

# Adjoint Sensitivity Equation and NeuralODE

A/Prof Richard Yi Da Xu

`richardxu.com`

University of Technology Sydney (UTS)

July 26, 2020

This notes is in an elaborated attempt to explain:

- ▶ Qiqi Wang's YouTube lecture  
<https://www.youtube.com/watch?v=7CZP6dHIkNE>
- ▶ NeuralODE paper: <https://arxiv.org/abs/1806.07366>
- ▶ I did **not** try to unify notations, i.e., I keep notations identical to original reference
- ▶ but I will have a page to explain how they convert from one to the other

# Motivation: Solving ODE by Separation of Variables

- ▶ looking at the simplest ODE example:

$$\frac{dy}{dt} = f(t, y(t)) \quad y(t_0) = y_0$$

- ▶ solution

$$y(t) = y_0 + \int_{t=t_0}^t f(t, y(t)) dt$$

- ▶ example  $\frac{dy}{dt} = y(t) \quad y(0) = 0$

$$\frac{dy}{dt} = y(t)$$

$$\frac{dy}{y(t)} = dt$$

$$\int \frac{1}{y(t)} dy = \int dt$$

$$\ln(y) + C_Y = t + C_t$$

$$\ln(y) = t + C_1$$

$$y = \exp(t + C_1)$$

$$y = C \exp(t)$$

- ▶ substitute  $y(0) = 1$ :

$$1 = C \exp(0)$$

$$\implies C = 1$$

- ▶ solution:

$$y = \exp(t)$$

# Motivation: Solving ODE by approximation: Euler's method

$$\frac{dy}{dt} = f(t, y(t)), \quad y(t_0) = y_0$$

- Euler's method:

$$y_{n+1} = y_n + hf(t_n, y_n)$$

- compare with solutions on  $\frac{dy}{dt} = f(t, y(t)) = y \quad y(t_0) = y_0$  and let  $h = 1$ :

$$\begin{aligned} y_1 &= y_0 + hf(y_0) = 1 + 1 \cdot 1 = 2 & \exp(1) &= 2.71 \\ y_2 &= y_1 + hf(y_1) = 2 + 1 \cdot 2 = 4 & \exp(2) &= 7.38 \\ y_3 &= y_2 + hf(y_2) = 4 + 1 \cdot 4 = 8 & \exp(3) &= 20.08 \\ y_4 &= y_3 + hf(y_3) = 8 + 1 \cdot 8 = 16 & \exp(4) &= 54.59 \end{aligned}$$

- if we substitute  $f(t_n, y_n) \equiv \sigma(W_n^\top y_n + B_n)$

$$y_{n+1} = y_n + \sigma(W_n^\top y_n + B_n)$$

this is ResNet!

# More efficient ODE solver: Adjoint method

- ▶ Euler's method is inefficient and can be replaced by modern solver
- ▶ also, Neural networks's gradient descent requires to compute  $\frac{\partial \mathcal{L}}{\partial \theta}$
- ▶ how can we do that when we are given:

$$\frac{dh}{dt} \equiv f_{\theta}(t, h(t))$$

- ▶ we need to use **adjoint** method

# Adjoint method: motivation through simple example

- ▶ look at the problem:

$$\begin{aligned} \min_{s \in \mathcal{S}} C^\top x \\ \text{s.t. } Ax = b(s) \end{aligned}$$

- ▶ naive way to solve:

1. somehow get a collection of all  $s \in \mathcal{S}$ :
2. then for each  $s$ , one finds a corresponding  $x$  using:

$$Ax = b(s)$$

3. substitute  $\{x\}$  into  $C^\top x$  to see which is smallest

## Using Adjoint equation/method:

$$\begin{aligned} \min_s C^\top x \\ \text{s.t. } Ax = b(s) \end{aligned}$$

- ▶ Adjoint solution

$$\begin{aligned} x &= A^{-1}b(s) \\ C^\top x &= \underbrace{C^\top A^{-1}}_{\hat{x}} b(s) = \hat{x}b(s) \end{aligned}$$

- ▶  $\hat{x}$  is the Adjoint:

$$\hat{x} = C^\top A^{-1} \quad \text{which is easy, as it does not depend on } s$$

