

Method of Adjoints

Cheng Soon Ong
Marc Peter Deisenroth

December 2020



Motivation

Automatic differentiation

Augment each variable (for example a) with an **adjoint** variable \overleftarrow{a} to form an adjoint pair (a, \overleftarrow{a}) .

Adjoint, I've heard that somewhere else...

Linear dynamical system

Consider a dynamical system with state variable \mathbf{x}_t and input variable \mathbf{u}_t . Assume that state evolves according to linear dynamics

$$\mathbf{x}_{t+1} = \mathbf{A}\mathbf{x}_t + \mathbf{B}\mathbf{u}_t \quad \text{for } t = 0, 1, \dots$$

where (\mathbf{A}, \mathbf{B}) are known state evolution matrices.

Find good control inputs

Find a sequence of inputs u_t that minimizes some quadratic cost over the trajectory, for a given x_0 :

$$\min_{u_t, x_t} \quad \frac{1}{2} x_N^\top S x_{N+1} + \frac{1}{2} \sum_{t=0}^N x_t^\top Q x_t + u_t^\top R u_t,$$

$$\text{subject to} \quad x_{t+1} = Ax_t + Bu_t, \quad \text{for } t = 0, 1, \dots, N.$$

Lagrangian

The Lagrangian has the form (with Lagrange multiplier λ)

$$\mathcal{L}(x, u, \lambda) = \frac{1}{2} x_N^\top S x_{N+1} + \frac{1}{2} \sum_{t=0}^N x_t^\top Q x_t + u_t^\top R u_t - \lambda_t^\top (x_{t+1} - Ax_t + Bu_t).$$

To satisfy $\nabla_{x_t} \mathcal{L} = 0$ we solve for the update equation for λ ,



$$\lambda_{t-1} = A^\top \lambda_t + Q x_t.$$

This is called the **adjoint dynamics** (with initial condition $\lambda_N = S x_{N+1}$).

Caching. In reverse.

We build λ_{t-1} by multiplying λ_t with A

What is in a name?

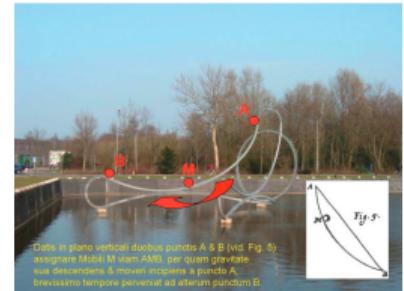
In dynamical systems, the Lagrange multipliers are called **costates** or **adjoint variables** and the dual optimization problem is called the **adjoint problem**.

“Adjoint functors arise everywhere”,

– Saunders Mac Lane, 1998, Categories for the Working Mathematician.

Pointers to literature

- ▶ Carathéodory's royal road (Pesch, 2012)
- ▶ Pontryagin's maximum principle (Pontryagin et al., 1964; Ohsawa, 2015)
- ▶ Book on optimal control (Anderson and Moore, 2007)
- ▶ Neural network view (Chen et al., 2018; Finlay et al., 2020)
- ▶ ICERM workshop on Scientific Machine Learning
<https://icerm.brown.edu/events/ht19-1-sml/>



Sculpture at Groningen,
via Pesch & Plail 6

Find good control inputs

Find a sequence of inputs u_t that minimizes some quadratic cost over the trajectory, for a given x_0 :

$$\min_{u_t, x_t} \quad \frac{1}{2} x_N^\top S x_{N+1} + \frac{1}{2} \sum_{t=0}^N x_t^\top Q x_t + u_t^\top R u_t,$$

$$\text{subject to} \quad x_{t+1} = Ax_t + Bu_t, \quad \text{for } t = 0, 1, \dots, N.$$

We consider the more general function

$$\begin{aligned} x^* = & \quad \operatorname{argmin}_x F(u, x) & = & \quad \operatorname{argmin}_x F(g(x), x) \\ & \text{subject to } G(u, x) = 0 & & \text{subject to } G(u, x) = 0. \end{aligned}$$

Recall: Jacobian

Definition (Jacobian)

The collection of all first-order partial derivatives of a vector-valued function

$$\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

is called the **Jacobian**.

$$\begin{aligned} J = \nabla_{\mathbf{x}} \mathbf{f} &= \frac{d\mathbf{f}(\mathbf{x})}{d\mathbf{x}} = \left[\frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_1} \quad \dots \quad \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_n} \right] \\ &= \left[\begin{array}{ccc} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_1(\mathbf{x})}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_m(\mathbf{x})}{\partial x_n} \end{array} \right] \in \mathbb{R}^{m \times n} \qquad \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \\ &\qquad \qquad \qquad J(i, j) = \frac{\partial f_i}{\partial x_j}. \end{aligned}$$

Recall: Automatic differentiation

Backward pass of reverse mode automatic differentiation

Given a function $G(w) := g(f(e(w)))$, with intermediate variables x, y and final output z .



