

Adjoint methods and sensitivity analysis for recurrence relations

Steven G. Johnson

October 22, 2007

1 Introduction

In this note, we derive an adjoint method for sensitivity analysis of the solution of recurrence relations. In particular, we suppose that we have a K -component vector \mathbf{x} that is determined by iterating a recurrence relation

$$\mathbf{x}^n = \mathbf{f}(\mathbf{x}^{n-1}, \mathbf{p}, n) \triangleq \mathbf{f}^n$$

for some function \mathbf{f} depending on the previous \mathbf{x} ,¹ a vector \mathbf{p} of P parameters, and the step index n . The initial condition is

$$\mathbf{x}^0 = \mathbf{b}(\mathbf{p})$$

for some given function \mathbf{b} of the parameters. Furthermore, we have some function g of \mathbf{x} :

$$g^n \triangleq g(\mathbf{x}^n, \mathbf{p}, n)$$

and we wish to compute the gradient $\frac{dg^N}{d\mathbf{p}}$ of g^N , for some N , with respect to the parameters \mathbf{p} .

2 The explicit gradient

The gradient of g^N can be written explicitly as:

$$\frac{dg^N}{d\mathbf{p}} = g_{\mathbf{p}}^N + g_{\mathbf{x}}^N (\mathbf{f}_{\mathbf{p}}^N + \mathbf{f}_{\mathbf{x}}^N [\mathbf{f}_{\mathbf{p}}^{N-1} + \mathbf{f}_{\mathbf{x}}^{N-1} \{ \mathbf{f}_{\mathbf{p}}^{N-2} + \dots \}]), \quad (1)$$

where subscripts denote partial derivatives, and should be thought of as row vectors, vs. column vectors \mathbf{x} and \mathbf{p} . So, for example, $g_{\mathbf{x}}^N$ is a $1 \times K$ matrix, and $\mathbf{f}_{\mathbf{p}}^N$ is a $K \times P$ matrix. Equation (1) is derived simply by applying the chain rule to $g^N = g(\mathbf{x}^n, \mathbf{p}) = g(f(\mathbf{x}^{n-1}, \mathbf{p}), \mathbf{p}) = g(f(f(\mathbf{x}^{n-2}, \mathbf{p}), \mathbf{p}), \mathbf{p}) = \dots$.

¹Note that if \mathbf{x}^n depends on $\mathbf{x}^{n-\ell}$ for $\ell = 1, \dots, L$, the recurrence can still be cast in terms of \mathbf{x}^{n-1} alone by expanding \mathbf{x} into a vector of length KL , in much the same way that an L th-order ODE can be converted into L 1st-order ODEs.

The natural way to evaluate eq. (1) might seem to be starting at the innermost parentheses and working outwards, but this is inefficient. Each parenthesized expression is a $K \times P$ matrix that must be multiplied by \mathbf{f}_x^n , a $K \times K$ matrix, requiring $O(K^2P)$ time for each multiplication assuming dense matrices. There are $O(N)$ such multiplications, so evaluating the whole expression in this fashion requires $O(NK^2P)$ time. However, for dense matrices, the evaluation of g^N itself requires $O(NK^2)$ time, which means that the gradient (calculated this way) is as expensive as $O(P)$ evaluations of g^N .

Similarly, evaluating gradients by finite-difference approximations or similar numerical tricks requires $O(P)$ evaluations of the function being differentiated (e.g. center-difference approximations require two function evaluations per dimension). So, direct evaluation of the gradient by the above technique, while it may be more accurate than numerical approximations, is not substantially more efficient. This is a problem if P is large.

3 The gradient by adjoints

Instead of computing the gradient explicitly (by “forward” differentiation), *adjoint methods* typically allow one to compute gradients with the same cost as evaluating the function roughly twice, regardless of the number P of parameters [1, 2, 3, 4]. A very general technique for constructing adjoint methods involves something similar to Lagrange multipliers, where one adds zero to g^N in a way cleverly chosen to make computing the gradient easier, and in a previous note I derived the adjoint gradient for recurrence relations by this technique, analogous to work by Cao and Petzold on adjoint methods for differential-algebraic equations [2]. However, Gil Strang has pointed out to me that in many cases adjoint methods can be derived much more simply just by parenthesizing the gradient equation in a different way [4], and this turns out to be the case for the recurrence problem above.

The key fact is that, in the gradient equation (1), we are evaluating lots of expressions like $g_x^N(\mathbf{f}_x^N \mathbf{f}_p^{N-1})$, $g_x^N(\mathbf{f}_x^N [\mathbf{f}_x^{N-1} \mathbf{f}_p^{N-2}])$, and so on. Parenthesized this way, these expressions require $O(K^2P)$ operations each, because they involve matrix-matrix multiplications. However, we can parenthesize them a different way, so that they involve only vector-matrix multiplications, in order to reduce the complexity to $O(K^2 + KP)$, which is obviously a huge improvement for large K and P . In particular, parenthesize them as $(g_x^N \mathbf{f}_x^N) \mathbf{f}_p^{N-1}$, $[(g_x^N \mathbf{f}_x^N) \mathbf{f}_x^{N-1}] \mathbf{f}_p^{N-2}$, and so on, involving repeated multiplication of a row vector on the left (starting with g_x^N) by a matrix \mathbf{f}_x^n on the right. This repeated multiplication defines an *adjoint recurrence* relation for a K -component column vector λ^n , recurring backwards from $n = N$ to $n = 0$:

$$\lambda^{n-1} = (\mathbf{f}_x^n)^T \lambda^n,$$

where T is the transpose, with “initial” condition

$$\lambda^N = (g_x^N)^T.$$

In terms of this adjoint vector (so-called because of the transposes in the expressions above), the gradient becomes:

$$\frac{dg^N}{d\mathbf{p}} = g_{\mathbf{p}}^N + \sum_{n=1}^N (\boldsymbol{\lambda}^n)^T \mathbf{f}_{\mathbf{p}}^n + (\boldsymbol{\lambda}^0)^T \mathbf{b}_{\mathbf{p}}. \quad (2)$$