- ▶ finally,

$$\min_s C^\top x \rightarrow \min_s \hat{x}b(s)$$

# Apply to discrete difference equation:

$$\begin{aligned} \min_s J(u_n) \\ \text{s.t. } u_n = u_n(u_{n-1}, s) \end{aligned}$$

- ▶ use Lagrange multiplier  $\{\hat{u}_n\}$  make the evolution of  $u_1 \rightarrow u_2 \rightarrow \dots \rightarrow u_n$  satisfy:

$$J = J(u_n) + \hat{u}_n^\top \underbrace{(u_n - u_n(u_{n-1}, s))}_{u_n = u_n(u_{n-1}, s)} + \hat{u}_{n-1}^\top \underbrace{(u_{n-1} - u_{n-1}(u_{n-2}, s))}_{u_{n-1} = u_{n-1}(u_{n-2}, s)} + \dots + \hat{u}_0^\top \underbrace{(u_0 - u_0(s))}_{u_0 = u_0(s)}$$

- ▶ define  $\delta_s F \equiv \delta F$  is response of  $F$  to an infinitesimal perturbation of  $s$
- ▶ solution is:

$$\begin{aligned} \delta J = \frac{\partial J}{\partial u_n} \delta u_n + \hat{u}_n^\top \left( \delta u_n - \frac{\partial u_n}{\partial u_{n-1}} \delta u_{n-1} - \frac{\partial u_n}{\partial s} \delta s \right) \quad \text{both } u_n \text{ and } s \text{ is function of } s \\ + \hat{u}_{n-1}^\top \left( \delta u_{n-1} - \frac{\partial u_{n-1}}{\partial u_{n-2}} \delta u_{n-2} - \frac{\partial u_{n-1}}{\partial s} \delta s \right) \\ + \dots \\ + \hat{u}_0^\top \left( \delta u_{n-1} - \frac{\partial u_0}{\partial s} \delta s \right) \end{aligned}$$

- ▶  $s$  is parameter, i.e.,  $\theta$



# notes on response to an infinitesimal perturbation

- ▶ one way to look at:

$$\begin{aligned}\delta J &= \frac{\partial J}{\partial u_n} \delta u_n + \hat{u}_n^\top \left( \delta u_n - \frac{\partial u_n}{\partial u_{n-1}} \delta u_{n-1} - \frac{\partial u_n}{\partial s} \delta s \right) \quad \text{both } u_n \text{ and } s \text{ is function of } s \\ &\quad + \hat{u}_{n-1}^\top \left( \delta u_{n-1} - \frac{\partial u_{n-1}}{\partial u_{n-2}} \delta u_{n-2} - \frac{\partial u_{n-1}}{\partial s} \delta s \right) \\ &\quad + \dots \\ &\quad + \hat{u}_0^\top \left( \delta u_0 - \frac{\partial u_0}{\partial s} \delta s \right)\end{aligned}$$

- ▶ following expression is the usual expression for  $\frac{\partial J}{\partial s}$ :

$$\begin{aligned}\frac{\partial J}{\partial s} &= \frac{\partial J}{\partial u_n} \frac{\partial u_n}{\partial s} + \hat{u}_n^\top \left( \frac{\partial u_n}{\partial s} - \frac{\partial u_n}{\partial u_{n-1}} \frac{\partial u_{n-1}}{\partial s} - \frac{\partial u_n}{\partial s} \right) \quad \text{both } u_n \text{ and } s \text{ is function of } s \\ &\quad + \hat{u}_{n-1}^\top \left( \frac{\partial u_{n-1}}{\partial s} - \frac{\partial u_{n-1}}{\partial u_{n-2}} \frac{\partial u_{n-2}}{\partial s} - \frac{\partial u_{n-1}}{\partial s} \right) \\ &\quad + \dots \\ &\quad + \hat{u}_0^\top \left( \frac{\partial u_0}{\partial s} - \frac{\partial u_0}{\partial s} \right)\end{aligned}$$