Input adjoint = output adjoint \times forward partials

$$\underbrace{\frac{dz}{dw}}_{\overleftarrow{w}} = \underbrace{\left(\frac{dz}{dy} \frac{dy}{dx} \right)}_{\overleftarrow{x}} \frac{dx}{dw} \quad (\text{reverse mode})$$

Vectorized computation

Efficient compute

Exploit computational architecture using the **Vector Jacobian Product**.

Assume that function $e : \mathbb{R}^n \rightarrow \mathbb{R}^m$

Input adjoint =	output adjoint \times	forward partials
$\overleftarrow{w} =$	$\overleftarrow{x} \times$	$\frac{\partial x}{\partial w}$
$\mathbb{R}^{1 \times n}$	$\mathbb{R}^{1 \times m}$	$\mathbb{R}^{m \times n}$

Constrained optimization

Consider the general optimization problem with two (potentially vector) variables:

$$x^* = \operatorname{argmin}_x F(g(x), x)$$

subject to $G(u, x) = 0$.

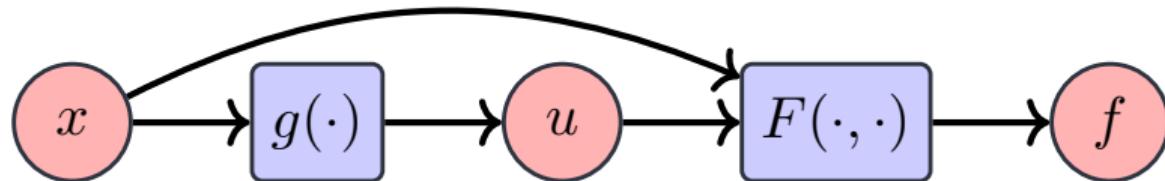
Constrained optimization

Consider the general optimization problem with two (potentially vector) variables:

$$x^* = \operatorname{argmin}_x F(g(x), x)$$

subject to $G(u, x) = 0$.

Think about how the information flows



Recall: Implicit function theorem

The implicit function theorem **provides conditions** under which we can write $G(x, y) = 0$ as $y = g(x)$, and also conditions when $g'(x) = -\frac{G_x}{G_y}$.

Hand wavy intuition

Given three variables (could be any topology) x, y, z and

$$G(x, y) = z,$$

What are the conditions for the following to be well behaved?

$$\begin{aligned} g(x) &= y \\ G(x, g(x)) &= z. \end{aligned}$$

Solve for gradient of constraint

Differentiate G with respect to x , using the chain rule,

$$\frac{dG}{dx} = \frac{\partial G}{\partial u} \frac{\partial u}{\partial x} + \frac{\partial G}{\partial x} \frac{\partial x}{\partial x}.$$

Solve for stationarity conditions, setting $\frac{dG}{dx} = 0$,

$$\frac{\partial u}{\partial x} = - \left(\frac{\partial G}{\partial u} \right)^{-1} \frac{\partial G}{\partial x}.$$

Chain rule on objective function F

Differentiate F with respect to x , using the chain rule,

$$\frac{dF}{dx} = \frac{\partial F}{\partial x} \frac{\partial x}{\partial x} + \frac{\partial F}{\partial u} \frac{\partial u}{\partial x}.$$

Substituting the solution of $\frac{\partial u}{\partial x}$ into $\frac{dF}{dx}$, we obtain

