt}\mathbf{u}(t, \mathbf{a}, \mathbf{b}) = \mathbf{r}(\mathbf{u}, t, \mathbf{a}), \quad (2b)$$

$$\mathbf{u}(t_0, \mathbf{a}, \mathbf{b}) = \mathbf{b}, \quad (2c)$$

where $\mathcal{J} : \mathcal{U} \rightarrow \mathbb{R}$ is the objective functional, $h : \mathbb{R}^k \rightarrow \mathbb{R}$, and $g : \mathbb{R}^k \rightarrow \mathbb{R}$.

Since \mathbf{u} is implicitly determined by \mathbf{a} and \mathbf{b} via the ODE system, for a fixed set of \mathbf{a} and \mathbf{b} we have a deterministic \mathcal{J} . Therefore, it makes sense to consider the total derivative of \mathcal{J} with respect to \mathbf{a} and \mathbf{b} . In fact, the central goal of this tutorial is to show how to compute

$$\frac{d\mathcal{J}}{d\mathbf{a}} \quad \text{and} \quad \frac{d\mathcal{J}}{d\mathbf{b}}, \quad (3)$$

so that gradient-based optimization algorithms can be employed to solve the ODE-CO problem. We will introduce two approaches to compute $\frac{d\mathcal{J}}{d\mathbf{a}}$ and $\frac{d\mathcal{J}}{d\mathbf{b}}$. Depending on when the time discretization happens, we refer to these two approaches as “early discretization” and “late discretization”, respectively.

3 Early discretization (discrete adjoint method)

This part is based on [1]. Let us immediately discretize the ODE system (1) in time so that

$$\mathbf{u}^n = \mathbf{f}(\mathbf{u}^{n-1}, n, \mathbf{a}) := \mathbf{f}^n, \quad \mathbf{u}^0 = \mathbf{b}, \quad (4)$$

where n is the time step number, and \mathbf{f} is any explicit ODE integrator, e.g., explicit Euler method or Runge-Kutta method.

The objective function is

$$J = h(\mathbf{u}^N), \quad (5)$$

where we have omitted the g term for simplicity.

With chain rule, we have

$$\frac{dJ}{d\mathbf{a}} = \frac{dh}{d\mathbf{u}^N} \left\{ \frac{\partial \mathbf{f}^N}{\partial \mathbf{a}} + \frac{\partial \mathbf{f}^N}{\partial \mathbf{u}^{N-1}} \left[\frac{\partial \mathbf{f}^{N-1}}{\partial \mathbf{a}} + \frac{\partial \mathbf{f}^{N-1}}{\partial \mathbf{u}^{N-2}} \left(\frac{\partial \mathbf{f}^{N-2}}{\partial \mathbf{a}} + \dots \right) \right] \right\}, \quad (6)$$

and

$$\frac{dJ}{d\mathbf{b}} = \frac{dh}{d\mathbf{u}^N} \left\{ \frac{\partial \mathbf{f}^N}{\partial \mathbf{b}} + \frac{\partial \mathbf{f}^N}{\partial \mathbf{u}^{N-1}} \left[\frac{\partial \mathbf{f}^{N-1}}{\partial \mathbf{b}} + \frac{\partial \mathbf{f}^{N-1}}{\partial \mathbf{u}^{N-2}} \left(\frac{\partial \mathbf{f}^{N-2}}{\partial \mathbf{b}} + \dots \right) \right] \right\} \quad (7)$$

$$= \frac{dh}{d\mathbf{u}^N} \frac{\partial \mathbf{f}^N}{\partial \mathbf{u}^{N-1}} \frac{\partial \mathbf{f}^{N-1}}{\partial \mathbf{u}^{N-2}} \dots \frac{\partial \mathbf{f}^1}{\partial \mathbf{u}^0} \frac{d\mathbf{u}^0}{d\mathbf{b}}. \quad (8)$$

Let us define adjoint variable $\boldsymbol{\lambda}_{\mathbf{b}} \in \mathbb{R}^k$ associated with initial condition \mathbf{b} so that

$$\boldsymbol{\lambda}_{\mathbf{b}}^{n-1} = \boldsymbol{\lambda}_{\mathbf{b}}^n \frac{\partial \mathbf{f}^n}{\partial \mathbf{u}^{n-1}}, \quad \boldsymbol{\lambda}_{\mathbf{b}}^N = \frac{dh}{d\mathbf{u}^N}, \quad (9)$$

and similarly the adjoint variable $\boldsymbol{\lambda}_{\mathbf{a}} \in \mathbb{R}^m$ associated with ODE parameters \mathbf{a} so that

$$\boldsymbol{\lambda}_{\mathbf{a}}^{n-1} = \boldsymbol{\lambda}_{\mathbf{b}}^n \frac{\partial \mathbf{f}^n}{\partial \mathbf{a}} + \boldsymbol{\lambda}_{\mathbf{a}}^n, \quad \boldsymbol{\lambda}_{\mathbf{a}}^N = \mathbf{0}. \quad (10)$$

With some tedious algebraic operations, we have

$$\frac{dJ}{d\mathbf{a}} = \boldsymbol{\lambda}_{\mathbf{a}}^0 \quad \text{and} \quad \frac{dJ}{d\mathbf{b}} = \boldsymbol{\lambda}_{\mathbf{b}}^0. \quad (11)$$

It is important to note that since the forward problem, i.e., Eq. (4) is not reversible (at least not so explicit to compute \mathbf{u}^{n-1} from \mathbf{u}^n), we have to store the entire trajectories of $\mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}^N$ in memory because they are required in the adjoint updates of Eq. (9) and Eq. (10). Therefore this is a linear memory approach, and does not scale well if the ODE system is large.