# Clever assignment of Lagrange/Adjoint variables: $\hat{u}_n^\top = -\frac{\partial J}{\partial u_n}$

$$\begin{aligned}
 \delta J &= \frac{\partial J}{\partial u_n} \delta u_n + \hat{u}_n^\top \left( \delta u_n - \frac{\partial u_n}{\partial u_{n-1}} \delta u_{n-1} - \frac{\partial u_n}{\partial s} \delta s \right) + \dots \\
 &= \frac{\partial J}{\partial u_n} \delta u_n - \frac{\partial J}{\partial u_n} \left( \delta u_n - \frac{\partial u_n}{\partial u_{n-1}} \delta u_{n-1} - \frac{\partial u_n}{\partial s} \delta s \right) + \dots \quad \text{let } \hat{u}_n^\top = -\frac{\partial J}{\partial u_n} \\
 &= \frac{\partial J}{\partial u_n} \delta u_n - \frac{\partial J}{\partial u_n} \delta u_n + \frac{\partial J}{\partial u_n} \frac{\partial u_n}{\partial u_{n-1}} \delta u_{n-1} + \frac{\partial J}{\partial u_n} \frac{\partial u_n}{\partial s} \delta s + \dots \\
 &= \underbrace{-\hat{u}_n^\top \frac{\partial u_n}{\partial s} \delta s + \frac{\partial J}{\partial u_{n-1}} \delta u_{n-1}}_{\text{bring second term}} + \dots \\
 &= -\hat{u}_n^\top \frac{\partial u_n}{\partial s} \delta s + \frac{\partial J}{\partial u_{n-1}} \delta u_{n-1} + \hat{u}_{n-1}^\top \left( \delta u_{n-1} - \frac{\partial u_{n-1}}{\partial u_{n-2}} \delta u_{n-2} - \frac{\partial u_{n-1}}{\partial s} \delta s \right) + \dots \quad \text{bring second term} \\
 &= -\hat{u}_n^\top \frac{\partial u_n}{\partial s} \delta s - \hat{u}_{n-1}^\top \delta u_{n-1} + \hat{u}_{n-1}^\top \delta u_{n-1} - \hat{u}_{n-1}^\top \frac{\partial u_{n-1}}{\partial u_{n-2}} \delta u_{n-2} - \hat{u}_{n-1}^\top \frac{\partial u_{n-1}}{\partial s} \delta s + \dots \quad \frac{\partial J}{\partial u_{n-1}} = -\hat{u}_{n-1}^\top \\
 &= \underbrace{-\hat{u}_n^\top \frac{\partial u_n}{\partial s} \delta s - \hat{u}_{n-1}^\top \frac{\partial u_{n-1}}{\partial s} \delta s - \hat{u}_{n-1}^\top \frac{\partial u_{n-1}}{\partial u_{n-2}} \delta u_{n-2}}_{\text{bring second term}} + \dots \\
 &= \sum_{i=0}^n -\hat{u}_i^\top \frac{\partial u_i}{\partial s} \delta s
 \end{aligned}$$

- ▶ so we have an expression for  $\delta J = \sum_{i=0}^n -\hat{u}_i^\top \frac{\partial u_i}{\partial s} \delta s$
- ▶ we must solve for  $\{\hat{u}_i\}$

# Order to solve for Adjoint equations

$$\begin{aligned}\hat{u}_n^\top &= -\frac{\partial J}{\partial u_n} & \frac{\partial J}{\partial u_{n-1}} &= -\hat{u}_{n-1}^\top \\ \implies \frac{\partial J}{\partial u_{n-1}} &= \frac{\partial J}{\partial u_n} \frac{\partial u_n}{\partial u_{n-1}} = -\hat{u}_{n-1}^\top \\ \implies -\hat{u}_n^\top \frac{\partial u_n}{\partial u_{n-1}} &= -\hat{u}_{n-1}^\top \\ \implies \hat{u}_{n-1}^\top &= \hat{u}_n^\top \frac{\partial u_n}{\partial u_{n-1}}\end{aligned}$$

