

Differentiation through iterative solvers

Nicolas Nytko

October 31, 2022

1 Background

Consider some function $f : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ that describes the iterative process

$$\mathbf{x}^{(k+1)} = f(\mathbf{x}^{(k)}, \boldsymbol{\theta}), \quad (1)$$

where $\mathbf{x}^{(k)}$ is the state at iteration k and $\boldsymbol{\theta}$ is some set of parameters used at each step of the iteration that we would like to differentiate this process with respect to. We will assume that as $k \rightarrow \infty$, $\mathbf{x}^{(k)}$ converges to a *fixed point* such that

$$\mathbf{x}^* = f(\mathbf{x}^*, \boldsymbol{\theta}), \quad (2)$$

at least to within some small numerical tolerance. There are two ways to look at this problem depending on what information we are trying to obtain:

1. We care about the end result only. We assume that f is some reasonable method and we converge to \mathbf{x}^* . Lets call this the *convergent* analysis.
2. We want to look at intermediate results. There is no assumption that f actually converges to \mathbf{x}^* — we might be optimizing the method itself here. Lets call this one the *nonconvergent* analysis.

2 Convergent Methods

We assume that we run the method f for a sufficient number of iterations such that we have a reasonable approximation to \mathbf{x}^* , which we will denote by $\hat{\mathbf{x}}$. We therefore have

$$\hat{\mathbf{x}} = f(\hat{\mathbf{x}}, \boldsymbol{\theta}) + \boldsymbol{\varepsilon}(\boldsymbol{\theta}), \quad (3)$$

for some (hopefully) small $\boldsymbol{\varepsilon}$. Taking the derivative of both sides with respect to the parameter $\boldsymbol{\theta}$, we get

$$\frac{\partial \hat{\mathbf{x}}}{\partial \boldsymbol{\theta}} = \frac{\partial}{\partial \boldsymbol{\theta}} f(\hat{\mathbf{x}}, \boldsymbol{\theta}) \quad (4)$$

$$= \frac{\partial f}{\partial \boldsymbol{\theta}}(\hat{\mathbf{x}}, \boldsymbol{\theta}) + \frac{\partial f}{\partial \mathbf{x}}(\hat{\mathbf{x}}, \boldsymbol{\theta}) \frac{\partial \hat{\mathbf{x}}}{\partial \boldsymbol{\theta}} + \frac{\partial \boldsymbol{\varepsilon}}{\partial \boldsymbol{\theta}}. \quad (5)$$

Rearranging terms gives

$$\frac{\partial \hat{\mathbf{x}}}{\partial \boldsymbol{\theta}} = \left(\mathbf{I} - \frac{\partial f}{\partial \mathbf{x}}(\hat{\mathbf{x}}, \boldsymbol{\theta}) \right)^{-1} \left(\frac{\partial f}{\partial \boldsymbol{\theta}}(\hat{\mathbf{x}}, \boldsymbol{\theta}) + \frac{\partial \boldsymbol{\varepsilon}}{\partial \boldsymbol{\theta}} \right) \quad (6)$$

$$\approx \left(\mathbf{I} - \frac{\partial f}{\partial \mathbf{x}}(\hat{\mathbf{x}}, \boldsymbol{\theta}) \right)^{-1} \frac{\partial f}{\partial \boldsymbol{\theta}}(\hat{\mathbf{x}}, \boldsymbol{\theta}). \quad (7)$$

Of course, in an autograd setting we are interested in computing the vector-Jacobian product instead of the Jacobian itself; left multiplying eq. (7) by some vector \mathbf{v} (and dropping the \approx to simplify notation) results in

$$\mathbf{v}^T \frac{\partial \hat{\mathbf{x}}}{\partial \boldsymbol{\theta}} = \mathbf{v}^T \left(\mathbf{I} - \frac{\partial f}{\partial \mathbf{x}}(\hat{\mathbf{x}}, \boldsymbol{\theta}) \right)^{-1} \frac{\partial f}{\partial \boldsymbol{\theta}}(\hat{\mathbf{x}}, \boldsymbol{\theta}). \quad (8)$$