4 Late discretization (continuous adjoint method)

This part is based on the Neural-ODE paper [2], the associated [tutorial](#), and the nice [Youtube video](#) by Chris H. Rycroft at Harvard.

We will delay the discretization of the ODE system. In stead, first notice that by Eq. 2b we have

$$\frac{d}{dt} \left(\frac{d\mathbf{u}}{d\mathbf{a}} \right) = \frac{d\mathbf{r}}{d\mathbf{a}}, \quad \frac{d}{dt} \left(\frac{d\mathbf{u}}{d\mathbf{b}} \right) = \frac{d\mathbf{r}}{d\mathbf{b}} = \mathbf{0}. \quad (12)$$

Then consider the adjoint variable $\lambda_{\mathbf{a}}(t) \in \mathbb{R}^m$ and $\lambda_{\mathbf{b}}(t) \in \mathbb{R}^k$ that satisfy the adjoint ODE system:

$$\frac{d}{dt} \begin{bmatrix} \lambda_{\mathbf{a}} \\ \lambda_{\mathbf{b}} \end{bmatrix} = \begin{bmatrix} -\lambda_{\mathbf{b}} \frac{\partial r}{\partial \mathbf{a}} \\ -\lambda_{\mathbf{b}} \frac{\partial r}{\partial \mathbf{u}} - \frac{\partial g}{\partial \mathbf{u}} \end{bmatrix}, \quad (13)$$

$$\begin{bmatrix} \lambda_{\mathbf{a}}(t_f) \\ \lambda_{\mathbf{b}}(t_f) \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \left. \frac{\partial h}{\partial \mathbf{u}} \right|_{t=t_f} \end{bmatrix}, \quad (14)$$

where $\lambda_{\mathbf{a}}$ and $\lambda_{\mathbf{b}}$ are solved backwards in time.

We have

$$\frac{d\mathcal{J}}{d\mathbf{a}} = \left(\frac{\partial h}{\partial \mathbf{u}} \frac{d\mathbf{u}}{d\mathbf{a}} \right) \Big|_{t=t_f} + \int_{t_0}^{t_f} \frac{\partial g}{\partial \mathbf{u}} \frac{d\mathbf{u}}{d\mathbf{a}} dt \quad (15)$$

$$= \lambda_{\mathbf{b}}(t_f) \frac{d\mathbf{u}}{d\mathbf{a}} \Big|_{t=t_f} + \int_{t_0}^{t_f} \frac{\partial g}{\partial \mathbf{u}} \frac{d\mathbf{u}}{d\mathbf{a}} dt$$

$$= \lambda_{\mathbf{b}}(t_0) \frac{d\mathbf{u}}{d\mathbf{a}} \Big|_{t=t_0} + \int_{t_0}^{t_f} \left(\frac{d\lambda_{\mathbf{b}}}{dt} \frac{d\mathbf{u}}{d\mathbf{a}} + \lambda_{\mathbf{b}} \frac{d}{dt} \left(\frac{d\mathbf{u}}{d\mathbf{a}} \right) + \frac{\partial g}{\partial \mathbf{u}} \frac{d\mathbf{u}}{d\mathbf{a}} \right) dt$$

$$= \lambda_{\mathbf{b}}(t_0) \frac{d\mathbf{b}}{d\mathbf{a}} + \int_{t_0}^{t_f} \lambda_{\mathbf{b}} \frac{\partial r}{\partial \mathbf{a}} dt \quad (16)$$

$$= \lambda_{\mathbf{a}}(t_0), \quad (17)$$

and similarly,

$$\frac{d\mathcal{J}}{d\mathbf{b}} = \left(\frac{\partial h}{\partial \mathbf{u}} \frac{d\mathbf{u}}{d\mathbf{b}} \right) \Big|_{t=t_f} + \int_{t_0}^{t_f} \frac{\partial g}{\partial \mathbf{u}} \frac{d\mathbf{u}}{d\mathbf{b}} dt \quad (18)$$

$$= \lambda_{\mathbf{b}}(t_f) \frac{d\mathbf{u}}{d\mathbf{b}} \Big|_{t=t_f} + \int_{t_0}^{t_f} \frac{\partial g}{\partial \mathbf{u}} \frac{d\mathbf{u}}{d\mathbf{b}} dt$$

$$= \lambda_{\mathbf{b}}(t_0) \frac{d\mathbf{u}}{d\mathbf{b}} \Big|_{t=t_0} + \int_{t_0}^{t_f} \left(\frac{d\lambda_{\mathbf{b}}}{dt} \frac{d\mathbf{u}}{d\mathbf{b}} + \lambda_{\mathbf{b}} \frac{d}{dt} \left(\frac{d\mathbf{u}}{d\mathbf{b}} \right) + \frac{\partial g}{\partial \mathbf{u}} \frac{d\mathbf{u}}{d\mathbf{b}} \right) dt$$

$$= \lambda_{\mathbf{b}}(t_0) \frac{d\mathbf{b}}{d\mathbf{b}} + \int_{t_0}^{t_f} \lambda_{\mathbf{b}} \frac{\partial r}{\partial \mathbf{b}} dt \quad (19)$$

$$= \lambda_{\mathbf{b}}(t_0), \quad (20)$$

where integration by part has been used.