- ▶ meaning we use back-propagation to solve for  $\{\hat{u}_i\}$

# Continuous case with ODE constraint (1)

- Difference equations:

$$\begin{aligned} & \min_s J(u_n) \\ & \text{s.t. } u_n = u_n(u_{n-1}, s) \\ \implies J = & \textcolor{red}{J}(u_n) + \hat{u}_n^\top (u_n - u_n(u_{n-1}, s)) + \hat{u}_{n-1}^\top (u_{n-1} - u_{n-1}(u_{n-2}, s)) + \cdots + \hat{u}_0^\top (u_0 - u_0(s)) \end{aligned}$$

- Continuous, ODE constrained problem:

$$\begin{aligned} & \min_{\textcolor{red}{s}} J_f(u(T)) + \int_0^T J_c(u(t)) dt \\ & \text{s.t. } \frac{du}{dt} = f(u, \textcolor{red}{s}(t)) \\ \implies J = & \textcolor{red}{J}_f(u(T)) + \int_0^T \textcolor{red}{J}_c(u(t)) dt + \int_0^T \hat{u}(t) \left( \frac{du}{dt} - f(u, s) \right) dt \quad \text{continuous Lagrange} \end{aligned}$$

- difficulty arise as having a derivative function in the Lagrange

# Simplify objective equation

- ▶ instead of including  $\int_0^T J_c(u(t))dt$  (for aerodynamics application)

$$\begin{aligned} & \min_{\mathbf{s}} J_f(u(T)) + \int_0^T J_c(u(t))dt \\ & \text{s.t. } \frac{du}{dt} = f(u, \mathbf{s}(t)) \\ \Rightarrow J &= J_f(u(T)) + \int_0^T J_c(u(t))dt + \int_0^T \hat{u}(t) \left( \frac{du}{dt} - f(u, \mathbf{s}) \right) dt \end{aligned}$$

- ▶ we simplify by removing  $\int_0^T J_c(u(t))dt$ :

$$\begin{aligned} & \min_{\mathbf{s}} J_f(u(T)) \\ & \text{s.t. } \frac{du}{dt} = f(u, \mathbf{s}(t)) \\ \Rightarrow J &= J_f(u(T)) + \int_0^T \hat{u}(t) \left( \frac{du}{dt} - f(u, \mathbf{s}) \right) dt \end{aligned}$$

# Adjoint equation for ODE (1):

- substitute:  $\hat{u}(T) = -\frac{\partial J_f}{\partial u}$ :

$$\begin{aligned}
 J &= J_f(u(T)) + \int_0^T \hat{u}(t) \left( \frac{du}{dt} - \underbrace{f(u, s)} \right) dt \\
 \implies \delta J &= \frac{\partial J_f}{\partial u} \delta u(T) dt + \int_0^T \hat{u}(t) \left( \frac{d\delta u}{dt} - \underbrace{\frac{\partial f}{\partial u} \delta u - \frac{\partial f}{\partial s} \delta s} \right) dt \\
 &= \frac{\partial J_f}{\partial u} \delta u(T) + \int_0^T \hat{u}(t) \frac{d\delta u}{dt} dt - \int_0^T \left( \frac{\partial f}{\partial u} \right)^\top \hat{u}(t) \delta u dt - \int_0^T \hat{u}(t) \frac{\partial f}{\partial s} \delta s dt \\
 &= \frac{\partial J_f}{\partial u} \delta u(T) + \underbrace{\left[ \hat{u}(t) \delta u \right]_0^T}_{\text{blue}} - \int_0^T \delta u \frac{d\hat{u}}{dt} dt - \int_0^T \left( \frac{\partial f}{\partial u} \right)^\top \hat{u}(t) \delta u dt - \int_0^T \hat{u}(t) \frac{\partial f}{\partial s} \delta s dt \\
 &= \frac{\partial J_f}{\partial u} \delta u(T) + \underbrace{\hat{u}(T) \delta u(T)}_{\text{red}} - \int_0^T \delta u \frac{d\hat{u}}{dt} dt - \int_0^T \left( \frac{\partial f}{\partial u} \right)^\top \hat{u}(t) \delta u dt - \int_0^T \hat{u}(t) \frac{\partial f}{\partial s} \delta s dt \\
 &= \cancel{\frac{\partial J_f}{\partial u} \delta u(T)} + \cancel{-\frac{\partial J_f}{\partial u} \delta u(T)}_{\text{red}} - \int_0^T \delta u \frac{d\hat{u}}{dt} dt - \int_0^T \left( \frac{\partial f}{\partial u} \right)^\top \hat{u}(t) \delta u dt - \int_0^T \hat{u}(t) \frac{\partial f}{\partial s} \delta s dt \\
 &= - \int_0^T \delta u \frac{d\hat{u}}{dt} dt - \int_0^T \hat{u}(t) \frac{\partial f}{\partial u} \delta u dt - \int_0^T \hat{u}(t) \frac{\partial f}{\partial s} \delta s dt
 \end{aligned}$$