Consider the computational cost to evaluate the gradient in this way. Evaluating g^N and the \mathbf{x}^n costs $O(NK^2)$ time, assuming dense matrices, and evaluating $\boldsymbol{\lambda}^n$ also takes $O(NK^2)$ time. Finally evaluating equation (2) takes $O(NKP)$ time in the worst case, dominated by the time to evaluate the summation assuming $\mathbf{f}_{\mathbf{p}}^n$ is a dense matrix. So, the total is $O(NK^2 + NKP)$, much better than $O(NK^2P)$ for large K and P .

In practice, the situation is likely to be even better than this, because often $\mathbf{f}_{\mathbf{p}}^n$ will be a sparse matrix: each component of \mathbf{p} will appear only for certain components of \mathbf{x} and/or for certain steps n . In this case the $O(NKP)$ cost will be greatly reduced, e.g. to $O(NK)$ or $O(KP)$ or similar. Then the cost of the gradient will be dominated by the two $O(NK^2)$ recurrences—i.e., as is characteristic of adjoint methods, the cost of finding the gradient will be comparable to the cost of finding the function value twice.

Note that there is, however, at least one drawback of the adjoint method (2) in comparison to the direct method (1): the adjoint method may require more storage. For the direct method, $O(K)$ storage is required for the current \mathbf{x}^n (which can be discarded once \mathbf{x}^{n+1} is computed) and $O(PK)$ storage is required for the $K \times P$ matrix being accumulated, to be multiplied by $g_{\mathbf{x}}^N$ at the end, for $O(PK)$ storage total. In the adjoint method, all of the \mathbf{x}^n must be stored, because they are used in the backwards recurrence for $\boldsymbol{\lambda}^n$ once \mathbf{x}^N is reached, requiring $O(NK)$ storage. [The $\boldsymbol{\lambda}^n$ vectors, on the other hand, can be discarded once $\boldsymbol{\lambda}^{n-1}$ is computed, assuming the summation in eq. (2) is computed on the fly. Only $O(K)$ storage is needed for this summation, assuming $\mathbf{f}_{\mathbf{p}}^n$ can be computed on the fly (or is sparse).] Whether the $O(PK)$ storage for the direct method is better or worse than the $O(NK)$ storage for the adjoint method obviously depends on how P compares to N .

4 A simple example

Finally, let us consider a simple example of a $K = 2$ linear recurrence:

$$\mathbf{x}^n = A\mathbf{x}^{n-1} + \begin{pmatrix} 0 \\ p_n \end{pmatrix}$$

with an initial condition

$$\mathbf{x}^0 = \mathbf{b} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

and some 2×2 matrix A , e.g.

$$A = \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix}$$

for $\theta = 0.1$. Here, $P = N$: there are N parameters p_n , one per step n , acting as “source” terms in the recurrence (which otherwise has oscillating solutions since A is unitary). Let us also pick a simple function g to differentiate, e.g.

$$g(\mathbf{x}) = (x_2)^2.$$

The adjoint recurrence for λ^n is then:

$$\lambda^{n-1} = (\mathbf{f}_{\mathbf{x}}^n)^T \lambda^n = A^T \lambda^n,$$

with “initial” condition:

$$\lambda^N = (g_{\mathbf{x}}^N)^T = \begin{pmatrix} 0 \\ 2x_2^N \end{pmatrix}.$$

Notice that this case is rather simple: since our recurrence is linear, the adjoint recurrence does not depend on \mathbf{x}^n except in the initial condition.

The gradient is also greatly simplified because $\mathbf{f}_{\mathbf{p}}^n$ is sparse: it is a $2 \times N$ matrix of zeros, except for the n -th column which is $(0, 1)^T$. That means that the gradient (2) becomes:

$$\frac{dg^N}{dp_k} = \lambda_2^k,$$

requiring $O(N)$ work to find the whole gradient.

As a quick test, I implemented this example in GNU Octave (a Matlab clone) and checked it against the numerical center-difference gradient; it only takes a few minutes to implement and is worthwhile to try if you are not clear on how this works. For extra credit, try modifying the recurrence, e.g. to make \mathbf{f} nonlinear in \mathbf{x} and/or \mathbf{p} .

References

- [1] R. M. Errico, “What is an adjoint model?,” *Bulletin Am. Meteorological Soc.*, vol. 78, pp. 2577–2591, 1997.
- [2] Y. Cao, S. Li, L. Petzold, and R. Serban, “Adjoint sensitivity analysis for differential-algebraic equations: The adjoint DAE system and its numerical solution,” *SIAM J. Sci. Comput.*, vol. 24, no. 3, pp. 1076–1089, 2003.
- [3] S. G. Johnson, “Notes on adjoint methods for 18.336.” October 2007.
- [4] G. Strang, *Computational Science and Engineering*. Wellesley, MA: Wellesley-Cambridge Press, 2007.

MIT OpenCourseWare
<http://ocw.mit.edu>

18.335J / 6.337J Introduction to Numerical Methods
Fall 2010

For information about citing these materials or our Terms of Use, visit: <http://ocw.mit.edu/terms>.