Because the ODE system is reversible, we do not need to store the forward solutions \mathbf{u} in memory. Instead, when solving the adjoint system (13), we can simultaneously solve \mathbf{u} backward in time starting from $\mathbf{u}(t_f)$. Therefore the continuous adjoint method can be trivially implemented with constant memory, which is important for large scale problems.

5 The role of automatic differentiation

In both discrete and continuous adjoint approaches, we need to compute something like $\lambda \frac{\partial f}{\partial \mathbf{u}}$ which is essentially vector-Jacobian product (VJP). We use JAX to compute these VJPs automatically for us.

5.1 Implementing JVP and VJP

The ODE-constrained optimization problem:

$$\min_{\mathbf{a} \in \mathbb{R}^l, \mathbf{b} \in \mathbb{R}^k} J(\mathbf{a}, \mathbf{b}) = h(\mathbf{u}(t_f, \mathbf{a}, \mathbf{b})), \quad (21a)$$

$$\text{s.t. } \frac{d}{dt} \mathbf{u}(t, \mathbf{a}, \mathbf{b}) = \mathbf{r}(\mathbf{u}, \mathbf{a}), \quad (21b)$$

$$\mathbf{u}(t_0, \mathbf{a}, \mathbf{b}) = \mathbf{b}, \quad (21c)$$

The goal is to compute the gradient $\frac{\partial J}{\partial \mathbf{a}}$ and $\frac{\partial J}{\partial \mathbf{b}}$.

5.1.1 JVP

Given perturbations $\Delta \mathbf{a} \in \mathbb{R}^l$ and $\Delta \mathbf{b} \in \mathbb{R}^k$, we define a new state variable $\mathbf{z} = \frac{\partial \mathbf{u}}{\partial \mathbf{a}} \Delta \mathbf{a} + \frac{\partial \mathbf{u}}{\partial \mathbf{b}} \Delta \mathbf{b}$. We have $\mathbf{z} \in \mathbb{R}^k$. The goal is to compute $\mathbf{z}(t_f, \mathbf{a}, \mathbf{b})$. The ODE for \mathbf{z} is

$$\frac{d}{dt} \mathbf{z}(t, \mathbf{a}, \mathbf{b}) = \frac{\partial \mathbf{r}}{\partial \mathbf{u}} \mathbf{z} + \frac{\partial \mathbf{r}}{\partial \mathbf{a}} \Delta \mathbf{a}, \quad (22)$$

$$\mathbf{z}(t_0, \mathbf{a}, \mathbf{b}) = \Delta \mathbf{b}. \quad (23)$$

5.1.2 VJP

Given perturbation $\mathbf{g} \in \mathbb{R}^k$, we define the adjoint variables $\boldsymbol{\lambda}_a \in \mathbb{R}^l$ and $\boldsymbol{\lambda}_b \in \mathbb{R}^k$. The goal is to compute $\mathbf{g} \frac{\partial \mathbf{u}(t_f, \mathbf{a}, \mathbf{b})}{\partial \mathbf{a}}$ and $\mathbf{g} \frac{\partial \mathbf{u}(t_f, \mathbf{a}, \mathbf{b})}{\partial \mathbf{b}}$. The ODE for $[\boldsymbol{\lambda}_a, \boldsymbol{\lambda}_b]$ is

$$\frac{d}{dt} \begin{bmatrix} \boldsymbol{\lambda}_a \\ \boldsymbol{\lambda}_b \end{bmatrix} = \begin{bmatrix} -\boldsymbol{\lambda}_b \frac{\partial \mathbf{r}}{\partial \mathbf{a}} \\ -\boldsymbol{\lambda}_b \frac{\partial \mathbf{r}}{\partial \mathbf{u}} \end{bmatrix}, \quad (24)$$

$$\begin{bmatrix} \boldsymbol{\lambda}_a(t_f) \\ \boldsymbol{\lambda}_b(t_f) \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{g} \end{bmatrix}. \quad (25)$$

The adjoint variables do have interpretations [2]. It can be helpful to think about an objective function defined as $L = \mathbf{g} \cdot \mathbf{u}(t_f, \mathbf{a}, \mathbf{b})$.

References

- [1] S. G. Johnson, “Adjoint methods and sensitivity analysis for recurrence relations,” *Course notes for MIT*, vol. 18, pp. 2007–2011, 2007.
- [2] R. T. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud, “Neural ordinary differential equations,” *arXiv preprint arXiv:1806.07366*, 2018.