- integration by parts:  $\int_a^b u(x) v'(x) dx = [u(x) v(x)]_a^b - \int_a^b u'(x) v(x) dx$

## Adjoint equation for ODE (2)

- ▶ we already let  $\hat{u}(T) = -\frac{\partial J_f}{\partial u(T)}$
- ▶ further, we let  $\hat{u}(t)$ : instead of let  $\hat{u}(t) = -\frac{\partial J_f}{\partial u(t)} \equiv -\frac{\partial J_f}{\partial u}(t)$ :
- ▶ **important: if we can prove:**

$$\hat{u}(t) = -\frac{\partial J_f}{\partial u(t)} \implies \frac{d\hat{u}}{dt} = -\left(\frac{\partial f}{\partial u(t)}\right)^\top \hat{u}(t)$$

- ▶ then we can solve  $\delta J$ :

$$\begin{aligned}\delta J &= -\int_0^T \delta u \frac{d\hat{u}}{dt} dt - \int_0^T \left(\frac{\partial f}{\partial u(t)}\right)^\top \hat{u}(t) \delta u dt - \int_0^T \hat{u}(t) \frac{\partial f}{\partial s} \delta s dt \\&= \int_0^T \left(\frac{\partial f}{\partial u(t)}\right)^\top \hat{u}(t) \delta u(t) dt - \int_0^T \left(\frac{\partial f}{\partial u}\right)^\top \hat{u}(t) \delta u dt - \int_0^T \hat{u}(t) \frac{\partial f}{\partial s} \delta s dt \\&= -\int_0^T \hat{u}(t) \frac{\partial f}{\partial s} \delta s dt\end{aligned}$$

- ▶ by running another ODE to solve  $\hat{u}(t)$ , from dynamic equation  $\frac{d\hat{u}}{dt}$

# Compare with discrete/difference equation

- ▶ continuous case:

$$\delta J = - \int_0^T \hat{u}(t) \frac{\partial f}{\partial s} \delta s \, dt$$

- ▶ compare with discrete case:

$$\delta J = \sum_{i=0}^n -\hat{u}_i^\top \frac{\partial u_i}{\partial s} \delta s$$



$$\text{Proof for } \hat{u}(t) = -\frac{\partial J_f}{\partial u(t)} \implies \frac{d\hat{u}}{dt} = -\left(\frac{\partial f}{\partial u(t)}\right)^\top \hat{u}(t) \quad (1)$$

- here we change notions to use Neural ODE:

$$\begin{aligned} \hat{u}(t) &= -\frac{\partial J_f}{\partial u(t)} \implies \frac{d\hat{u}}{dt} = -\left(\frac{\partial f}{\partial u(t)}\right)^\top \hat{u}(t) \\ \mathbf{a}(t) &= \frac{\partial L}{\partial \mathbf{z}(t)} \implies \frac{d\mathbf{a}}{dt} = -\left(\frac{\partial f}{\partial \mathbf{z}(t)}\right)^\top \mathbf{a}(t) \end{aligned}$$

- similar

$$\begin{aligned} \delta J &= -\int_0^T \hat{u}(t) \frac{\partial f}{\partial s} \delta s \, dt \\ \frac{\partial L}{\partial \theta} &\equiv \frac{\partial L}{\partial \theta}(t_0) = -\int_{t_N}^{t_0} \mathbf{a}(t) \frac{\partial f}{\partial \theta} \, dt \end{aligned}$$