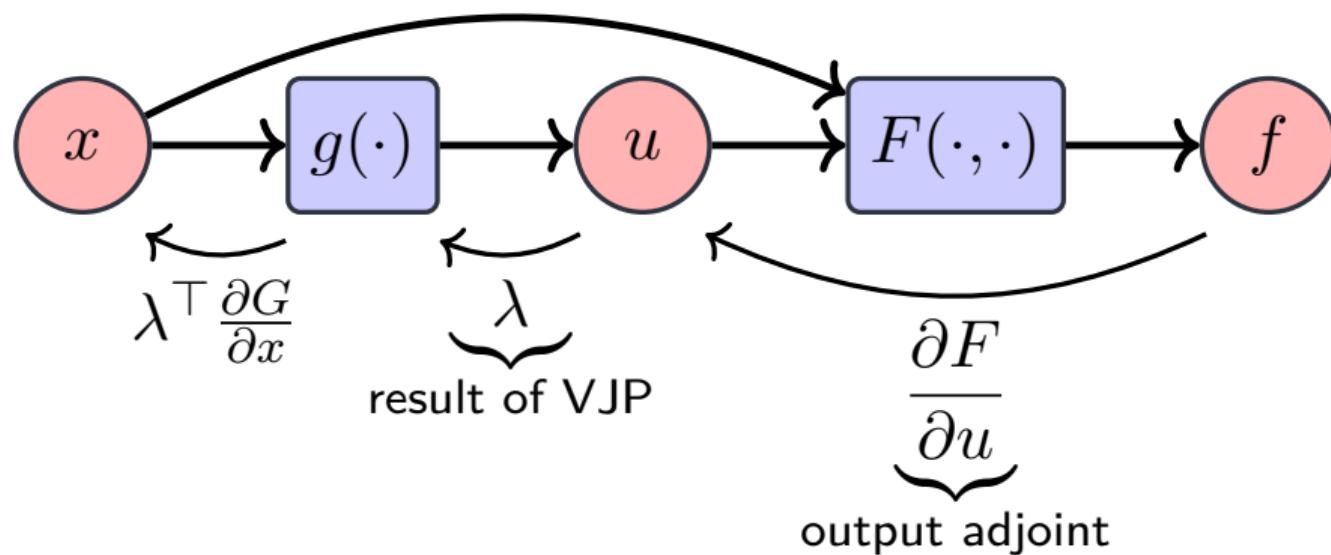
$$\frac{dF}{dx} = \frac{\partial F}{\partial x} - \underbrace{\frac{\partial F}{\partial u} \left(\frac{\partial G}{\partial u} \right)^{-1} \frac{\partial G}{\partial x}}_{= -\lambda^\top}.$$

Observe that λ can be computed by a **vector Jacobian** product.

Identify Lagrange multiplier λ

Rewrite in terms of λ .

$$\frac{dF}{dx} = \frac{\partial F}{\partial x} + \lambda^\top \frac{\partial G}{\partial x} \quad \text{where} \quad \lambda = \frac{\partial F}{\partial u} \left(\frac{\partial G}{\partial u} \right)^{-1}$$



Adjoint = Lagrange multiplier in method of adjoints

We consider the general function

$$\begin{aligned} x^* = \quad & \operatorname{argmin}_x F(u, x) \\ & \text{subject to } G(u, x) = 0 \end{aligned} \quad = \quad \begin{aligned} & \operatorname{argmin}_x F(g(x), x) \\ & \text{subject to } G(u, x) = 0. \end{aligned}$$

From the implicit function theorem

$$\frac{dF}{dx} = \frac{\partial F}{\partial x} + \lambda^\top \frac{\partial G}{\partial x} \quad \text{where} \quad \lambda = \frac{\partial F}{\partial u} \left(\frac{\partial G}{\partial u} \right)^{-1}$$

Alternatively, by the method of Lagrange

$$\nabla_x \mathcal{L} = \nabla_x F(x, u) + \lambda^\top \nabla_x G(x, u)$$

►► **adjoint dynamics** ($\lambda_{t-1} = A^\top \lambda_t + Q x_t$).

Summary

- ▶ Method of adjoints studied in optimal control
- ▶ ... adjoint state method, Pontryagin's principle
- ▶ Vector Jacobian products for efficient computation
- ▶ Adjoint variable in autodiff = Lagrange multiplier

Backpropagation is just ...

the method of adjoints

References

- Anderson, B. D. O. and Moore, J. B. (2007). *Optimal Control: Linear Quadratic Methods*. Dover.
- Chen, T. Q., Rubanova, Y., Bettencourt, J., and Duvenaud, D. K. (2018). Neural Ordinary Differential Equations. In *Advances in Neural Information Processing Systems*.
- Finlay, C., Jacobsen, J.-H., Nurbekyan, L., and Oberman, A. M. (2020). How to train your neural ODE: the world of Jacobian and kinetic regularization.
- Lane, S. M. (1998). *Categories for the Working Mathematician*. Springer, 2nd edition.
- Ohsawa, T. (2015). Contact Geometry of the Pontryagin Maximum Principle. *Automatica*, 55(2015):1–5.
- Pesch, H. J. (2012). In *Optimization Stories*.
- Pontryagin, L. S., Boltyanskii, V. G., Gamkrelidze, R. V., and Mishchenko, E. F. (1964). *The Mathematical Theory of Optimal Processes*. Pergamon Press.