Defining the intermediate vector \mathbf{w} like

$$\mathbf{w}^T = \mathbf{v}^T \left(\mathbf{I} - \frac{\partial f}{\partial \mathbf{x}}(\hat{\mathbf{x}}, \boldsymbol{\theta}) \right)^{-1}, \quad (9)$$

we can rearrange terms to get

$$\mathbf{w}^T \left(\mathbf{I} - \frac{\partial f}{\partial \mathbf{x}}(\hat{\mathbf{x}}, \boldsymbol{\theta}) \right) = \mathbf{v}^T \quad (10)$$

$$\mathbf{w}^T - \mathbf{w}^T \frac{\partial f}{\partial \mathbf{x}}(\hat{\mathbf{x}}, \boldsymbol{\theta}) = \mathbf{v}^T \quad (11)$$

$$\mathbf{w}^T = \mathbf{v}^T + \mathbf{w}^T \frac{\partial f}{\partial \mathbf{x}}(\hat{\mathbf{x}}, \boldsymbol{\theta}), \quad (12)$$

which is itself a fixed-point iteration that we can find with intermediate vector-Jacobian products of f . Once we find \mathbf{w}^T , eq. (7) can be computed by VJP of f wrt $\boldsymbol{\theta}$ and \mathbf{w} .

3 Nonconvergent Methods

In the case that we do not run f to convergence, we will cast the iteration as an ODE like

$$\frac{d\mathbf{x}}{dt} = f(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{x} = g(\mathbf{x}, \boldsymbol{\theta}). \quad (13)$$

Observe that if we integrate eq. (13) with an Euler forward timestepper with $\Delta t = 1$ and let $\mathbf{x}^{(k)} = \mathbf{x}(k)$, we recover the iteration exactly:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + (\Delta t) g(\mathbf{x}^{(k)}, \boldsymbol{\theta}) \quad (14)$$

$$= \mathbf{x}^{(k)} + f(\mathbf{x}^{(k)}, \boldsymbol{\theta}) - \mathbf{x}^{(k)} \quad (15)$$

$$= f(\mathbf{x}^{(k)}, \boldsymbol{\theta}). \quad (16)$$

Assume we have run the above for j iterations and have obtained $\mathbf{x}^{(j)}$ (but not necessarily the intermediate steps). To differentiate some scalar loss (ℓ) with respect to $\mathbf{x}^{(0)}$ and $\boldsymbol{\theta}$, we will introduce the adjoint equation

$$\mathbf{a}(t) = \frac{\partial \ell}{\partial \mathbf{x}(t)}, \quad (17)$$

with derivative

$$\frac{d\mathbf{a}(t)}{dt} = -\mathbf{a}(t)^T \frac{\partial g}{\partial \mathbf{x}}(\mathbf{x}(t), \boldsymbol{\theta}). \quad (18)$$

Which we can integrate to get

$$\frac{d\ell}{d\mathbf{x}^{(0)}} = - \int_j^0 \mathbf{a}(t)^T \frac{\partial g(\mathbf{x}(t), \boldsymbol{\theta})}{\partial \mathbf{x}} \quad (19)$$

$$\frac{d\ell}{d\boldsymbol{\theta}} = - \int_j^0 \mathbf{a}(t)^T \frac{\partial g(\mathbf{x}(t), \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}. \quad (20)$$

Because we have discarded the intermediate results of $\mathbf{x}^{(i)}$, we can also integrate backwards from $\mathbf{x}^{(j)}$ at the same time using backwards Euler (see appendix A.1) to re-create the sequence of values.

3.1 Adjoint computation of gradient

We can now iteratively solve for the gradients $\frac{\partial \ell}{\partial \mathbf{x}^{(0)}}$ and $\frac{\partial \ell}{\partial \boldsymbol{\theta}}$. We begin with $\mathbf{x}^{(j)}$ and $\mathbf{a}^{(j)} = \mathbf{a}(j) = \frac{\partial \ell}{\partial \mathbf{x}^{(j)}}$. We first compute $\mathbf{x}^{(j-1)}$ by implicit Euler, which we will perform the nonlinear solve using a gradient descent.

Define the residual like

$$\mathbf{r} := \mathbf{x}^{(k)} - \mathbf{f}(\mathbf{y}, \boldsymbol{\theta}), \quad (21)$$

clearly if $\mathbf{y} = \mathbf{x}^{(k-1)}$ then $\mathbf{r} = \mathbf{0}$. To minimize this, we can optimize $\mathbf{r}^T \mathbf{r}$ by first taking the gradient wrt \mathbf{y} ,

$$\nabla_{\mathbf{y}} (\mathbf{r}^T \mathbf{r}) = -2\mathbf{x}^T J_{\mathbf{y}} (\mathbf{f}(\mathbf{y}, \boldsymbol{\theta})) + J_{\mathbf{y}} (\mathbf{f}(\mathbf{y}, \boldsymbol{\theta})) \mathbf{f}(\mathbf{y}, \boldsymbol{\theta}) + \mathbf{f}(\mathbf{y}, \boldsymbol{\theta})^T J_{\mathbf{y}} (\mathbf{f}(\mathbf{y}, \boldsymbol{\theta})) \quad (22)$$