- sign change may be caused by swapping the integrand  $t_0$  and  $t_N$

$$\text{Proof for } \mathbf{a}(t) = \frac{\partial L}{\partial \mathbf{z}(t)} \implies \frac{d\mathbf{a}}{dt} = - \left( \frac{\partial f}{\partial \mathbf{z}(t)} \right)^\top \mathbf{a}(t)$$

since  $\frac{d\mathbf{z}(t)}{dt} = f(\mathbf{z}(t), t)$ :

$$\mathbf{z}(t + \epsilon) = \int_t^{t+\epsilon} f(\mathbf{z}(t), t) dt = T_\epsilon(\mathbf{z}(t), t)$$

$$\frac{dL}{\partial \mathbf{z}(t)} = \frac{dL}{d\mathbf{z}(t + \epsilon)} \frac{d\mathbf{z}(t + \epsilon)}{d\mathbf{z}(t)} \implies \mathbf{a}(t) = \mathbf{a}(t + \epsilon) \frac{\partial T_\epsilon(\mathbf{z}(t), t)}{\partial \mathbf{z}(t)}$$

$$\begin{aligned} \frac{d\mathbf{a}(t)}{dt} &= \lim_{\epsilon \rightarrow 0+} \frac{\mathbf{a}(t + \epsilon) - \mathbf{a}(t)}{\epsilon} \\ &= \lim_{\epsilon \rightarrow 0+} \frac{\mathbf{a}(t + \epsilon) - \mathbf{a}(t + \epsilon) \frac{\partial}{\partial \mathbf{z}(t)} T_\epsilon(\mathbf{z}(t), t)}{\epsilon} \\ &\approx \lim_{\epsilon \rightarrow 0+} \frac{\mathbf{a}(t + \epsilon) - \mathbf{a}(t + \epsilon) \frac{\partial}{\partial \mathbf{z}(t)} (\mathbf{z}(t) + \epsilon f(\mathbf{z}(t), t))}{\epsilon} \quad \text{Taylor } T_\epsilon(\mathbf{z}(t), t) \approx \mathbf{z}(t) + \epsilon f(\mathbf{z}(t), t) \\ &= \lim_{\epsilon \rightarrow 0+} \frac{\mathbf{a}(t + \epsilon) - \mathbf{a}(t + \epsilon) \left( \frac{\partial \mathbf{z}(t)}{\partial \mathbf{z}(t)} + \frac{\partial \epsilon f(\mathbf{z}(t), t)}{\partial \mathbf{z}(t)} \right)}{\epsilon} = \lim_{\epsilon \rightarrow 0+} \frac{\mathbf{a}(t + \epsilon) - \mathbf{a}(t + \epsilon) \epsilon \left( I + \frac{\partial f(\mathbf{z}(t), t)}{\partial \mathbf{z}(t)} \right)}{\epsilon} \\ &= \lim_{\epsilon \rightarrow 0+} \frac{-\mathbf{a}(t + \epsilon) \epsilon \frac{\partial f(\mathbf{z}(t), t)}{\partial \mathbf{z}(t)}}{\epsilon} = \lim_{\epsilon \rightarrow 0+} -\mathbf{a}(t + \epsilon) \frac{\partial f(\mathbf{z}(t), t)}{\partial \mathbf{z}(t)} \\ &= -\mathbf{a}(t) \frac{\partial f(\mathbf{z}(t), t)}{\partial \mathbf{z}(t)} \\ &= - \left( \frac{\partial f}{\partial \mathbf{z}(t)} \right)^\top \mathbf{a}(t) \quad \text{multi-dimensional case} \end{aligned}$$

# How to compute adjoint $\mathbf{a}(t)$ from its dynamics $\frac{d\mathbf{a}(t)}{dt}$

- ▶ just like  $\hat{x} = C^T x$ , but complicated.
- ▶ initial condition:  $\mathbf{a}(t_N) = \frac{dL}{dz(t_N)}$ , and run backwards to obtain  $\mathbf{a}(t_0)$ :

$$\begin{aligned}\mathbf{a}(t_0) &= \mathbf{a}(t_N) + \int_{t_N}^{t_0} \frac{d\mathbf{a}(t)}{dt} dt \\ &= \mathbf{a}(t_N) - \int_{t_N}^{t_0} \left( \frac{\partial f}{\partial \mathbf{z}(t)} \right)^T \mathbf{a}(t) dt \quad \text{substitute}\end{aligned}$$

# main learning task $\frac{\partial L}{\partial \theta}$ and $\frac{\partial L}{\partial t}$ (1)

►  $\frac{\partial L}{\partial \theta}$ :

$$\frac{\partial L}{\partial \theta} \equiv \frac{\partial L}{\partial \theta}(t_0) = \frac{\partial L}{\partial \theta}(t_N) - \int_{t_N}^{t_0} \mathbf{a}(t) \frac{\partial f}{\partial \theta}(t) dt$$

dynamic equation:  $-\mathbf{a}(t) \frac{\partial f}{\partial \theta}(t)$       initial condition:  $\frac{\partial L}{\partial \theta}(t_N) = \mathbf{0}$

►  $\frac{dL}{dt_0}$ :

$$\frac{dL}{dt_0} \equiv \frac{dL}{dt}(t_0) = \frac{dL}{dt}(t_N) - \int_{t_N}^{t_0} \mathbf{a}(t) \frac{df}{dt}(t) dt$$

dynamic equation:  $-\mathbf{a}(t) \frac{df}{dt}(t)$       initial condition:  $\frac{dL}{dt}(t_N) = \mathbf{a}(t_N) f(\mathbf{z}(t_N, t_N))$

- combine with  $\mathbf{a}(t_0)$ :

$$\mathbf{a}(t_0) = \mathbf{a}(t_N) - \int_{t_N}^{t_0} \mathbf{a}(t) \frac{\partial f}{\partial \mathbf{z}}(t) dt$$

dynamic equation:  $-\mathbf{a}(t) \frac{\partial f}{\partial \mathbf{z}}(t)$

initial condition:  $\frac{\partial L}{\partial \mathbf{z}(t_N)} = \mathbf{a}(t_N)$

- combine with  $\mathbf{z}(t_0)$ :

$$\mathbf{z}(t_0) = \mathbf{z}(t_N) - \int_{t_N}^{t_0} f(\mathbf{z}, t) dt$$

dynamic equation:  $f(\mathbf{z}, t)$

initial condition:  $\mathbf{z}(t_N)$

# Combine everything in a backward $t_N \rightarrow t_0$ pass

- combined back dynamics:

$$- \begin{bmatrix} f(\mathbf{z}, t) & \mathbf{a} \frac{\partial f}{\partial \mathbf{z}} & \mathbf{a} \frac{\partial f}{\partial \theta} & \mathbf{a} \frac{df}{dt} \end{bmatrix} (t)$$

- to solve:

$$\begin{bmatrix} \mathbf{z}(t_0) & \mathbf{a}(t_0) & \frac{\partial L}{\partial \theta}(t_0) & \frac{dL}{dt}(t_0) \end{bmatrix}$$

- with initial condition at  $t_N$ :

$$\begin{bmatrix} \mathbf{z}(t_N) & \mathbf{a}(t_N) & \mathbf{0} & \frac{df}{dt}(t_N) \end{bmatrix}$$

# Neural ODE: feed-forward much easier

- Inference using feed-forward:

$$\mathcal{L}(\mathbf{z}(t_1)) = \mathcal{L}\left(\mathbf{z}(t_0) + \int_{t=0}^{t_1} \underbrace{f(\mathbf{z}(t), t, \theta)}_{\text{NN}^t(\mathbf{z})} dt\right)$$

- $f(\mathbf{z}(t), t, \theta) \equiv \text{NN}^t(\mathbf{z})$  is Neural Network at  $t^{\text{th}}$  infinitesimal layer
- we can solve this by in a Solver, to specify:
  1. boundary condition:  $\mathbf{z}(t_0)$
  2. start and end time  $t_0$  and  $t_1$
  3. dynamic equation:  $f(\mathbf{z}(t), t, \theta)$
- the solver looks like:

$$\mathcal{L}(\mathbf{z}(t_1)) = \mathcal{L}\left(\text{ODESolve}(\mathbf{z}(t_0), t_0, t_1, f(\mathbf{z}(t), t, \theta))\right)$$