$$= (\mathbf{f}(\mathbf{y}, \boldsymbol{\theta}) - 2\mathbf{x})^T J_{\mathbf{y}} (\mathbf{f}(\mathbf{y}, \boldsymbol{\theta})) + J_{\mathbf{y}} (\mathbf{f}(\mathbf{y}, \boldsymbol{\theta})) \mathbf{f}(\mathbf{y}, \boldsymbol{\theta}), \quad (23)$$

then descending on \mathbf{y} until we find the previous iterate. At each gradient calculation we must compute 3 vector-Jacobian products: most autograd software do not natively do Jacobian-vector products and instead implement it as two VJP.

$$\mathbf{a}^{(k-1)} = \mathbf{a}^{(k)} + \mathbf{a}^{(k-1)} J(x, \theta) \quad (24)$$

$$\mathbf{a}^{(k-1)} - \mathbf{a}^{(k-1)} J(x, \theta) - \mathbf{a}^{(k)} = 0 \quad (25)$$

$$\mathbf{a}^{(k-1)} (I - J(x, \theta)) - \mathbf{a}^{(k)} = 0 \quad (26)$$

$$\left(\mathbf{a}^{(k-1)} (I - J(x, \theta)) - \mathbf{a}^{(k)} \right)^T \left(\mathbf{a}^{(k-1)} (I - J(x, \theta)) - \mathbf{a}^{(k)} \right) \quad (27)$$

$$\left((I - J(x, \theta))^T \left(\mathbf{a}^{(k-1)} \right)^T - \left(\mathbf{a}^{(k)} \right)^T \right) \left(\mathbf{a}^{(k-1)} (I - J(x, \theta)) - \mathbf{a}^{(k)} \right) \quad (28)$$

$$(I - J(x, \theta))^T \left(\mathbf{a}^{(k-1)} \right)^T \mathbf{a}^{(k-1)} (I - J(x, \theta)) - 2 \left(\mathbf{a}^{(k)} \right)^T \mathbf{a}^{(k-1)} (I - J(x, \theta)) + \left(\mathbf{a}^{(k)} \right)^T \mathbf{a}^{(k)} \quad (29)$$

$$\left(\mathbf{a}^{(k-1)} \right)^T \mathbf{a}^{(k-1)} (I - J(x, \theta))^T (I - J(x, \theta)) - 2 \left(\mathbf{a}^{(k)} \right)^T \mathbf{a}^{(k-1)} (I - J(x, \theta)) + \left(\mathbf{a}^{(k)} \right)^T \mathbf{a}^{(k)} \quad (30)$$

A Misc

A.1 Backward Euler, backwards is forward Euler, forwards, but backwards

Let $x(t)$ be an ODE whose derivative is defined by

$$\frac{dx}{dt} = f(x, t). \quad (31)$$

If we time-step forward from $t^{(0)}$ to $t^{(1)}$ starting from $x^{(0)} = x(0)$ using forward Euler and obtain the intermediate sequence of $x(t)$ values

$$x^{(1)}, x^{(1)}, \dots, x^{(j)}, \quad (32)$$

we will obtain the same sequence of values if we time-step backward from $t^{(1)}$ to $t^{(0)}$ starting from $x^{(j)} = x(t^{(1)})$ by using backward Euler.

Proof. Denote the step size taken in both forward and backward Euler by Δt . Each forward iteration is defined by

$$x_{\text{fwd}}^{(k+1)} = x_{\text{fwd}}^{(k)} + \Delta t f \left(x_{\text{fwd}}^{(k)} \right), \quad (33)$$

with each backward iteration being equivalently defined by

$$x_{\text{bwd}}^{(k)} = x_{\text{bwd}}^{(k+1)} - \Delta t f \left(x_{\text{bwd}}^{(k)} \right). \quad (34)$$

Assuming we start the backward iteration from the final result obtained by the forward iteration, we have $x_{\text{bwd}}^{(j)} = x_{\text{fwd}}^{(j)}$. Inductively, we get

$$x_{\text{bwd}}^{(k)} = x_{\text{fwd}}^{(k+1)} - \Delta t f \left(x_{\text{bwd}}^{(k)} \right), \quad (35)$$

or, with some rearranging and from eq. (33),

$$x_{\text{bwd}}^{(k)} + \Delta t f \left(x_{\text{bwd}}^{(k)} \right) = x_{\text{fwd}}^{(k+1)} = x_{\text{fwd}}^{(k)} + \Delta t f \left(x_{\text{fwd}}^{(k)} \right). \quad (36)$$

□

References